

11-2016

On the comparison of the strength of morphological integration across morphometric datasets

Dean C. Adams
Iowa State University, dcadams@iastate.edu

Michael L. Collyer
Western Kentucky University

Follow this and additional works at: http://lib.dr.iastate.edu/eeob_ag_pubs



Part of the [Evolution Commons](#), and the [Statistical Methodology Commons](#)

The complete bibliographic information for this item can be found at http://lib.dr.iastate.edu/eeob_ag_pubs/201. For information on how to cite this item, please visit <http://lib.dr.iastate.edu/howtocite.html>.

1 Brief Communication

2

3 ON THE COMPARISON OF THE STRENGTH OF MORPHOLOGICAL INTEGRATION ACROSS

4

MORPHOMETRIC DATASETS

5

6

Dean C. Adams^{1,3} and Michael L. Collyer²

7

8 ¹*Department of Ecology, Evolution, and Organismal Biology, Iowa State University, Ames IA, USA*

9

Department of Statistics, Iowa State University, Ames IA, USA

10

²*Department of Biology, Western Kentucky University, Bowling Green, KY, USA*

11

³*Corresponding author email: dcadams@iastate.edu*

12

13

14

15 Short title: Comparing integration strength across datasets

This is the peer reviewed version of the following article: Adams, D. C. and Collyer, M. L. (2016), On the comparison of the strength of morphological integration across morphometric datasets. *Evolution*, 70: 2623–2631, which has been published in final form at doi:10.1111/evo.13045. This article may be used for non-commercial purposes in accordance With Wiley Terms and Conditions for self-archiving.

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

Abstract

Evolutionary morphologists frequently wish to understand the extent to which organisms are integrated, and whether the strength of morphological integration among subsets of phenotypic variables differ among taxa or other groups. However, comparisons of the strength of integration across datasets are difficult, in part because the summary measures that characterize these patterns (RV and r_{PLS}) are dependent both on sample size and on the number of variables. As a solution to this issue we propose a standardized test statistic (a z -score) for measuring the degree of morphological integration between sets of variables. The approach is based on a partial least squares analysis of trait covariation, and its permutation-based sampling distribution. Under the null hypothesis of a random association of variables, the method displays a constant expected value and confidence intervals for datasets of differing sample sizes and variable number, thereby providing a consistent measure of integration suitable for comparisons across datasets. A two-sample test is also proposed to statistically determine whether levels of integration differ between datasets, and an empirical example examining cranial shape integration in Mediterranean wall lizards illustrates its use. Some extensions of the procedure are also discussed.

Accepted for Evolution

32

Introduction

33 Over the past several decades, evaluating the degree to which morphological traits covary
34 (*morphological integration*: sensu Olson and Miller 1958), has become a prominent subject in
35 evolutionary biology (Cheverud 1982; Cheverud 1996; Bookstein et al. 2003; Pigliucci 2003;
36 Klingenberg 2008; Goswami and Polly 2010). Myriad studies have characterized patterns of
37 morphological integration in a variety of organisms, in an effort to decipher the genetic, developmental,
38 and functional mechanisms that generate such patterns (e.g., Hallgrímsson et al. 2002; Mitteroecker et al.
39 2004; Monteiro et al. 2005; Young and Badyaev 2006; Gómez-Robles and Polly 2012). These empirical
40 studies have been facilitated in part by the development of quantitative approaches for characterizing
41 patterns of morphological integration in high-dimensional data (e.g., Magwene 2001; Bookstein et al.
42 2003; Mitteroecker and Bookstein 2007; Márquez 2008; Klingenberg 2009; Adams and Felice 2014;
43 Bookstein 2015). In particular, methods that evaluate covariance patterns across *a priori* subsets of
44 variables have received considerable attention.

45 In the field of geometric morphometrics, a number of approaches are utilized for characterizing
46 the integration among subsets of variables. One approach is based on Escoffier's (1973) RV coefficient
47 (Klingenberg 2009), which is a ratio describing the degree of covariation between sets of variables
48 relative to the variation and covariation within sets of variables. The RV coefficient ranges between zero
49 and one, and larger values describe greater covariation between sets of variables relative to within them,
50 which may provide evidence that there is higher integration among subsets than expected by chance. An
51 alternative measure is based on partial least squares (PLS), where a singular value decomposition of the
52 covariance matrix between two sets of variables (\mathbf{S}_{12}) is used to describe the maximal covariation between
53 them (Bookstein et al. 2003; Mitteroecker and Bookstein 2007). The dominant singular value of \mathbf{S}_{12}
54 explains the maximal covariation between the two sets of variables, whose pattern of covariation is
55 described by the first set of linear combinations (singular vectors) in each of the two datasets (Bookstein
56 et al. 2003). Scores projected on these axes are routinely used to estimate the maximal correlation among
57 sets of variables (r_{PLS} : Rohlf and Corti 2000), with higher correlations indicating a greater level of

58 covariation. For both the RV and PLS approaches, statistical evaluation of the observed pattern is
59 accomplished using permutation, where the rows (individuals) are shuffled in one subset of variables
60 while leaving the rows in the other subset constant, thereby disassociating the covariation between subsets
61 and generating a distribution of possible outcomes under the null hypothesis of no association between
62 variable subsets. The observed statistic is then compared to a distribution of random statistics obtained
63 from this procedure to evaluate its significance (see Rohlf and Corti 2000; Bookstein et al. 2003;
64 Klingenberg 2009).

