# Repeatability of glucocorticoid hormones in vertebrates: A meta-analysis

**Kelsey L Schoenemann** [Corresp., 1] , **Frances Bonier** [1]

[1] Biology Department, Queen's University, Kingston, Ontario, Canada

Corresponding Author: Kelsey L Schoenemann
Email address: kelsey.schoene@gmail.com

We often expect that investigations of the patterns, causes, and consequences of among-individual variation in a trait of interest will reveal how selective pressures or ecological conditions influence that trait. However, many endocrine traits, such as concentrations of glucocorticoid (GC) hormones, exhibit adaptive plasticity and, therefore, do not necessarily respond to these pressures as predicted by among-individual phenotypic correlations. To improve our interpretations of among-individual variation in GC concentrations, we need more information about the repeatability of these traits within individuals. Many studies have already estimated the repeatability of baseline, stress-induced, and integrated GC measures, which provides an opportunity to use meta-analytic techniques to investigate 1) whether GC titers are generally repeatable across taxa, and 2) which biological or methodological factors may impact these estimates. From an intensive search of the literature, we collected 91 GC repeatability estimates from 47 studies. Overall, we found evidence that GC levels are repeatable, with mean repeatability estimates across studies ranging from 0.230 for baseline levels to 0.386 for stress-induced levels. We also noted several factors that predicted the magnitude of these estimates, including taxon, sampling season, and lab technique. Amphibians had significantly higher repeatability in baseline and stress-induced GCs than birds, mammals, reptiles, or bony fish. The repeatability of stress-induced GCs was higher when measured within, rather than across, life history stages. Finally, estimates of repeatability in stress-induced and integrated GC measures tended to be lower when GC concentrations were quantified using commercial kit assays rather than in-house assays. The extent to which among-individual variation in GCs may explain variation in organismal performance or fitness (and thereby inform our understanding of the ecological and evolutionary processes driving that variation) depends on whether measures of GC titers accurately reflect how individuals differ overall. Our findings suggest that while GC titers can reflect some degree of consistent differences among individuals, they frequently may not. We discuss how our findings contribute to interpretations of variation in GCs, and suggest routes for the design and analysis of future

research.

**Title:** Repeatability of glucocorticoid hormones in vertebrates: a meta-analysis
**Authors:** Kelsey L. Schoenemann[a], Frances Bonier[a]
**Affiliations:** [a]Department of Biology, Queen's University, Kingston, Ontario, Canada K7L 3N6
**Corresponding author:** Kelsey Schoenemann; kelsey.schoene@gmail.com

1  **Abstract**
2  We often expect that investigations of the patterns, causes, and consequences of among-
3  individual variation in a trait of interest will reveal how selective pressures or ecological
4  conditions influence that trait. However, many endocrine traits, such as concentrations of
5  glucocorticoid (GC) hormones, exhibit adaptive plasticity and, therefore, do not necessarily
6  respond to these pressures as predicted by among-individual phenotypic correlations. To improve
7  our interpretations of among-individual variation in GC concentrations, we need more
8  information about the repeatability of these traits within individuals. Many studies have already
9  estimated the repeatability of baseline, stress-induced, and integrated GC measures, which
10 provides an opportunity to use meta-analytic techniques to investigate 1) whether GC titers are
11 generally repeatable across taxa, and 2) which biological or methodological factors may impact
12 these estimates. From an intensive search of the literature, we collected 91 GC repeatability
13 estimates from 47 studies. Overall, we found evidence that GC levels are repeatable, with mean
14 repeatability estimates across studies ranging from 0.230 for baseline levels to 0.386 for stress-
15 induced levels. We also noted several factors that predicted the magnitude of these estimates,
16 including taxon, sampling season, and lab technique. Amphibians had significantly higher
17 repeatability in baseline and stress-induced GCs than birds, mammals, reptiles, or bony fish. The
18 repeatability of stress-induced GCs was higher when measured within, rather than across, life
19 history stages. Finally, estimates of repeatability in stress-induced and integrated GC measures
20 tended to be lower when GC concentrations were quantified using commercial kit assays rather
21 than in-house assays. The extent to which among-individual variation in GCs may explain
22 variation in organismal performance or fitness (and thereby inform our understanding of the
23 ecological and evolutionary processes driving that variation) depends on whether measures of
24 GC titers accurately reflect how individuals differ overall. Our findings suggest that while GC
25 titers can reflect some degree of consistent differences among individuals, they frequently may
26 not. We discuss how our findings contribute to interpretations of variation in GCs, and suggest
27 routes for the design and analysis of future research.
28
29 **Keywords:** glucocorticoid; cortisol; corticosterone; repeatability; heritability; intraclass
30 correlation coefficient; individual variation

## 1. Introduction

Since the development of immunoassays that allow the measurement of hormones in relatively small-volume tissue samples (Ekins, 1960; Yalow & Berson, 1960), the number of studies investigating the patterns, causes, and consequences of endocrine trait variation has soared. Early work in this field described variation in hormone concentrations across species, populations, and life history stages (e.g., Boswell et al., 1994; Klosterman et al., 1986; Pancak and Taylor, 1983), while more recent work often measures among-individual variation in multiple endocrine traits, including hormone concentration, receptor density, binding protein concentration, and endocrine axis responsiveness (e.g., Breuner et al., 2006; Bizon et al., 2001; Lattin & Romero, 2014; Liebl, Shimizu & Martin, 2013). Thus, much of our understanding of how selection has shaped these traits derives from comparative studies that determine how conserved or variable hormones, receptors, or their effects are across taxa, or how those traits vary among individuals with geography, phylogeny, or other traits of interest (e.g., Bókony et al., 2009; Eikenaar et al., 2014; Heidinger et al., 2006). Yet, traits that exhibit adaptive plasticity, such as hormone titers, might not respond to selective pressures or ecological conditions as predicted by among-individual phenotype-fitness correlations (Stinchcombe et al., 2002; Bonier et al., 2009; Bonier & Martin, 2016). Moving beyond this comparative approach to better understand endocrine trait evolution requires knowledge about heritable individual differences in evolutionarily-important traits because natural selection acts upon this heritable variation at the individual level (Bennett, 1987; Williams, 2008). However, the extent to which variation in hormone levels can be attributed to fixed individual differences is poorly understood.

Concentrations of glucocorticoid (GC) hormones, for example, exhibit plasticity, here defined as the ability of a single genotype to produce multiple phenotypes in response to

55   environmental changes, and referred to as flexibility in some contexts (*sensu* Bonier and Martin,

56   2016). The plasticity of GC titers helps organisms maintain allostasis, despite changing energetic

57   needs. Rapid secretion of GCs promotes behavioral and physiological changes that enable an

58   organism to respond to and recover from acute energetic challenges, while modulation of

59   baseline circulating GCs supports responses to predictable changes in energetic demands across

60   daily or seasonal cycles (Sapolsky, Romero & Munck, 2000; Romero, 2004; Wingfield, 2005;

61   Romero, Dickens & Cyr, 2009). Failure to acknowledge, measure, or control for these sources of

62   within-individual variation can diminish our ability to detect biologically significant patterns in

63   GC secretion among individuals.

64        Estimating the repeatability (i.e., consistency over time or across contexts) of GC titers is

65   one technique for assessing and potentially avoiding this pitfall. Multiple test statistics have been

66   used to estimate the repeatability of a trait in a population (e.g., Spearman rank and Pearson

67   correlation coefficients), but the intraclass correlation coefficient (ICC) is the most prevalent in

68   recent literature (Sokal and Rohlf 1995; Nakagawa and Schielzeth, 2010). The repeatability of

69   GCs within individuals can be used to determine the degree to which inferences made about GC

70   measures may be generalized beyond providing information about the individuals at the time of

71   sampling (e.g., Bosson et al., 2009; Harris et al., 2016; Wada et al., 2008). Moreover,

72   repeatability itself may reflect the ability or strategy of an individual to cope with a challenge

73   and, thus, is worthy of study in its own right (Careau, Buttemer & Buchanan, 2014; Roche,

74   Careau & Binning, 2016). Finally, estimates of repeatability can approximate the upper limit of

75   heritability of individual variation and, thereby, the extent to which natural selection can shape a

76   trait (Falconer and Mackay 1996; but see Dohm, 2002). Perhaps in recognition of these points,

77   many studies have estimated the repeatability of GC measures (e.g., Cook et al., 2012; Narayan

78  et al., 2013; Romero and Reed, 2008; Wada et al., 2008). The availability of these estimates

79  provides an opportunity to investigate whether GCs are generally repeatable across taxa, and

80  how biological or methodological factors may impact these estimates.