65 Recently, there has been increased interest in understanding the extent to which patterns of
66 morphological integration are consistent across levels of biological organization, and whether levels of
67 integration change over evolutionary time ([Armbruster et al. 2014](#); [Goswami et al. 2014](#); Klingenberg
68 2014). To this end researchers have characterized levels of integration across traits and species using one
69 or more of the methods mentioned above for subsequent qualitative or quantitative comparison (for recent
70 examples see: Drake and Klingenberg 2010; [Goswami et al. 2014](#); [Lazic et al. 2015](#); Martin-Serra et al.
71 2015; [Neaux et al. 2015](#)). However, for such comparisons to be meaningful requires that the evaluated
72 test measures are unaffected by other attributes of the data. Unfortunately this is not the case. For
73 instance, the RV coefficient has been shown to be sensitive to both the sample size (n) and the number of
74 variables examined (p), rendering comparisons of RV measures across datasets uninformative (Adams
75 2016; also: Smilde et al. 2009; [Fruciano et al. 2013](#); for an extended critique of the RV coefficient see:
76 Bookstein 2016). Additionally, as shown in part by Mitteroecker and Bookstein (2007), and
77 comprehensively below, the PLS correlation coefficient (r_{PLS}) suffers from the same inherent issues. The
78 objective of the current manuscript is to provide a standardized test statistic (a z -score) for measuring the
79 degree of morphological integration between sets of variables. Our procedure is developed for and is used
80 on a PLS correlation of among-partition trait covariation. However, it is sufficiently flexible that it may
81 be applied to any meaningful measure that captures the degree of integration in a dataset, and is thus a
82 useful approach for comparing the degree of integration as new analytical approaches are developed (see
83 Discussion).

84

85

Sample Size and Variable Dependency of r_{PLS}

86 To understand the properties of r_{PLS} we conducted simulations similar to those of Adams (2016).

87 Specifically, simulated datasets were obtained by generating random variables drawn from a normal

88 distribution $\sim N(0,1)$, and variables were randomly assigned to one of two subsets with the constraint that

89 the number of variables was the same in each subset. Thus, each simulated dataset represented what was

90 expected under the null hypothesis of a random association of variables. Using this procedure, we

91 generated 100 datasets for differing levels of sample size (n), where the total number of variables was the92 same ($p = 30$). Next we performed the reciprocal simulation where all datasets contained the same93 number of specimens ($n = 100$), but where the total number of variables differed. From each simulated94 dataset r_{PLS} was estimated, and at each level of n and p , the mean and 95% confidence intervals across the

95 100 datasets were calculated. All simulations were performed in R 3.2.0 (R Core Team 2015).

96 As is clear from Figure 1, values of r_{PLS} vary between zero and one, with larger values attained

97 under smaller sample sizes, as well as with a larger number of variables (Fig. 1A, B: see also

98 Mitteroecker and Bookstein 2007). Thus, like the RV coefficient, estimates of morphological integration

99 using partial least squares are also sensitive to n and p , rendering comparisons of these values across

100 datasets challenging (see also Mitteroecker and Bookstein 2007). Thus, for this purpose an alternative

101 estimate of the degree of integration across sets of variables is required.

102

The Z-Score for Comparing the Strength of Integration

104 Although studies of integration rarely state a null hypothesis in terms of the parameters tested, the

105 permutation-based procedure described above evaluates the observed measure against a distribution of

106 values obtained under a null hypothesis of no association between subsets of variables. Thus, some

107 generalization of the Pearson product-moment null hypothesis, $\rho = 0$, could be implied. However,108 whereas Pearson's r , as an estimate of ρ , has an expected value of 0 under the null hypothesis for109 univariate tests, r_{PLS} has a lower limit of 0 and an expected value that varies with n and p (Fig. 1 A, B).

110 For single-sample hypotheses, stating the null hypothesis as “no association” between matrices is
 111 sufficient; the expected value is simply the mean of the sampling distribution of r_{PLS} from the permutation
 112 procedure (as described above) and the percentile of the observed r_{PLS} value is the estimate of the P -value.
 113 To either qualitatively compare or actually test the dissimilarity of two measures of integration from two
 114 samples requires calculation of effect sizes in relation to expected values under the null hypothesis of no
 115 integration (e.g., Collyer et al. 2015), especially if the two samples have different expected values. This
 116 can be accomplished by calculating the standard deviates (effect sizes) of r_{PLS} for the different samples,

$$\hat{z} = \frac{r - \hat{\mu}_r}{\hat{\sigma}_r}, \quad (1)$$

117 where $\hat{\mu}_r$ is the estimated expected value of r_{PLS} under the null hypothesis, found as the mean of the
 118 sampling distribution, and $\hat{\sigma}_r$ is the standard deviation of the sampling distribution (i.e., standard error of
 119 the mean). Calculating effect sizes this way assumes that the sampling distribution is normally
 120 distributed, a property we demonstrate via simulation (below).

121 At first glance, the numerator of the standard deviate calculation might also seem sufficient for
 122 calculating a statistic that allows integration to be compared between sets of variables. Indeed, the
 123 numerator of the standard deviate calculation detrends r_{PLS} values that might have different expectations
 124 under the null hypothesis (Fig. 1 C, D). However, there are two important concerns with using detrended
 125 r_{PLS} as a comparable statistic. First, although the statistic is no longer n - or p -dependent, the standard error
 126 (and thus, confidence interval) of the statistic varies across n and p for the same number of random
 127 permutations, decreasing with either increased n or p (Fig. 1 C, D). Second, the value itself has little
 128 meaning as a correlation coefficient. For example, negative values do not indicate negative covariation,
 129 but rather less covariation than expected under the null hypothesis. Furthermore, because low n or high p
 130 can produce large expected r_{PLS} values (Fig. 1 A, B), a veritable strong correlation between two sets of
 131 variables with large sample size will have a small detrended value, as a maximum r_{PLS} value of 1.0
 132 precludes large detrended values for large samples. Therefore, standardization – dividing by the standard
 133 deviation of the sampling (null) distribution – is needed to generate a test statistic (z) that has a constant

134 expected value and same variance, standard deviation, or confidence interval across the entire spectrum of
 135 sample sizes and variable number (Fig 1. E, F).