81      To date, researchers have estimated the repeatability of GC levels in every class of

82  vertebrates, and across various environmental contexts and spans of time. A meta-analysis of

83  repeatability estimates across these studies could determine whether GCs are generally

84  repeatable, and whether variation in the magnitude of repeatability can be explained by

85  biological or methodological factors. For example, meta-analyses of behavior and metabolic rate

86  repeatabilities have provided evidence of significant trait repeatability, as well as differences in

87  repeatability according to sex, sampling interval, captive condition, and taxon (Nespolo &

88  Franco, 2007; Bell, Hankison & Laskowski, 2009; White, Schimpf & Cassey, 2013). Here, we

89  similarly seek to investigate sources of variation in estimates of repeatability of GCs.

90  Specifically, the aim of this meta-analysis is to: 1) summarize the available evidence of

91  repeatability of GC concentrations; and 2) identify biological and methodological factors that

92  predict variation in the magnitude of GC repeatability.

93

94  **2. Methods**

95  *2.1 Literature Search*

96      We performed literature searches on Google Scholar between March 2016 and November

97  2017 using the terms: "repeatab*," "consisten*," "glucocorticoid," "cortisol", "corticoster*",

98  "repeated measure," and "individual variation." We identified 716 records in these searches. We

99  screened the titles and abstracts of these records, looking for papers that estimated the

100 repeatability (or 'consistency' or 'individuality') of concentrations of glucocorticoid hormones in

101    a variety of tissues (e.g., blood, saliva, feces, feathers). To be selected for inclusion in this

102    analysis, a study needed to have assessed repeated measurements from the same individual and

103    estimated a repeatability coefficient (e.g., Spearman rank, Pearson, or ICC). We excluded

104    duplicate and irrelevant articles and those that did not meet our inclusion criteria (Fig. 1). We

105    also checked reference lists of selected papers to find additional studies that were not identified

106    in the initial search. Lastly, we included 3 studies that collected repeated measurements of

107    hormone concentrations from the same individuals but did not estimate repeatability, when we

108    could obtain the original data to calculate repeatability.

109

110    *2.2 Repeatability Estimates*

111            We extracted repeatability estimates from the selected studies and categorized them as

112    representing either *initial*, *response*, or *integrated* GC repeatability measures. We used the

113    category *initial* to group repeatability estimates of GC titers measured in circulation within a

114    time period expected not to reflect the acute stress of capture, *response* for repeatability

115    estimates of the elevated GC titers following an acute capture, handling, or confinement stress,

116    and *integrated* for repeatability estimates of GC titers that represent hormone secretion over a

117    relatively long period of time (e.g., GC concentrations in feces, feathers, and saliva). If the study

118    did not calculate repeatability, then, where possible, we obtained the original data and calculated

119    an ICC repeatability, using the 'rptR' package (version: 0.9.2) in R (version 3.4.0, 2017-04-21)

120    (Nakagawa & Schielzeth, 2010).

121

122    *2.3 Statistical Analysis*

123    We harvested information about several methodological and biological factors associated

124    with each repeatability estimate and categorized these data for analysis (Table 1). We used linear

125    mixed-effect models (LMMs) with the 'lme4' package (version: 1.1.13) to investigate variation

126    in repeatability estimates. We included study identity as a random effect to control for potential

127    bias arising from non-independence of multiple estimates derived from the same study

128    (Nakagawa & Santos, 2012). One study, however, was coded with two independent study

129    identities because the datasets included in this one study were collected by two different

130    researchers, on different species, in different field sites (Ouyang, Hau & Bonier, 2011). We

131    constructed separate LMMs to address each of the following questions with *initial*, *response*, or

132    *integrated* GC repeatability measures:

133    1. *Does sampling regime predict repeatability?* To answer this question, we evaluated the

134       following fixed effects: sample size, average time span between samples, and average

135       number of samples.

136    2. *Does subject biology or sampling environment predict repeatability?* We evaluated the fixed

137       effects taxonomic class, sex, whether samples were collected within or across life history

138       stage, captive condition, and experimental manipulation (whether or not some/all individuals

139       underwent a stressful manipulation intended to produce a response [not including routine

140       capture and handling stress] at some point during the course of the study). We lacked

141       sufficient power to evaluate the effect of age because we identified only two estimates of

142       repeatability that were measured solely in juveniles or immature individuals. We also

143       evaluated the fixed effect of life history stage (breeding, non-breeding, or pre-breeding) in a

144       subset of GC repeatability estimates measured within a single stage.

145  **3.** *Do laboratory or statistical techniques predict repeatability?* We evaluated the fixed effects

146      use of an in-house assay or commercial assay kit, use of a radioactive or enzymatic tracer,

147      and whether or not the statistical analysis incorporated confounding factors (i.e., if the

148      repeatability estimate controlled for correlations between GCs and factors such as the time or

149      year of sampling, and the breeding status, age, or body mass of the individuals sampled).

150

151          With the exception of models that included sample size as a fixed factor (question 1,

152  above), we weighted each estimate by its sample size to account for differences in statistical

153  power among studies. Thus, estimates from larger studies had a greater influence in the models.

154  We verified the normality of model residuals with a Shapiro test. When model residuals failed to

155  meet the assumption of normality, we square-root transformed the data. To identify important

156  predictors of repeatability, we coded global models with all candidate variables included as main

157  effects and used the *dredge* function from the 'MuMIn' package (version: 1.15.6) to rank

158  recombinant models with the Akaike's information criterion corrected for small sample sizes

159  (AICc). We did not include any interaction terms in our models, due to small sample sizes. We

160  report effect size and p-values from either the best-fit model or, when more than one model was

161  ranked within 2 ΔAICc of the best-fit model, from a conditional average of all top models. Due

162  to the small sample size of *integrated* measures available to address question 2, we compared the

163  saturated model to a null model using an F-test with Kenward-Roger approximation using the

164  'pbkrtest' package (version: 0.4-7) (Kenward & Roger, 1997; Halekoh & Højsgaard, 2014). For

165  some non-ordinal variables (e.g., taxonomic class, sampling interval), it is more informative to

166  consider the significance of the factor as a whole rather than at specific levels; therefore, in such

167  cases, we performed a Type III ANOVA with Satterthwaite approximation for degrees of

168    freedom using the 'lmerTest' package (version: 2.0-33) to obtain p-values (Kuznetsova,

169    Brockhoff & Christensen, 2016).

170        In addition to including study identity as a random effect, we employed several other

171    methods to address potential bias or pseudo-replication. First, we did not include redundant

172    estimates from the same study nor re-analyses of the same data. Second, we assessed the

173    independence of multiple repeatability estimates originating from the same study. If a single GC

174    measure is correlated among multiple groups of individuals (e.g., similarly low *initial* GC

175    repeatability in males and females from same population), then we might expect multiple

176    repeatability estimates of the same population to be non-independent. To test for this effect, we

177    performed a linear regression analysis with those studies that reported more than one estimate to

178    test whether the number of estimates of repeatability in a study was associated with repeatability

179    (Nespolo & Franco, 2007; Bell, Hankison & Laskowski, 2009). We did not find a relationship

180    between *initial* repeatability and the number of estimates reported in the study (linear model:

181    initial $n = 37$, $p = 0.127$), and no studies of *integrated* repeatability reported more than two

182    estimates. We did find a significant negative relationship between the number of estimates of the

183    repeatability of *response* GCs and their magnitude ($n = 31$, $\beta = -0.10$, $p = 0.002$), however, this

184    relationship was driven by a single study that reported multiple estimates of 0.00 repeatability.

185    Thus, our inclusion of study identity as a random effect in all models was deemed sufficient to

186    control for non-independence of multiple estimates from the same study.

187        Finally, to determine whether GCs are generally repeatable across all studies, we first

188    needed to assess whether the estimates we obtained from the literature represent a random

189    sample of the 'true' repeatability of GC titers. Given that the primary focus of most studies

190    included in this analysis was not to estimate repeatability, we expect publication bias is unlikely

191  to be an important source of bias for our results. Nevertheless, we assessed this and other

192  potential biases directly by plotting every estimate against its sample size in funnel plots. Upon

193  finding these plots symmetrical (Supplemental Fig. 1), we concluded that bias is unlikely (Egger

194  et al., 1997). Therefore, we calculated 95% confidence intervals around the mean repeatabilities

195  of *initial*, *response*, and *integrated* measures across all studies, regardless of taxon, using 1000

196  bootstrap samples of the data with replacement. We interpret a confidence interval that does not

197  overlap zero as indicating that the mean GC repeatability estimate is greater than zero (i.e., the

198  GC measure is, on average, somewhat repeatable), and interpret confidence intervals that do not

199  overlap each other as indicating different repeatabilities.

200

201  **3. Results**

202  *3.1 Summary of the data set*

203       We identified 47 studies that met our criteria for inclusion, from which we extracted 91

204  estimates of GC repeatability (summarized in Table 2, see Supplementary Information for

205  complete dataset). In brief, more estimates were made of *initial* or *response* measures than of

206  *integrated* measures. The repeatability estimates included data from 36 species; however, more

207  than two-thirds of the estimates originated from studies of birds. Free-ranging populations of

208  adults with both sexes combined were more often studied than captive populations, juveniles or

209  immatures, or separately for the sexes. About three-quarters of the estimates spanned a sampling

210  interval of less than one year. The majority of estimates came from repeated measurement within

211  the same life history stage and, of those measured within a stage, more were derived from

212  measurements taken during the breeding season. Finally, the ICC was the most common

213  repeatability estimate reported, with 42 studies reporting an ICC and only 4 reporting either

214    Pearson or Spearman correlations; in one study, the authors did not clearly report method used

215    nor respond to our requests for information.

216

217    **3.2 Repeatability of GCs**

218         Overall, GC levels were moderately repeatable, with mean repeatabilities ranging from

219    0.230 for *initial* measures, 0.320 for *integrated* measures, and 0.386 for *response* measures (Fig.

220    2). Moreover, the 95% confidence intervals around the mean repeatability of all three types of

221    measures did not overlap zero (initial: 0.230 [0.162, 0.294], response: 0.386 [0.318, 0.449],

222    integrated: 0.320 [0.235, 0.410]). As indicated by non-overlapping confidence intervals, the

223    mean repeatability of *response* measures were greater than those of *initial* measures.

224

225    **3.3 Relationships between repeatability and biological or methodological factors**

226    *3.3.1 Does sampling regime predict repeatability?*

227         We found little evidence that sample size, time span between samples, or number of

228    samples predicts GC repeatability. The null was the best-fit model for *integrated* measures and,

229    while number of measurements and sample size were retained in top models of *initial* and

230    *response* measures (Supplementary Table 1), we did not find evidence that *initial* or *response*

231    repeatability varied significantly with these factors (model average: all $p > 0.12$). Sampling

232    interval, however, was retained in top models of *response* measures and, on average,

233    repeatability was greater when repeated measurements were collected within 8-14 days of each

234    other (0.607, $n = 8$), compared to either shorter (0-7 days; 0.327, $n = 5$) or longer (15-365+ days;

235    0.324, $n = 24$) intervals (Type III ANOVA; $n = 37$, $F_{(5,35)} = 2.840$, $p = 0.030$).

236

237    *3.3.2 Does subject biology or sampling environment predict repeatability?*

238         Taxonomic class was retained in the top models explaining variation in repeatability

239    estimates for both *initial* and *response* measures (Supplemental Table 2). On average,

240    amphibians had higher *initial* and *response* repeatability (0.833, *n* = 4; 0.786, *n* = 4, respectively)

241    than birds (0.162, *n* = 35; 0.318, *n* = 21), mammals ([no *initial* GC repeatability estimates in

242    mammals]; 0.446, *n* = 5), reptiles (0.270, *n* = 1; 0.21, *n* = 2), or fish (0.201, *n* = 2; 0.359, *n* = 5)

243    (Fig. 3; Type III ANOVA; initial*: n* = 38, *F(3,38)* = 9.359, *p* < 0.0001; response: *n* = 27,

244    *F(4,23)* = 4.984, *p* = 0.005). While sex was retained in the top models of *initial* measures, we did

245    not find strong evidence that repeatabilities varied by sex (model average: all *p* > 0.15).

246    Estimates of *response* repeatability were higher when derived from measurements within a life

247    history stage (0.502, *n* = 22) than when derived from measurements across stages (0.072, *n* = 5)

248    (Supplemental Table 3; model average: *n* = 27, *β* = 0.235, *p* = 0.007). Neither experimental

249    manipulation nor captive condition was retained in any top models. The global model evaluating

250    *integrated* measures was not better-fit than the null (F-test: *n* = 10, *F(7,3023)* = 0.191,

251    *p* = 0.988).

252         Finally, in the subset analyses of repeatability estimates measured within a life history

253    stage, we found little evidence that life history stage (breeding, non-breeding, or pre-breeding)

254    predicts repeatability. The null model was the best-fit model for *initial* and *response* measures

255    (Supplemental Table 2). However, a univariate model including life history stage performed

256    better than the null for *integrated* measures, where repeatability was on average higher in the

257    non-breeding season (0.555, *n* = 3) compared to breeding (0.266, *n* = 5; F-test: *n* = 8, *F(1,2370)*

258    = 10.7, *p* = 0.001).

259

260    *3.3.3 Do laboratory or statistical techniques predict repeatability?*

261    Assay type (in-house or kit) was retained in top models of *initial*, *response*, and

262    *integrated* measures, while assay tracer was retained in the top models of *initial* and *integrated*

263    measures (Supplemental Table 4). Repeatabilities of *initial* and *integrated* hormone

264    concentrations measured with RIA were lower than those measured with EIA, although this

265    difference was not as evident for *initial* measures (Supplemental Table 5; model average *initial*:

266    $n = 40$, $\beta = -0.132$, $p = 0.071$; *integrated*: $n = 11$, $\beta = -0.194$, $p = 0.024$). In addition, the

267    repeatabilities of *response* measures were lower when measured with a kit than those measured

268    with an in-house assay, and tended to be lower for repeatability of *integrated* measures

269    (Supplemental Table 5; model average: *response*: $n = 35$, $\beta = -0.184$, $p = 0.040$; *integrated*:

270    $n = 11$, $\beta = -0.172$, $p = 0.062$). Finally, whether or not confounding factors were controlled was

271    retained in one top model of *response* measures, however, we did not find strong evidence that

272    repeatability varied with this factor (Supplemental Table 5; model average: $n = 35$, $\beta = 0.101$,

273    $p = 0.340$).

274

## Discussion

276    To better understand individual variation in GCs, we summarized published estimates of

277    GC repeatability and identified factors that predicted the magnitude of those estimates. We found

278    measures of *initial*, *response*, and *integrated* GCs had mean repeatabilities of 0.230, 0.386, and

279    0.320, respectively, with *response* repeatability estimates greater than *initial* repeatability. In

280    general, this finding suggests that measures of GC titers reflect a moderate degree of consistent

281    differences among individuals, however, some measures were more or less repeatable, depending

282    on how the biological sample was collected and analyzed or which individuals were sampled.

283    Specifically, we found that some estimates of GC repeatability were greater in amphibians, when

284    all samples from an individual were collected within a single life history stage, and when

285    samples collected within a life history stage came from the non-breeding season. We also found

286    some evidence that GC repeatability was greater when hormone concentrations were measured

287    using an in-house immunoassay, with an enzyme tracer, and when repeated measurements of the

288    same individuals were collected across a relatively short time span (i.e., a sampling interval of 8-

289    14 days).