136 We wish to reiterate that other than concern for n - or p -dependency for comparisons of
 137 integration across multiple datasets, r_{PLS} is sufficient as a single-sample test statistic. The effect size
 138 calculation in Equation 1 is merely descriptive, but allows qualitative comparison between two measures
 139 of integration that have different expected values. (Inferring the probability of the effect size from a
 140 standard normal distribution is not necessary, as it is already accomplished from the resampling
 141 experiment of the PLS analysis, as shown below.) For a statistical treatment of the comparisons of two or
 142 more datasets, however, Equation 1 can be modified to calculate the effect size of the difference between
 143 two integration effect sizes as

144

$$\hat{z}_{12} = \frac{|(r_1 - \hat{\mu}_1) - (r_2 - \hat{\mu}_2)|}{\sqrt{\hat{\sigma}_1^2 + \hat{\sigma}_2^2}}. \quad (2)$$

145

146 The probability of \hat{z}_{12} under the null hypothesis of equal integration (P -value) can be estimated from a
 147 standard normal distribution, provided the within-sample random r_{PLS} values are approximately normally
 148 distributed. The numerator of Equation 1 can be rewritten as $(r_1 - r_2) - (\hat{\mu}_1 - \hat{\mu}_2)$, which characterizes the
 149 effect size as a difference in levels of integration compared to their expected difference. We assert that a
 150 null hypothesis test for comparing levels of integration is naturally a two-tailed test (hence the absolute
 151 value in the numerator of Equation 2), as the direction of the difference between detrended r_{PLS} scores has
 152 no appreciable meaning, especially in the absence of standardization. For example, a smaller value of
 153 $(r_2 - \hat{\mu}_2)$ than $(r_1 - \hat{\mu}_1)$ might produce a positive difference, $(r_1 - \hat{\mu}_1) - (r_2 - \hat{\mu}_2)$, even if \hat{z}_2 is larger
 154 than \hat{z}_1 , because of a smaller standard error in the second sample (perhaps because of a larger expected
 155 value from the sampling distribution). Therefore, only the magnitude of the effect size of the integration
 156 difference is valuable as a test statistic, rendering the hypothesis inherently two-tailed. Likewise, a

157 confidence interval for the effect size can be estimated as

$$(1-\alpha)100\%CI = |(r_1 - \hat{\mu}_1) - (r_2 - \hat{\mu}_2)| \pm z_{\alpha/2} \sqrt{\hat{\sigma}_{r_1}^2 + \hat{\sigma}_{r_2}^2}, \quad (3)$$

158

159 where $z_{\alpha/2}$ is the quantile from a standard normal distribution corresponding to the two-tailed probability

160 for the level of significance, α , and $\sqrt{\hat{\sigma}_{r_1}^2 + \hat{\sigma}_{r_2}^2}$ is the pooled standard error. The null hypothesis is

161 rejected with $(1-\alpha)100\%$ confidence, if the confidence interval does not contain 0, the expected

162 difference in effect sizes under the null hypothesis. Thus, concerning hypothesis tests, confidence

163 intervals and the two-sample z -score are test statistics for the difference in r_{PLS} between different data sets,

164 relative to their expected difference based on n - or p -dependency. The effect sizes (equation 1) allow one

165 to infer which sample has greater morphological integration, when a significant difference is observed.

166

167

Statistical Properties

168

169 To determine if this test of difference in the strength of integration among samples has

170 appropriate statistical properties (type I and type II error rates, statistical power), we performed a

171 simulation experiment. This experiment consisted of 100 runs of generating two “populations” of data,

172 each with $N = 10,000$ individuals, where each individual comprised a vector of random values sampled

173 from a normal distribution, $\sim N(0,1)$, for a matrix, \mathbf{X} , comprising 10 variables and matrix, \mathbf{Y} , also

174 comprising 10 variables. In one population, \mathbf{X} and \mathbf{Y} were simulated with no linear association (the range

175 of “true” PLS correlations was 0.0459-0.0709 over 100 runs). In the other population, \mathbf{X} and \mathbf{Y} were

176 simulated with a fairly strong linear association (the range of true PLS correlations was 0.8239- 0.8340

177 over 100 runs). Within each simulation run, two sampling frames of $n = 30$ specimen numbers were

178 randomly generated, for 200 permutations. The first sampling frame was used to sample individuals from

179 each of the “no effect” population (no linear association) and “effect” population (linear association). In

each sample, r_{PLS} values were calculated and sampling distributions for r_{PLS} were generated, based on

1,000 random resampling permutations for the PLS analysis); \hat{Z}_{12} was calculated between samples, as in Equation 2. For this comparison, we expected the null hypothesis to be rejected. Therefore, the proportion of times it was rejected is a measure of statistical power for the difference in effects. (Likewise, the proportion of times it was not rejected is a measure of the type II error rate.)

The second sampling frame was used to draw a second sample from each population, and r_{PLS} values and their sampling distributions were calculated as above, but \hat{Z}_{12} values were calculated between samples from the same population. For these comparisons, we expected the null hypothesis to be supported, but in one case no integration (no effect) and in the other, substantial integration (effect) should be found from samples of each population. Irrespective of the level of integration, this procedure simulated multiple runs of a null model (no expected difference in integration between populations), allowing us to ascertain whether the type I error rates – the proportion of times the null hypothesis was rejected – was consistent in spite of different levels of morphological integration. In all cases, the significance level for the null hypothesis was assigned as $\alpha = 0.05$. Type I error rates and statistical power were calculated within runs as the proportion of 200 permutations that the null hypothesis was rejected. Means and 95% confidence intervals were calculated across the 100 runs to evaluate the tendencies of type I error rates and statistical power.

Additionally, because integration tests are often performed on geometric morphometric data, we also considered the contribution of generalized Procrustes analysis (GPA) to generate inherently correlated shape variables, and to what extent this would affect results of our proposed procedure. Briefly, GPA scales anatomical landmark configurations to unit size, centers them, and rotates them via least squares superimposition to render configurations invariant in size, position, or orientation. Procrustes residuals – the aligned landmarks following GPA – are inherently correlated shape variables produced by GPA. Therefore, even generating isotropic error on points of landmark configurations will produce correlated variables within modules prior to measuring the PLS correlation between modules. We considered the GPA-artifact by assigning the randomly generated data from the simulation experiment above as residuals on mean configurations. We generated 5-point configurations of two-dimensional

206 landmarks (a vector of 10 values) in each simulation run. These configurations served as mean vectors, to
207 which the row vectors of the previously generated random \mathbf{X} and \mathbf{Y} matrices were added as residuals to
208 generate individual configurations within populations. For the “no effect” cases described above, we
209 forced the 5-point configurations to be approximately uncorrelated ($r_{PLS} < 0.1$) between the mean
210 configurations for \mathbf{X} and \mathbf{Y} . For the “effect” cases, mean configurations were the same. The mean
211 configurations were also generated to have dispersion of several orders of magnitude greater than the
212 random within-point dispersion, to ensure resulting individual landmark configurations were reasonable.
213 In each permutation within each simulation run, GPA was performed and Procrustes residuals were used
214 for r_{PLS} calculations. Thus, our results allowed us to evaluate both GM (Procrustes residuals) and non-GM
215 (multivariate data) applications for consistently generated residuals, and whether GPA altered
216 interpretations.

217 We repeated the entire process for 5 simulation runs with 20 sampling events and PLS analyses
218 with 1,000 random permutations to consider whether the sampling distributions of r_{PLS} values were
219 normally distributed. In each sampling event, we calculated a Shapiro-Wilk W statistic (i.e., 100 values
220 total for each comparison), and the ranges of these values were qualitatively evaluated to determine if
221 sampling distributions were approximately normally distributed (i.e., a Shapiro-Wilk W statistic that tends
222 toward 1.0). All simulations were performed in R 3.2.0 (R Core Team 2015).

223 *Results.* The simulation experiment demonstrated that type I error rates were appropriate for both
224 random multivariate data and Procrustes residuals (Fig. 2). For “no effect” comparisons, the mean type I
225 error rates (0.0453 and 0.0494 for multivariate data and Procrustes residuals, respectively) were
226 approximately the same as the nominal significance level (0.05). Interestingly, the type I error rates were
227 lower for comparisons between equal but large effects (i.e., when comparing datasets where both
228 exhibited marked morphological integration of similar strength). A type I error was simulated in 0 and 1
229 permutation of the 100×200 total permutations, for multivariate data and Procrustes residuals,
230 respectively. The nearly non-existent type I errors imply that when both datasets contain integration but
231 their levels are similar, the method displays fewer false positives as compared to when comparing datasets

232 lacking integration (Fig.2). When comparing samples from the “no effect” and “effect” populations, the
233 mean statistical power was 0.8133 and 0.883 for multivariate data and Procrustes residuals, respectively.
234 Both of these values were, approximately equal to or greater than the generally desired value of 0.8 (Fig.
235 2), which represents a 4:1 trade-off between type II error risk and type I error risk (Cohen 1988).
236 However, the inherently generated within-module correlations from GPA appear to slightly increase
237 statistical power, perhaps owing to the additional correlation between mean configurations that was
238 simulated. These results suggest that the two-sample test of integration disparity presented here behaves
239 as expected, statistically, and GPA does not negatively impact results.

240 In terms of the appropriateness of the two-sample test of integration disparity for PLS analyses,
241 we found no evidence to suggest that the sampling distributions of r_{PLS} statistics were non-normally
242 distributed. For multivariate data, the ranges in W statistics for the “no effect” case (0.9822-0.9996) and
243 “effect” case (0.9885-0.9997) were quite similar and sufficiently close to 1.0 in each case. For the
244 Procrustes residuals, the ranges in W statistics for the “no effect” case (0.9910-0.9996) and “effect” case
245 (0.9895-0.9996) were also quite similar and sufficiently close to 1.0 in each case. Thus, whether GPA was
246 performed had no consequence, and any PLS analysis produced sufficiently normal sampling
247 distributions.

248 It should be noted that the sampling distribution of r_{PLS} values is arbitrarily based on the *a priori*
249 chosen number of PLS resampling permutations. One may wish to either choose a sufficiently large
250 number of permutations or confirm that the standard deviation of the sampling distribution remains
251 consistent under a range of permutations. For example, we found the standard deviation in our samples of
252 30 individuals in the simulation experiment was rather consistent between 200-10,000 random
253 permutations, suggesting the 1,000 permutations we used was adequate for measuring the difference in
254 strength of integration among datasets. We also found that when using a small number of PLS resampling
255 permutations (e.g., 100-200), a large effect for the observed r_{PLS} , could skew the sampling distribution
256 (which should be normally distributed). This problem is alleviated by simply removing the observed r_{PLS}
257 from the sampling distribution, as a small-sample bias adjustment. For sufficiently large numbers of PLS

258 permutations (e.g., 1,000), this step was not needed, but also did not affect results. We, therefore,
259 recommend removing the observed r_{PLS} as a procedural step to assure a more appropriate standard
260 deviation of sampling distributions, especially for large effect sizes.