290         The repeatability of GCs within individuals can be used to: 1) determine whether

291    inferences made about GC measures may be generalized beyond the time of sampling (e.g.,

292    Bosson et al., 2009; Harris et al., 2016; Wada et al., 2008), 2) describe the ability or strategy of

293    an individual to cope with a challenge (Careau, Buttemer & Buchanan, 2014; Roche, Careau &

294    Binning, 2016), and 3) approximate the upper limit of heritability of individual variation and,

295    thereby, the extent to which natural selection can shape a trait (Falconer and Mackay 1996; but

296    see Dohm, 2002). Below, we interpret our findings in light of each of these applications of

297    estimates of repeatability.

298         While we found that some measures of GCs were highly repeatable (i.e., >0.70; see

299    Angelier et al., 2010; Ferrari et al., 2013; and Narayan et al., 2013b) and, therefore, expected to

300    be reliable indicators of an individual's endocrine phenotype beyond the period of sampling,

301    many other measures were not. Low repeatability may be caused by high within-individual

302    variation, high measurement error, low among-individual variation, or a combination of all three.

303    Whether a population exhibits low repeatability due to high within-individual variation (rather

304    than low among-individual variation), or due to variation in trait consistency among individuals

305    has different implications for how to collect and interpret data from that population of

306    individuals (Jenkins, 2011; Biro & Stamps, 2015). When high within-individual variation is a

307    concern, a single measurement of GCs will best capture individual differences when collected

308    from all individuals instantaneously or while controlling for as many sources of environmental

309    variation as possible. In the case of variation among individuals in trait consistency, a single

310    measure of GCs will be unlikely to capture how individuals differ overall.

311           Whether or not an endocrine trait is repeatable for a given population, if individuals are

312    sampled across different physical or social environments, or if some individuals differ in

313    personality-related strategies, then the within-individual relationship between hormones and

314    another variable of interest can differ from the population-level response in unexpected ways

315    (Roche et al., 2016). For example, while a study found no relationship between brood size and

316    baseline GCs among female tree swallows (*Tachycineta bicolor*), baseline GCs increased within

317    individuals following an experimental increase in brood size (Bonier, Moore & Robertson,

318    2011). Additionally, olive flounder (*Paralichthys olivaceus*) with bold behavioral phenotypes

319    responded physiologically to an acute stress in a manner opposite that of shy types, and these

320    divergent responses were repeatable (Rupia et al., 2016). In both of these cases, failure to

321    measure within-individual changes in GCs, or to recognize among-individual variation in the

322    direction of those responses, would have obscured detection of the effects of the challenge of

323    interest (i.e., brood size, acute stress) at the population level. Our finding of relatively low GC

324    repeatability, particularly for *initial* GCs, strongly suggests that these measures frequently reflect

325    an individual's short-term response to the environment more so than fixed differences among

326    individuals.

327           Variation in GC repeatability can also be used to investigate differences in the ability or

328    strategy of individuals or populations to respond to environmental change. For example, our

329    finding of significantly greater repeatability in *response*, compared to *initial*, measures could

330    indicate relatively greater canalization in the acute activation of the HPA axis, and a reduced

331    plasticity of this trait within individuals. Consistent with this interpretation, previous studies have

332    estimated greater realized heritability of the GC response in genetic lines selected for high, rather

333    than low, stress responses (Brown & Nestor, 1973; Satterlee & Johnson, 1988). Additionally, the

334    greater repeatability of both *initial* and *response* GCs in amphibians could indicate different

335    functions and/or responsiveness of the HPA axis in amphibians compared to other taxonomic

336    classes (Narayan et al., 2013a). Finally, our finding greater repeatability of *response*, but not

337    *initial*, GCs measured within a life history stage somewhat aligns with previous work, which has

338    shown greater seasonal variation in baseline, rather than stress-induced, GC titers (Romero,

339    2002). And although our sample size was small ($n = 8$), our finding of greater repeatability of

340    *integrated* GC measures during the non-breeding season seems to suggest less variation within

341    individuals in the total secretion of GCs during that period, which could reflect a broader pattern

342    of seasonal GC secretion across taxa.

343         If one aims to compare repeatability or trait consistency among individuals, populations,

344    or even species, as described above, then an important consideration is whether variation among

345    repeatability estimates is due to laboratory or statistical methodologies impacting within- or

346    among-individual variation in the trait of interest. We found that some repeatability estimates

347    were lower when measured with a commercial kit compared to an in-house assay, and when

348    measured with an RIA as compared to an EIA. Commercial assay kits can be less precise (as

349    well as less accurate) in measuring GC concentrations if they are not carefully validated for the

350    study system (Buchanan & Goldsmith, 2004; Sheriff et al., 2011), which may explain lower

351    repeatability estimates for GCs measured with kits. Further, the ease of use of commercial kits

352    might lend itself to less precise lab practices than the more involved in-house assays. However, it

353    is not clear why RIAs would be associated with lower repeatability. Brown et al. (2010) found

354    that, while urinary cortisol assessed with either RIA or EIA exhibited qualitatively-similar

355    temporal profiles, the RIA detected proportionally lower hormone concentrations (i.e., decreased

356    among-individual variation) (Brown et al., 2010). This lower among-individual variation could

357    lead to lower repeatability, if it is not counteracted by simultaneously lower within-individual

358    variation. Previous work has documented large inter-laboratory variation in measurements of

359    absolute steroid hormone concentrations (Bókony et al., 2009; Fanson et al., 2017; Feswick et

360    al., 2014; Ganswindt et al., 2012), indicating that across-study comparisons of absolute values of

361    individuals' GC titers are not valid. Finally, while we also found some evidence that *response*

362    GC repeatability was greater when repeated measurements were collected over a relatively short

363    time span (i.e., 8-14 days apart), even shorter time spans did not show a consistent pattern, and

364    we did not detect a similar effect in any of the other GC measures. Overall, if one seeks to

365    investigate the causes and consequences of variable GC repeatability among groups, to better

366    understand the ability or strategy of these groups to respond to environmental conditions,

367    methodological sources of variation must be considered and, ideally, controlled.

368          A final application of estimates of trait repeatability is to approximate the upper limit of

369    heritability. The average repeatability of *initial* and *response* GCs reported here align well with

370    the results of artificial selection and animal model approaches that estimate a similar degree of

371    heritability in GC titers and the GC response (Evans et al., 2006; Jenkins et al., 2014; Pottinger

372    & Carrick, 1999; Touma et al., 2008). These studies often find that the heritability of baseline

373    GCs is much lower than response GCs, if it is detectable at all (e.g., Brown & Nestor, 1973;

374    Satterlee & Johnson, 1988; Evans et al., 2006). Thus, we expect baseline concentrations will be

375    less likely to exhibit evolutionary change than stress-induced concentrations, when exposed to

376  similar selective pressures. Furthermore, Jenkins et al. (2014) failed to find phenotypic or genetic

377  correlations between baseline and stress-induced concentrations within individuals. This finding

378  suggests that different mechanisms may control GC secretion during normal activity versus

379  during challenging events, and that selection could affect variation in these traits independently

380  (Jenkins et al., 2014). As a result, selective or ecological pressures should be expected to produce

381  complex, context-dependent relationships between hormone titers and factors of interest.

382  Altogether, the low-to-moderate repeatability and heritability of GC titers underscores the extent

383  to which plasticity may generate individual variation, as well as the extent to which that variation

384  may be transmitted to future generations.

385          While our meta-analysis of GC repeatability estimates allowed us to look for patterns in

386  trait consistency across a range of methodological and biological factors, there are limitations to

387  our dataset and thus our ability to draw strong inferences from it. For example, many studies

388  calculated repeatability as a way to compliment or support their main results. If researchers are

389  more likely to report repeatability estimates that support their main findings, then repeatability

390  estimates available in the literature may overestimate true repeatability. In addition, our

391  categorization of the biological and methodological data associated with each repeatability

392  estimate could have over-simplified or otherwise misrepresented the reality of the study, which

393  could make real patterns more difficult to detect, or possibly cause spurious patterns (e.g., among

394  the more weakly-supported findings). Finally, sample size was limited for many categories

395  included in our analyses, thereby reducing our statistical power to detect real patterns.