261

262 **A Biological Example**

263 To illustrate the method described above we conducted a comparison of levels of integration in
264 cranial shape between rural and urban populations of the Mediterranean wall lizard *Podarcis muralis*. The
265 data were part of a series of studies that evaluated the effects of urbanization on various aspects of
266 phenotypic variation, including patterns of allometry, developmental stability, and integration in juvenile
267 and adult lizards (see Lazic et al. 2015, 2016: available on Dryad and from the original authors). For this
268 example, landmark-based geometric morphometric methods (Bookstein 1991; Mitteroecker and Gunz
269 2009; Adams et al. 2013) were used to characterize head shape, based on the positions of 28 homologous
270 locations (Fig. 3A) collected from the dorsal view of 482 juvenile and 359 adult lizards from several
271 localities. Of these, approximately half of the specimens were collected from rural locations (218
272 juveniles and 191 adults), and the remainder from urban sites (264 juveniles and 168 adults). For each
273 specimen, a mirror image of their landmark locations was obtained by reflecting the coordinates about the
274 mid-line, and a generalized Procrustes analysis was then performed to remove the effects of non-shape
275 variation from the dataset (Fig. 3B). The symmetric component of shape was subsequently obtained by
276 averaging landmark locations for each specimen and its mirror image. From the symmetric component of
277 shape variation, we evaluated the degree to which the anterior and posterior regions of the head were
278 integrated with one another. Landmarks were classified as anterior or posterior (Fig. 3A) based on the
279 timing of ossification events during development (Lazic et al. 2015). For both juveniles and adults from
280 rural and urban populations, the degree of integration between modules was estimated using partial least
281 squares. Here the maximal covariation between modules was characterized using the PLS correlation
282 (r_{PLS}), which was statistically evaluated using 1,000 random permutations. Additionally, z -scores were
283 obtained for all four groups, and were statistically compared to one another using the procedure described

284 above. All analyses were performed in R 3.2.0 (R Core Team 2015) using the package *geomorph* (Adams
285 and Otárola-Castillo 2013; Adams et al. 2016), including the function, `compare.pls`, which performs
286 the method introduced here.

287 *Results.* For both juvenile and adults from rural and urban lizard populations, the degree of
288 morphological integration between anterior and posterior regions of the head was large and highly
289 significant (rural_{juv}: $r_{PLS} = 0.770$, $P_{rand} = 0.001$; rural_{adult}: $r_{PLS} = 0.826$, $P_{rand} = 0.001$; urban_{juv}: $r_{PLS} = 0.761$,
290 $P_{rand} = 0.001$; urban_{adult}: $r_{PLS} = 0.858$, $P_{rand} = 0.001$). Generally, the observed integration of the frontal and
291 distal regions of the head was represented by a relative shortening of the snout accompanied by an
292 enlargement of the posterior area of the head (Fig. 3C); a pattern broadly observed in all groups. When
293 converted to effect sizes, all z -score values were very large (Fig. 3D), implying that the degree of
294 integration in each group greatly exceeded that which was expected by chance. Interestingly, when
295 compared using equation 2 above, we found no evidence of differences in levels of integration between
296 rural and urban juveniles or rural and urban adults, but significant differences existed in levels of
297 integration between juveniles and adults within each population; with adults displaying significantly
298 greater integration relative to juveniles (Table 1). Thus, while there was no evidence that environmental
299 disturbances have affected the strength of morphological integration in urban populations, the results
300 demonstrate that morphology is more integrated through ontogenetic time. Further, when combined with
301 the prior observation that morphological variation among adult specimens was reduced when compared to
302 that of juveniles (see Lazic et al. 2016), the pattern identified here is consistent with what is expected
303 under the hypothesis of developmental canalization (Hallgrímsson et al. 2002).

304

305

Discussion

306 An important question in evolutionary biology is whether different taxa or traits display similar
307 levels of morphological integration. Unfortunately, direct quantitative comparisons of the level of
308 integration across datasets have been hampered by the fact that the measures that characterize these
309 patterns (RV and r_{PLS}) are dependent on both sample size and the number of variables. Here we described

310 an unbiased effect size for quantifying the strength of morphological integration between sets of
311 variables, utilizing partial least squares analysis and its permutation sampling distribution. We
312 demonstrated that under the null hypothesis of a random association of variables, the approach displays a
313 constant expected value, and the same variance, standard deviation, and confidence intervals across the
314 entire spectrum of sample sizes and variable number. We further proposed a two-sample test to
315 statistically evaluate the difference in effect sizes across two datasets. The approach displays appropriate
316 type I error and statistical power, thereby providing a rigorous means of evaluating whether the degree of
317 morphological integration differs between them. Thus the approach provides evolutionary morphologists
318 with a consistent means of deciphering whether levels of integration are similar to one another in two or
319 more datasets.

320 Extensions to the approach developed here can be envisioned that address a wider array of
321 empirical challenges than the ones presented. For example, if the integration among three sets of variables
322 is of biological interest, a three-block partial least squares approach may be utilized to characterize the
323 degree of covariation among sets of variables (see Bookstein et al. 2003). Alternatively, one may evaluate
324 the mean of the pairwise r_{PLS} values as a general test measure, as has been proposed for evaluating
325 hypotheses of modularity (Klingenberg 2009; Adams 2016). In either case it is important to recognize that
326 the method developed here simply provides a quantitative estimate on the magnitude, or strength of
327 morphological integration among sets of variables. The method provides no description of the type of
328 integration, or *how* traits covary with one another and in what manner. For this, understanding patterns of
329 morphological integration and the set of coordinated shape changes it embodies must still be
330 accomplished via a thorough examination of the singular vectors from the PLS and a visual inspection of
331 the shape changes associated with the singular vectors (see Bookstein 2016 for discussion). Thus, a
332 proper biological understanding of patterns of morphological integration is accomplished via the
333 combination of a quantitative characterization and statistical assessment of the magnitude of integration,
334 along with its anatomical interpretation.

335 Finally, while the biological concept of morphological integration has been embraced in

336 evolutionary biology for decades (Olson and Miller 1958), formal statistical tests of these patterns are still
337 rather novel. As the theory of morphological integration develops, so too will the analytical methods for
338 measuring integration, and their associated hypothesis tests (see e.g., Bookstein 2015). An important
339 advance made here is that for any two measures of morphological integration – irrespective of the number
340 of variables, number of specimens studied, or expected values in a null distribution – effect sizes
341 calculated as standard deviates in sampling distributions are values that can be compared in a general two-
342 sample hypothesis test. We chose to use the correlation coefficient from two-block PLS as the basis for
343 this test, but one could have likewise used maximum singular values, or vector correlations (or angles)
344 between left and right singular vectors found through PLS, or other statistically-relevant summaries.
345 Since standard deviates can be calculated using any test statistic with a sampling distribution – essentially
346 any statistic if resampling experiments are used – the hypothesis testing framework proposed here is
347 merely a methodological extension of two sample Z-tests, but using resampling experiments to generate
348 sampling distributions rather than requiring *a priori* knowledge of population standard deviations. Thus,
349 as new analytical approaches are developed for evaluating patterns of morphological integration, the test
350 procedure described here may be utilized for comparing the degree of integration observed in different
351 datasets based on those measures. Ultimately however, the choice of which test statistic to utilize in this
352 procedure must be driven by biology. In the case of morphological integration, biological interpretations
353 of such patterns depend on a deep understanding of how covariation patterns are embodied in terms of
354 their singular vectors (see Bookstein 2016). Nevertheless, as the field of evolutionary morphological
355 integration evolves, and better conceptual measures of morphological integration are developed, a
356 hypothesis-testing framework developed here is already in place, and ready for such advances.

357

358 **Acknowledgments**

359 We thank A. Kaliontzopoulou for her comments on the manuscript. A Kaliontzopoulou and M.
360 Lazić kindly provided the data for the empirical example and the image for Fig. 2a. This work was
361 sponsored in part by National Science Foundation grants DEB-1556379 (to DCA) and DEB-1556540 (to

362 MLC).

Accepted for Evolution

References

- 363
- 364 Adams, D. C. 2016. Evaluating modularity in morphometric data: challenges with the RV coefficient and
365 a new test measure. *Methods Ecol. Evol.* 7:565-572.
- 366 Adams, D. C., M. Collyer, and E. Sherratt. 2016. geomorph 3.0.1: Software for geometric morphometric
367 analyses. R package version 3.0.1. <http://CRAN.R-project.org/package=geomorph>.
- 368 Adams, D. C. and R. N. Felice. 2014. [Assessing trait covariation and morphological integration on
369 phylogenies using evolutionary covariance matrices. PLoS ONE 9:e94335.](#)
- 370 Adams, D. C. and E. Otárola-Castillo. 2013. [geomorph: an R package for the collection and analysis of
371 geometric morphometric shape data. Methods Ecol. Evol.](#) 4:393-399.
- 372 Adams, D. C., F. J. Rohlf, and D. E. Slice. 2013. [A field comes of age: Geometric morphometrics in the
373 21st century. Hystrix 24:7-14.](#)
- 374 Armbruster, W. S., C. Pelabon, G. H. Bolstad, and T. F. Hansen. 2014. [Integrated phenotypes:
375 understanding trait covariation in plants and animals. Phil. Trans. Roy. Soc. London B.
376 369:20130245.](#)
- 377 Bookstein, F. L. 1991. [Morphometric tools for landmark data: geometry and biology. Cambridge
378 University Press, Cambridge.](#)
- 379 Bookstein, F. L. 2015. [Integration, disintegration, and self-similarity: characterizing the scales of shape
380 variation in landmark data. Evol. Biol.](#) 42:395-426.
- 381 Bookstein, F. L. 2016. [The inappropriate symmetries of multivariate statistical analysis in geometric
382 morphometrics. Evol. Biol.](#) 43:DOI 10.1007/s11692-11016-19382-11697.
- 383 Bookstein, F. L., P. Gunz, P. Mitteroecker, H. Prossinger, K. Schaefer, and H. Seidler. 2003. [Cranial
384 integration in Homo: singular warps analysis of the midsagittal plane in ontogeny and evolution. J.
385 Hum. Evol.](#) 44:167-187.
- 386 Cheverud, J. M. 1982. [Phenotypic, genetic, and environmental morphological integration in the cranium.
387 Evolution 36:499-516.](#)
- 388 Cheverud, J. M. 1996. [Developmental integration and the evolution of pleiotropy. Am. Zool.](#) 36:44-50.

- 389 [Cohen, J. 1988. Statistical power analysis for the behavioral sciences. Lawrence Erlbaum, Hillsdale, New](#)
390 [Jersey. .](#)
- 391 [Drake, A. G. and C. P. Klingenberg. 2010. Large scale diversification of skull shape in domestic dogs:](#)
392 [disparity and modularity. Am. Nat. 175:289-301.](#)
- 393 [Escoufier, Y. 1973. Le traitement des variables vectorielles. Biometrics 29:751–760.](#)
- 394 [Fruciano, C., P. Franchini, and A. Meyer. 2013. Resampling-Based Approaches to Study Variation in](#)
395 [Morphological Modularity. PLoS ONE 8:e69376.](#)
- 396 [Gómez-Robles, A. and P. D. Polly. 2012. Morphological integration in the hominin dentition:](#)
397 [evolutionary, developmental, and functional factors. Evolution 66:1024-1043.](#)
- 398 [Goswami, A. and P. D. Polly. 2010. Methods for studying morphological integration and modularity. Pp.](#)
399 [213-243 in J. Alroy, and G. Hunt, eds. Quantitative Methods in Paleobiology.](#)
- 400 [Goswami, A., J. B. Smaers, C. Soligo, and P. D. Polly. 2014. The macroevolutionary consequences of](#)
401 [phenotypic integration: from development to deep time. Phil. Trans. Roy. Soc. London B.](#)
402 [369:20130254.](#)
- 403 [Hallgrímsson, B., K. Willmore, and B. K. Hall. 2002. Canalization, developmental stability, and](#)
404 [morphological integration in primate limbs. Am. J. Phys. Anthropol. 119:131–158.](#)
- 405 [Klingenberg, C. P. 2008. Morphological integration and developmental modularity. Ann. Rev. Ecol. Evol.](#)
406 [Syst. 39:115-132.](#)
- 407 [Klingenberg, C. P. 2009. Morphometric integration and modularity in configurations of landmarks: tools](#)
408 [for evaluating a priori hypotheses. Evol. Develop. 11:405-421.](#)
- 409 [Klingenberg, C. P. 2014. Studying morphological integration and modularity at multiple levels: concepts](#)
410 [and analysis. Phil. Trans. Roy. Soc. London B. 369:20130249.](#)
- 411 [Lazic, M. M., M. A. Carretero, J. Crnobrnja-Isailovic, and A. Kaliontzopoulou. 2015. Effects of](#)
412 [environmental disturbance on phenotypic variation: an integrated assessment of canalization,](#)
413 [developmental stability, modularity, and allometry in lizard head shape. Am. Nat. 185:44-58.](#)