396

397  **Conclusion**

398    Overall, this meta-analysis provides new insights into individual variation in GC titers,

399    and highlights the importance of repeatability estimation to improve methods for collecting and

400    interpreting biological data. We found that GCs were moderately repeatable, on average, but

401    these estimates were also highly variable. Additionally, *initial* and *response* GC measures were

402    more repeatable in amphibians than any other taxonomic class, while *response* GCs were more

403    repeatable when measured within the same life history stage and *integrated* GC were more

404    repeatable during the non-breeding season. We look forward to new research that further

405    investigates how and why repeatability differs with these factors. However, our finding that

406    laboratory techniques were also associated with variation in repeatability could serve as a

407    reminder to be meticulous in monitoring for issues with the reproducibility of hormone data.

408    Moving forward, a better understanding of endocrine trait evolution requires knowledge about

409    heritable individual differences in evolutionarily-important traits. Our analysis shows that a

410    single measure of individual variation in GC titers may not reflect how those individuals differ in

411    general, and suggests different approaches to capture that signal, including repeated

412    measurements of individuals both within and across environments.

413
414 **Acknowledgments**
415 We thank Y. Aharon-Rotman, A. Gladbach, B. Dantzer, J. Riechert, C. Vleck, H. Wada, and S.
416 Winberg for providing data or additional information about their published studies included in
417 the analyses. We also thank R. Montgomerie for advice on the statistical analyses.
418
419 **Funding**
420 Funding from an NSERC Discovery Grant and Queen's University supported KS's graduate
421 stipend.

## References

Angelier F., Wingfield JC., Weimerskirch H., Chastel O. 2010. Hormonal correlates of
        individual quality in a long-lived bird: a test of the "corticosterone-fitness hypothesis".
        *Biology letters* 6:846–849. DOI: 10.1098/rsbl.2010.0376.

Bell AM., Hankison SJ., Laskowski KL. 2009. The repeatability of behaviour: a meta-analysis.
        *Animal Behaviour* 77:771–783. DOI: 10.1016/j.anbehav.2008.12.022.

Bennett AF. 1987. Inter-individual variability: an under-utilized resource. In: Feder ME, Bennett
        AF, Burggren WW & Huey RB, eds. *New directions in ecological physiology*. Cambridge
        University Press: Cambridge, 147–169

Biro PA., Stamps JA. 2015. Using repeatability to study physiological and behavioural traits:
        ignore time-related change at your peril. *Animal Behaviour* 105:223–230. DOI:
        10.1016/j.anbehav.2015.04.008.

Bizon JL., Helm KA., Han JS., Chun HJ., Pucilowska J., Lund PK., & Gallagher M. (2001).
        Hypothalamic–pituitary–adrenal axis function and corticosterone receptor expression in
        behaviourally characterized young and aged Long–Evans rats. *European Journal of
        Neuroscience 14*:1739-1751.

Bókony V., Lendvai ÁZ., Liker A., Angelier F., Wingfield JC., Chastel O. 2009. Stress response
        and the value of reproduction: are birds prudent parents? *The American naturalist* 173:589–
        598. DOI: 10.1086/597610.

Bonier F., Martin PR. 2016. How can we estimate natural selection on endocrine traits? Lessons
        from evolutionary biology. *Proceedings of the Royal Society B* 283:20161887. DOI:
        10.1098/rspb.2016.1887.

Bonier F., Martin PR., Moore IT., Wingfield JC. 2009. Do baseline glucocorticoids predict
        fitness? *Trends in Ecology and Evolution* 24:634–642. DOI: 10.1016/j.tree.2009.04.013.

Bonier F., Moore IT., Robertson RJ. 2011. The stress of parenthood? Increased glucocorticoids
        in birds with experimentally enlarged broods. *Biology Letters* 7:944–946. DOI:
        10.1098/rsbl.2011.0391.

Bosson CO., Palme R., Boonstra R. 2009. Assessment of the stress response in Columbian
        ground squirrels: laboratory and field validation of an enzyme immunoassay for fecal
        cortisol metabolites. *Physiological and biochemical zoology* 82:291–301. DOI:
        10.1086/597530.

Boswell T., Woods SC., Kenagy GJ. 1994. Seasonal changes in body mass, insulin, and
        glucocorticoids of free-living golden-mantled ground squirrels. *General and comparative
        endocrinology* 96:339–346. DOI: 10.1006/gcen.1994.1189.

Breuner CW., Lynn SE., Julian GE., Cornelius JM., Heidinger BJ., Love OP., Sprague RS.,
        Wada H., Whitman BA. (2006). Plasma-binding globulins and acute stress
        response. *Hormone and Metabolic Research*, *38*:260-268.

Brown JL., Kersey DC., Freeman EW., Wagener T. 2010. Assessment of diurnal urinary cortisol
        excretion in Asian and African elephants using different endocrine methods. *Zoo Biology*
        29:274–283. DOI: 10.1002/zoo.20268.

Brown KI., Nestor KE. 1973. Some physiological responses of turkeys selected for high and low
        adrenal response to cold stress . *Poultry science*  52:1948.

Buchanan KL., Goldsmith AR. 2004. Noninvasive endocrine data for behavioural studies: The
        importance of validation. *Animal Behaviour* 67:183–185. DOI:
        10.1016/j.anbehav.2003.09.002.

Careau V., Buttemer WA., Buchanan KL. 2014. Early-developmental stress, repeatability, and

468        canalization in a suite of physiological and behavioral traits in female zebra finches.
469        *Integrative and comparative biology* 54:539–554. DOI: 10.1093/icb/icu095.
470   Cook K V., O'Connor CM., McConnachie SH., Gilmour KM., Cooke SJ. 2012. Condition
471        dependent intra-individual repeatability of stress-induced cortisol in a freshwater fish.
472        *Comparative Biochemistry and Physiology, Part A* 161:337–343. DOI:
473        10.1016/j.cbpa.2011.12.002.
474   Dohm MR. 2002. Repeatability estimates do not always set an upper limit to heritability.
475        *Functional Ecology* 16:273–280. DOI: 10.1046/j.1365-2435.2002.00621.x.
476   Donham RS. 1979. Annual cycle of plasma luteinizing hormone and sex hormones in male and
477        female mallards (Anas platyrhynchos). *Biology of Reproduction* 21:1273-1285.
478   Egger M., Smith GD., Schneider M., Minder C. 1997. Bias in meta-analysis detected by a simple
479        , graphical test measures of funnel plot asymmetry. *Bmj* 315:629–34. DOI:
480        10.1136/bmj.315.7109.629.
481   Eikenaar C., Klinner T., Stöwe M. 2014. Corticosterone predicts nocturnal restlessness in a long-
482        distance migrant. *Hormones and Behavior* 66:324–329. DOI: 10.1016/j.yhbeh.2014.06.013.
483   Ekins RP. 1960. The estimation of thyroxine in human plasma by an electrophoretic technique.
484        *Clinica chimica acta; international journal of clinical chemistry* 5:453–459. DOI:
485        10.1016/0009-8981(60)90051-6.
486   Evans MR., Roberts ML., Buchanan KL., Goldsmith AR. 2006. Heritability of corticosterone
487        response and changes in life history traits during selection in the zebra finch. *Journal of
488        Evolutionary Biology* 19:343–352. DOI: 10.1111/j.1420-9101.2005.01034.x.
489   Falconer DS., Mackay TF. 1996. *Introduction to Quantitative Genetics*, 4th edn. Longman,
490        Harlow.
491   Ferrari C., Pasquaretta C., Carere C., Cavallone E., von Hardenberg A., Réale D. 2013. Testing
492        for the presence of coping styles in a wild mammal. *Animal Behaviour* 85:1385–1396. DOI:
493        10.1016/j.anbehav.2013.03.030.
494   Halekoh U., Højsgaard S. 2014. A Kenward-Roger Approximation and Parametric Bootstrap
495        Methods for Tests in Linear Mixed Models - The R Package pbkrtest. *Journal of Statistical
496        Software* 59:1–32. DOI: 10.18637/jss.v059.i09.
497   Harris CM., Madliger CL., Love OP. 2016. Temporal overlap and repeatability of feather
498        corticosterone levels: practical considerations for use as a biomarker. *Conservation
499        Physiology* 4:1–11. DOI: 10.1093/conphys/cow051.
500   Heidinger BJ., Nisbet ICT., Ketterson ED. 2006. Older parents are less responsive to a stressor in
501        a long-lived seabird: a mechanism for increased reproductive performance with age?
502        *Proceedings of the Royal Society B: Biological Sciences* 273:2227–2231. DOI:
503        10.1098/rspb.2006.3557.
504   Jenkins SH. 2011. Sex differences in repeatability of food-hoarding behaviour of kangaroo rats.
505        *Animal Behaviour* 81:1155–1162. DOI: 10.1016/j.anbehav.2011.02.021.
506   Jenkins BR., Vitousek MN., Hubbard JK., Safran RJ. 2014. An experimental analysis of the
507        heritability of variation in glucocorticoid concentrations in a wild avian population.
508        *Proceedings of the Royal Society B* 281:20141302. DOI: 10.1098/rspb.2014.1302.
509   Lattin CR., & Romero LM. (2014). Chronic stress alters concentrations of corticosterone
510        receptors in a tissue-specific manner in wild house sparrows (Passer domesticus). *Journal of
511        Experimental Biology* 217:2601-2608.
512   Liebl AL., Shimizu T., & Martin LB. (2013). Covariation among glucocorticoid regulatory
513        elements varies seasonally in house sparrows. *General and Comparative Endocrinology*

514      *183*:32-37.
515  Kenward MG., Roger JH. 1997. Small Sample Inference for Fixed Effects from Restricted
516      Maximum Likelihood. *Biometrics* 53:983. DOI: 10.2307/2533558.
517  Klosterman LL., Murai JT., Siiteri PK. 1986. Cortisol levels, binding, and properties of
518      corticosteroid-binding globulin in the serum of primates. *Endocrinology* 118:424–434. DOI:
519      10.1210/endo-118-1-424.
520  Künzl C., Sachser N. 1999. The behavioral endocrinology of domestication: a comparison
521      between the domestic guinea pig (Cavia apereaf. porcellus) and its wild ancestor, the
522      cavy (Cavia aperea). *Hormones and Behavior*, 35:28-37.
523  Kuznetsova A., Brockhoff PB., Christensen RHB. 2016. lmerTest: Tests in Linear Mixed Effects
524      Models. R package version 2.0-33. https://CRAN.R-project.org/package = lmerTest
525  Nakagawa S., Santos ESA. 2012. Methodological issues and advances in biological meta-
526      analysis. *Evolutionary Ecology* 26:1253–1274. DOI: 10.1007/s10682-012-9555-5.
527  Nakagawa S., Schielzeth H. 2010. Repeatability for Gaussian and non-Gaussian data: a practical
528      guide for biologists. *Biological Reviews* 85:935–956. DOI: 10.1111/j.1469-
529      185X.2010.00141.x.
530  Narayan EJ., Cockrem JF., Hero JM. 2013a. Are baseline and short-term corticosterone stress
531      responses in free-living amphibians repeatable? *Comparative Biochemistry and Physiology,
532      Part A* 164:21–28. DOI: 10.1016/j.cbpa.2012.10.001.
533  Narayan EJ., Cockrem JF., Hero JM. 2013b. Repeatability of baseline corticosterone and short-
534      term corticosterone stress responses, and their correlation with testosterone and body
535      condition in a terrestrial breeding anuran (Platymantis vitiana). *Comparative Biochemistry
536      and Physiology - A Molecular and Integrative Physiology* 165:304–312. DOI:
537      10.1016/j.cbpa.2013.03.033.
538  Nespolo RF., Franco M. 2007. Whole-animal metabolic rate is a repeatable trait: a meta-analysis.
539      *The Journal of Experimental Biology* 210:2000–5. DOI: 10.1242/jeb.02780.
540  Ouyang JQ., Hau M., Bonier F. 2011. Within seasons and among years: When are corticosterone
541      levels repeatable? *Hormones and Behavior* 60:559–564. DOI:
542      10.1016/j.yhbeh.2011.08.004.
543  Pancak MK., Taylor H. 1983. Seasonal and daily plasma corticosterone rhythms in American
544      toads, Bufo americanus. *General and Comparative Endocrinology* 50:490–497. DOI: 001-
545      6480/83.
546  Pottinger TG., Carrick TR. 1999. Modification of the Plasma Cortisol Response to Stress in
547      Rainbow Trout by Selective Breeding. *General and Comparative Endocrinology* 116:122–
548      132. DOI: http://dx.doi.org/10.1006/gcen.1999.7355.
549  Roche DG., Careau V., Binning SA. 2016. Demystifying animal "personality" (or not): why
550      individual variation matters to experimental biologists. *Journal of Experimental Biology*.
551      DOI: 10.1242/jeb.146712.
552  Romero LM. 2004. Physiological stress in ecology: lessons from biomedical research. *Trends in
553      Ecology and Evolution* 19:249–255. DOI: 10.1016/j.tree.2004.03.008.
554  Romero LM., Dickens MJ., Cyr NE. 2009. The reactive scope model: a new model integrating
555      homeostasis, allostasis, and stress. *Hormones and Behavior* 55:375–389. DOI:
556      10.1016/j.yhbeh.2008.12.009.
557  Romero LM., Reed JM. 2008. Repeatability of baseline corticosterone concentrations. *General
558      and Comparative Endocrinology* 156:27–33. DOI: 10.1016/j.ygcen.2007.10.001.
559  Sapolsky RM., Romero LM., Munck AU. 2000. How do glucocorticoids influence stress

560      responses? Integrating suppressive, stimulatory, and preparative actions. *Endocrine Reviews*
561      21:55–89. DOI: 10.1210/er.21.1.55.

562    Satterlee DG., Johnson WA. 1988. Selection of Japanese quail for contrasting blood
563      corticosterone response to immobilization. *Poultry Science* 67:25–32. DOI:
564      10.3382/ps.0670025.

565    Sheriff MJ., Dantzer B., Delehanty B., Palme R., Boonstra R. 2011. Measuring stress in wildlife:
566      Techniques for quantifying glucocorticoids. *Oecologia* 166:869–887. DOI:
567      10.1007/s00442-011-1943-y.

568    Sokal RR., Rohlf FJ. 1995. *Biometry: The principles and practice of statistics in biological*
569      *research*, 3rd edn.W.H. Freeman and Company, New York.

570    Sossinka R. 1982. Domestication in birds. In: Farner, DS, King JR, Parkes KC, eds. *Avian*
571      *biology*. Academic Press, 373-403

572    Stinchcombe JR., Rutter MT., Burdick DS., Tiffin P., Rausher MD., Mauricio R. 2002. Testing
573      for Environmentally Induced Bias in Phenotypic Estimates of Natural Selection: Theory and
574      Practice. *The American Naturalist* 160:511–523. DOI: 10.1086/342069.

575    Touma C., Bunck M., Glasl L., Nussbaumer M., Palme R., Stein H., Wolferstätter M., Zeh R.,
576      Zimbelmann M., Holsboer F., Landgraf R. 2008. Mice selected for high versus low stress
577      reactivity: a new animal model for affective disorders. *Psychoneuroendocrinology* 33:839–
578      862. DOI: 10.1016/j.psyneuen.2008.03.013.

579    Wada H., Salvante KG., Stables C., Wagner E., Williams TD., Breuner CW. 2008.
580      Adrenocortical responses in zebra finches (Taeniopygia guttata): individual variation,
581      repeatability, and relationship to phenotypic quality. *Hormones and Behavior* 53:472–480.
582      DOI: 10.1016/j.yhbeh.2007.11.018.