- 414 [Lazic, M. M., M. A. Carretero, J. Crnobrnja-Isailovic, and A. Kaliontzopoulou. 2016. Postnatal dynamics](#)
415 [of developmental stability and canalization of lizard head shape under different environmental](#)
416 [conditions. *Evol. Biol.* DOI 10.1007/s11692-016-9377-4.](#)
- 417 [Magwene, P. M. 2001. New tools for studying integration and modularity. *Evolution* 55:1734–1745.](#)
- 418 [Márquez, E. J. 2008. A statistical framework for testing modularity in multidimensional data. *Evolution*](#)
419 [62:2688–2708.](#)
- 420 [Martin-Serra, A., B. Figueirido, J. A. Perez-Claros, and P. Palmqvist. 2015. Patterns of morphological](#)
421 [integration in the appendicular skeleton of mammalian carnivores. *Evolution* 69:321-340.](#)
- 422 [Mitteroecker, P. and F. L. Bookstein. 2007. The conceptual and statistical relationship between](#)
423 [modularity and morphological integration. *Syst. Biol.* 56:818–836.](#)
- 424 [Mitteroecker, P. and P. Gunz. 2009. Advances in geometric morphometrics. *Evol. Biol.* 36:235-247.](#)
- 425 [Mitteroecker, P., P. Gunz, M. Bernhard, K. Schaefer, and F. L. Bookstein. 2004. Comparison of cranial](#)
426 [ontogenetic trajectories among great apes and humans. *J. Hum. Evol.* 46:679-698.](#)
- 427 [Monteiro, L. R., V. Bonato, and S. F. d. Reis. 2005. Evolutionary integration and morphological](#)
428 [diversification in complex morphological structures: Mandible shape divergence in spiny rats](#)
429 [\(Rodentia, Echimyidae\). *Evol. Develop.* 7:429-439.](#)
- 430 [Neaux, D., E. Gilissen, W. Coudyzer, and F. Guy. 2015. Integration Between the Face and the Mandible](#)
431 [of Pongo and the Evolution of the Craniofacial Morphology of Orangutans. *Am. J. Phys. Anthropol.*](#)
432 [158:475-486.](#)
- 433 [Olson, E. C. and R. L. Miller. 1958. *Morphological Integration*. University of Chicago Press, Chicago.](#)
- 434 [Pigliucci, M. 2003. Phenotypic integration: studying the ecology and evolution of complex phenotypes.](#)
435 [*Ecol. Lett.* 6:265-272.](#)
- 436 R Core Team. 2015. R: a language and environment for statistical computing. Version 3.2.0.
437 <http://cran.R-project.org>. R Foundation for Statistical Computing, Vienna.
- 438 [Rohlf, F. J. and M. Corti. 2000. The use of partial least-squares to study covariation in shape. *Syst. Biol.*](#)
439 [49:740-753.](#)

440 [Smilde, A. K., H. A. L. Kiers, S. Biklsma, C. M. Rubingh, and M. J. vanErk. 2009. Matrix correlations](#)
441 [for high-dimensional data: the modified RV-coefficient. *Bioinformatics* 25:401-405.](#)

442 [Young, R. L. and A. V. Badyaev. 2006. Evolutionary persistence of phenotypic integration: influence of](#)
443 [developmental and functional relationships on complex trait evolution. *Evolution* 60:1291–1299.](#)

444

Accepted for Evolution

445

446

447 Table 1. Results from empirical example comparing levels of morphological integration in juvenile and

448 adult lizards from urban and rural populations. A) Matrix of pairwise differences in PLS effect sizes, and

449 B) their associated significance levels. Biologically-relevant focal comparisons are underlined; significant

450 focal comparisons are shown in bold. Populations are designated as: UR: urban, RU: rural, Ad: adult, Juv:

451 juvenile.

452

Z	UR _{Ad}	RU _{Ad}	UR _{Juv}	RU _{Juv}		P	UR _{Ad}	RU _{Ad}	UR _{Juv}	RU _{Juv}
UR _{Ad}	0					UR _{Ad}	0			
RU _{Ad}	<u>0.105</u>	0				RU _{Ad}	<u>0.459</u>	0		
UR _{Juv}	<u>2.818</u>	2.777	0			UR _{Juv}	0.002	0.002	0	
RU _{Juv}	2.447	<u>2.399</u>	<u>0.296</u>	0		RU _{Juv}	0.007	0.008	<u>0.383</u>	0