583    White CR., Schimpf NG., Cassey P. 2013. The repeatability of metabolic rate declines with time.
584      *Journal of Experimental Biology* 216:1763–1765. DOI: 10.1242/jeb.076562.

585    Williams TD. 2008. Individual variation in endocrine systems: moving beyond the "tyranny of
586      the Golden Mean." *Philosophical Transactions of the Royal Society B: Biological Sciences*
587      363:1687–1698. DOI: 10.1098/rstb.2007.0003.

588    Wingfield JC. 2005. The concept of allostasis: coping with a capricious environment. *Journal of*
589      *Mammalogy* 86:248–254. DOI: 10.1644/BHE-004.1.

590    Yalow RS., Berson SA. 1960. Immunoassay of endogneous plasma insulin in man. *Journal of*
591      *Clinical Investigation* 39:1157–1175. DOI: 10.1172/JCI104130.

592

**Table 1**(on next page)

Table 1

Table 1. List describing how methodological and biological factors associated with each repeatability estimate were categorized for analysis.

1    Table 1. List describing how methodological and biological factors associated with each
2    repeatability estimate were categorized for analysis.
3

| FACTOR | CATEGORIES |
|---|---|
| *Time between measurements[1]* | 0-7d, 8-14d, 15-30d, 31-90d, 91-195d, or 365+ |
| *Number of measurements[1]* | Two, more than 2 |
| *Captive condition* | Free-ranging, captive, wild-caught captive |
| *Taxonomic class* | Bird, mammal, amphibian, bony fish, reptile |
| *Age* | Adult, juvenile, both |
| *Sex* | Male, female, both |
| *Life history stage (LHS)* | Breeding, non-breeding, pre-breeding, NA[2] |
| *Measured within LHS* | Yes, No |
| *Assay source* | In-house, commercial kit |
| *Assay tracer* | Radioactive, enzymatic |
| *Experimental manipulation[3]* | Yes, No |
| *Adjusted[4]* | Yes, No |

4    [1]Average, weighted by number of individuals when possible
5    [2]We categorized life history stage as "NA" for domesticated or captive-born species because
6    domestication can alter seasonal patterns in hormone physiology (Donham, 1979; Sossinka, 1982; Künzl
7    & Sachser, 1999). Estimates from these species were not included in analyses that examined the effect of
8    life history stage.
9    [3]Experimental manipulation refers to studies in which some or all individuals underwent a stressful
10   manipulation intended to produce a response (not including routine capture and handling stress) at some
11   point during the course of the study.
12   [4]Adjusted refers to whether or not estimates reflect GC repeatability after statistically controlling for
13   factors expected to explain some of the variation in GC titers (e.g., year, sex, weather).
14
15
16

# Table 2**(on next page)**

Table 2

Table 2. Summary of the data included in the meta-analysis. Except for sample size, numbers provided reflect the number of estimates in each category.

1  Table 2. Summary of the data included in the meta-analysis. Except for sample size, numbers
2  provided reflect the number of estimates in each category.

| GC measure | Initial[1] | Response[2] | Integrated[3] | | | |
|---|---|---|---|---|---|---|
| | 42 | 37 | 12 | | | |
| Sample size | Mean | Range | | | | |
| | 36 ± SE 4.5 | 8 - 352 | | | | |
| Sampling interval | 0-7d | 8-14d | 15-30d | 31-90d | 91-195d | 365+d |
| | 13 | 26 | 8 | 17 | 4 | 23 |
| Number of measurements | 2 | >2 | | | | |
| | 39 | 52 | | | | |
| Captive condition | Free-ranging | Captive-born | Wild-caught captive | | | |
| | 58 | 14 | 19 | | | |
| Taxonomic class | Bird | Mammal | Amphibian | Bony fish | Reptile | |
| | 60 | 11 | 8 | 9 | 3 | |
| Age | Adult | Juvenile | Both | | | |
| | 80 | 2 | 9 | | | |
| Sex | Male | Female | Both | | | |
| | 18 | 30 | 43 | | | |
| Life history stage (LHS)[4] | Breeding | Non-breeding | Pre-breeding | NA | | |
| | 36 | 21 | 9 | 25 | | |
| Within LHS | Y | N | NA[4] | | | |
| | 64 | 11 | 16 | | | |
| Assay source | In-house | Kit-based | | | | |
| | 51 | 35 | | | | |
| Assay tracer | Radioactive | Enzyme | | | | |
| | 42 | 44 | | | | |
| Experimental manipulation[5] | Y | N | | | | |
| | 21 | 70 | | | | |
| Adjusted[6] | Y | N | | | | |
| | 21 | 70 | | | | |

3
4  [1]Initial GCs refer to concentrations of GCs expected not to reflect the acute stress of capture.
5  [2]Response GCs refer to elevated GC titers following an acute capture, handling, or confinement stress.
6  [3]Integrated GCs refer to GC titers representing hormone secretion over a relatively long time.
7  [4] We categorized life history stage as "NA" for domesticated or captive-born species because
8  domestication can alter seasonal patterns in hormone physiology. Estimates from these species were not
9  included in analyses that examined the effect of life history stage.
10  [5]Experimental manipulation refers to studies in which some or all individuals underwent a stressful
11  manipulation intended to produce a response (not including routine capture and handling stress) at some
12  point during the course of the study.
13  [6]Adjusted refers to whether or not estimates reflect GC repeatability after statistically controlling for
14  factors expected to explain some of the variation in GC titers (e.g., year, sex, weather).
15
16

# Figure 1

PRISMA flow diagram

Figure 1. Preferred Reporting Items for Systematic Reviews and Meta-analysis (PRISMA) flowchart illustrating the process of study identification, screening, and inclusion in the meta-analysis.

Figure 1 Footnotes:

[1]We used the search terms: repeatab*, consisten*, glucocorticoid, cortisol, corticoster*, repeated measure, individual variation

[2]We included three studies that did not meet inclusion criteria (i.e., collected repeated within individuals, but did not estimate repeatability) because we were able obtain the original data from the study authors and calculate repeatability ourselves.

[3]We used the following inclusion criteria: the study had to assess repeated measurements of glucocorticoids within the same individual, and estimate a repeatability coefficient (e.g., Spearman rank, Pearson, or intraclass correlation coefficient).

```
┌─────────────────────────────┐      ┌─────────────────────────────┐
│ Records identified in Google │      │ Studies identified through   │
│  Scholar search¹ (Nov 2017): │      │  personal communication with │
│            n=716             │      │       authors²: n=3          │
└─────────────────────────────┘      └─────────────────────────────┘
```

Records identified in Google Scholar search[1] (Nov 2017): n=716

Studies identified through personal communication with authors[2]: n=3

Records remaining after removing inaccessible/duplicate records: n=664

Record titles and abstracts screened: n=664

Records excluded for being irrelevant/duplicate: n=528

Full-text studies assessed for eligibility: n=136

Studies excluded for not meeting inclusion criteria[3]: n=89

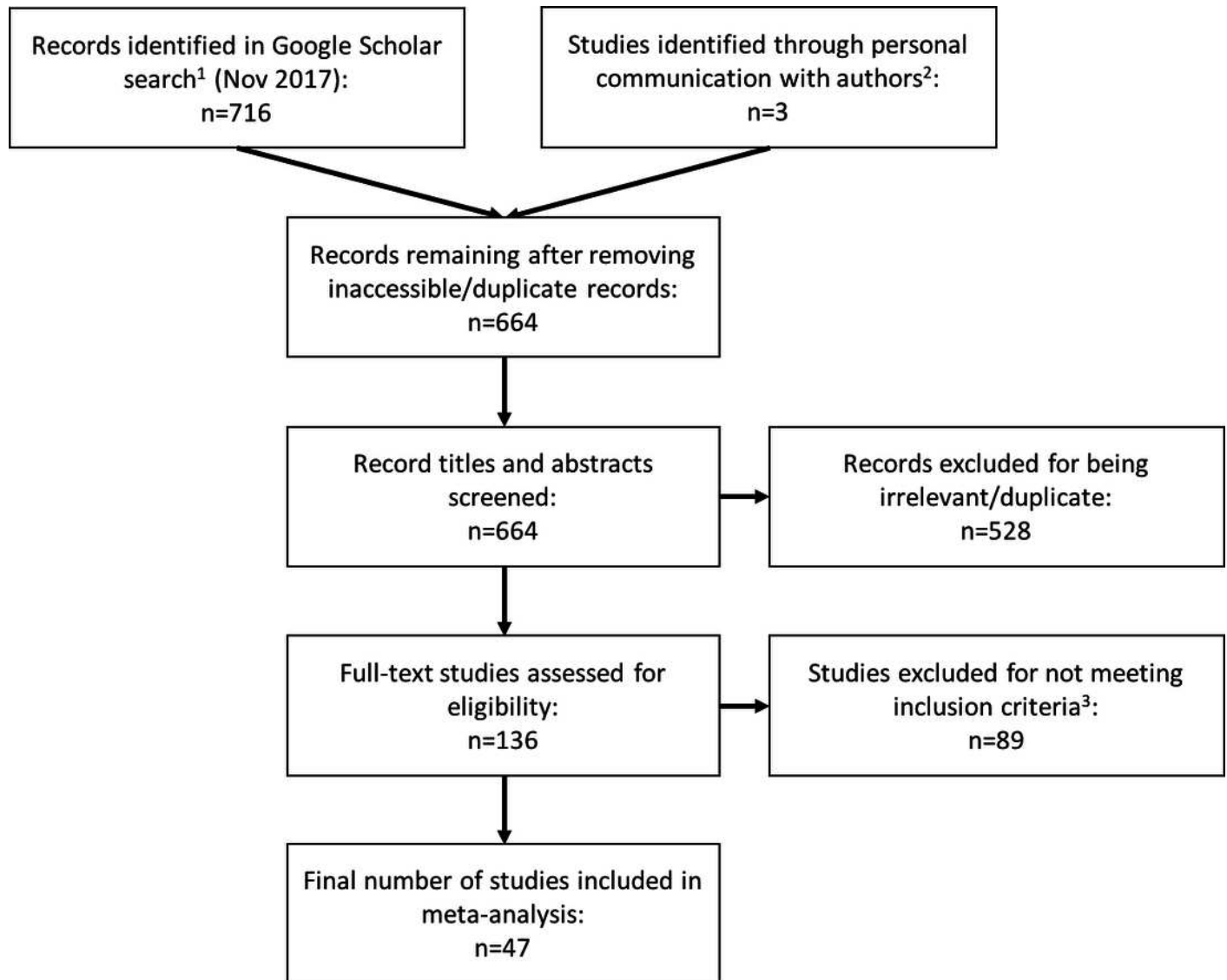Final number of studies included in meta-analysis: n=47

# Figure 2

Figure 2

Figure 2. Frequency distributions of all estimates of repeatabilities of A) *initial*, B) *response*, and C) *integrated* glucocorticoid (GC) measures included in the meta-analyses. The mean repeatability across all estimates of each category of GC is represented by a solid line, and the 95% CI (calculated from 1000 bootstrap samples of the data with replacement) is represented by a dashed line. In this study, we defined *initial* measures as those representing GCs in circulation within a time period expected not to reflect the acute stress of capture, *response* for elevated GC titers following an acute capture stress, and *integrated* for GC titers that represent hormone secretion over a relatively long period of time (e.g., GC concentrations in feces, feathers, and saliva).
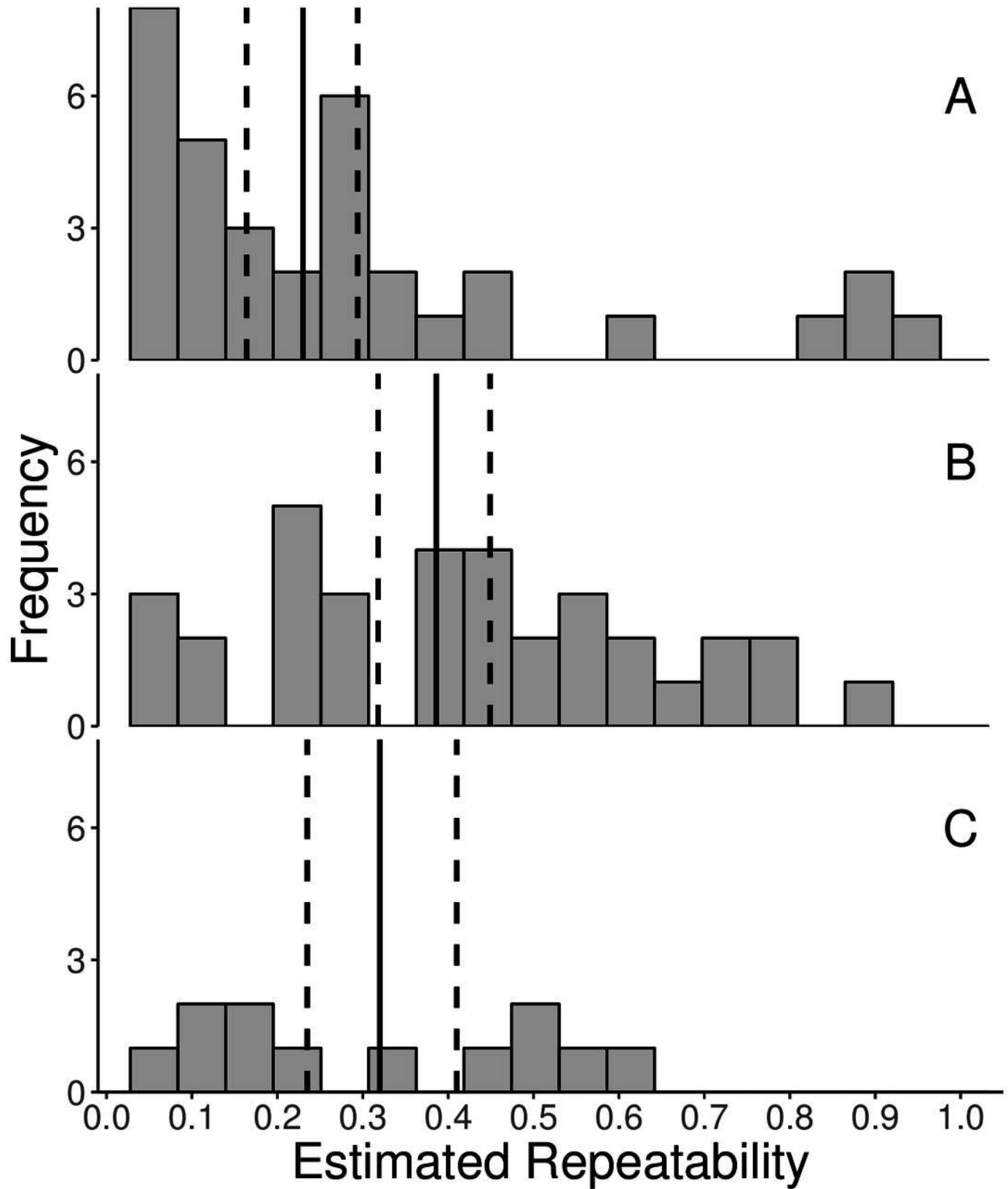
# Figure 3

Figure 3

Figure 3. Boxplots showing variation in the average repeatability of all glucocorticoid (GC) measures across taxonomic classes (data are jittered along x-axis for ease of interpretation). The plot's whiskers represent the 1.5 interquartile range, while the boxes represent the first and third quartiles, and the midline represents the median. Repeatability estimates for *initial* (open circles) and *response* (open triangles), but not *integrated* (closed squares), GC measures varied across taxonomic class (Type III ANOVA; initial*: n*=38, *F(3,38)*=9.359, *p*<0.0001; response: *n*=27, *F(4,23)*=4.984, *p*=0.005). In this study, we defined *initial* measures as those representing GCs in circulation within a time period expected not to reflect the acute stress of capture, *response* for elevated GC titers following an acute capture stress, and *integrated* for GC titers that represent hormone secretion over a relatively long period of time (e.g., GC concentrations in feces, feathers, and saliva).