453

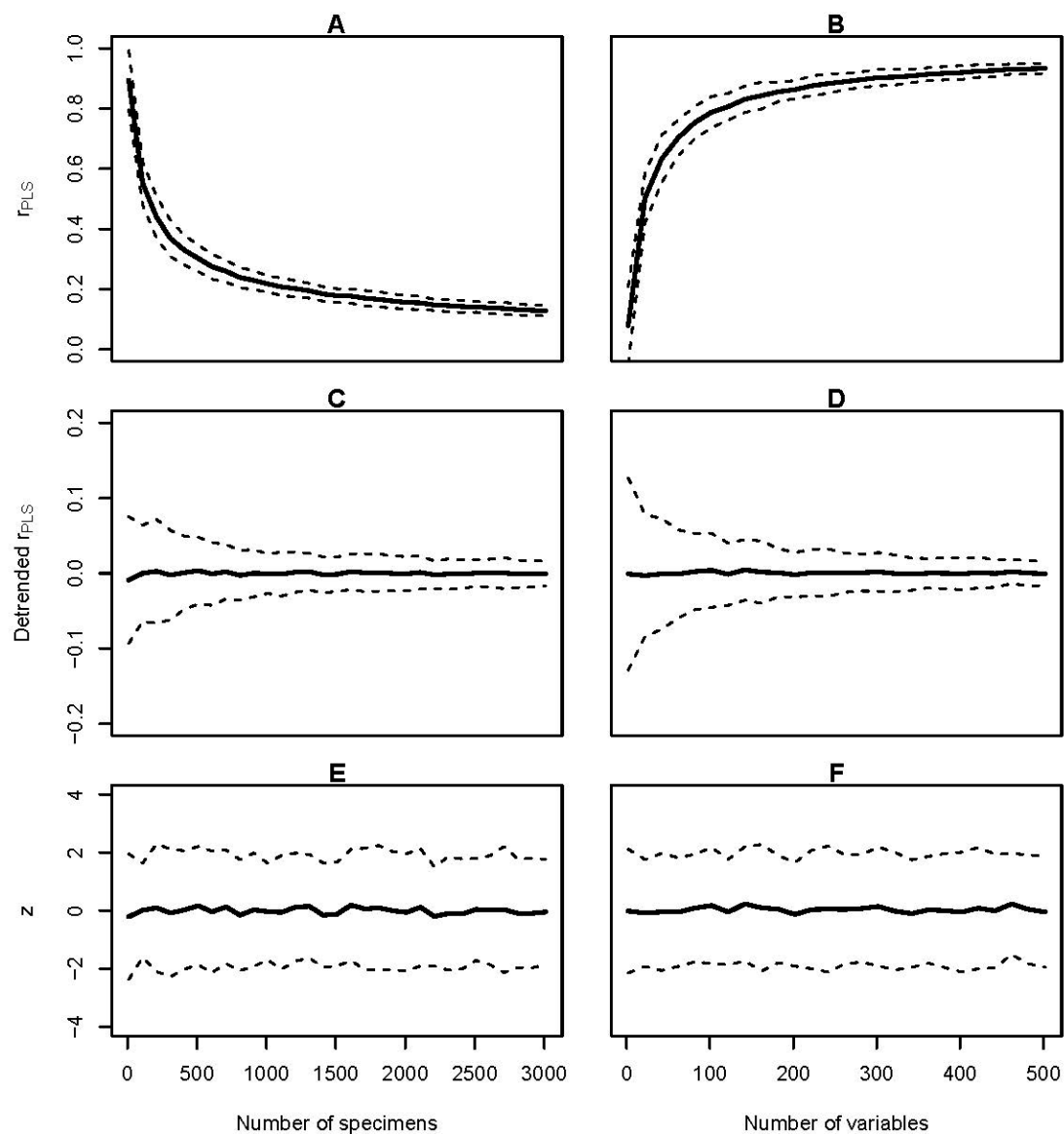
454

455

456

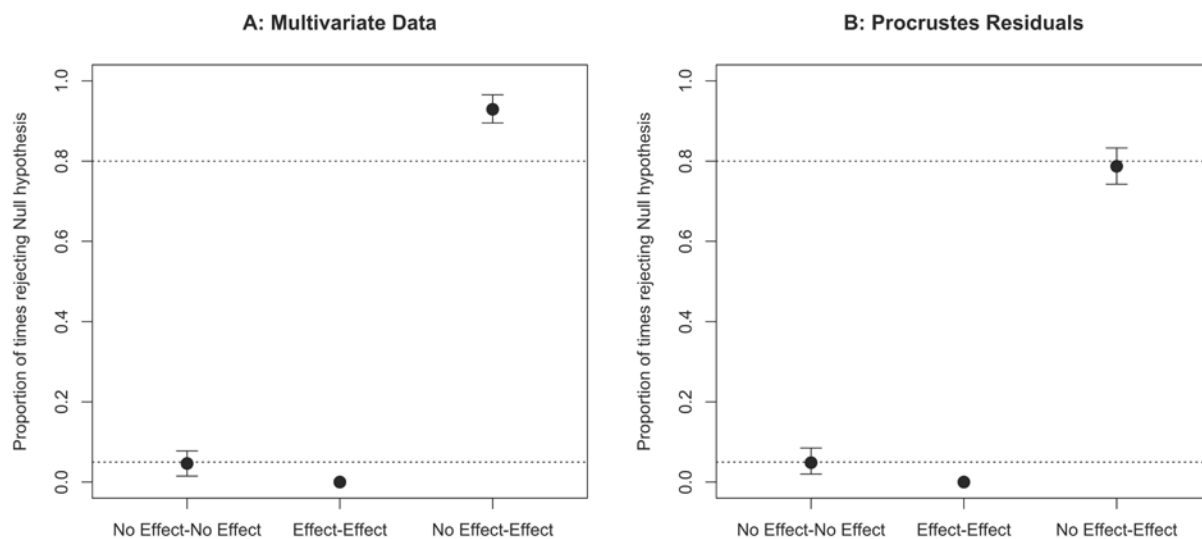
457

458



459
 460 **Fig. 1.** Evaluation of r_{PLS} under the hypothesis of random associations of variables. Mean and 95%
 461 confidence intervals of RV values obtained from (A) 100 datasets simulated across a range of sample
 462 sizes, and from (B) 100 datasets simulated across a range of variable number. Mean and 95% confidence
 463 intervals of detrended r_{PLS} for the same simulations of (C) 100 datasets simulated across a range of sample
 464 sizes, and from (D) 100 datasets simulated across a range of variable number. Mean and 95% confidence
 465 intervals of z -scores for the same simulations of (E) 100 datasets simulated across a range of sample sizes,
 466 and from (F) 100 datasets simulated across a range of variable number.
 467 .

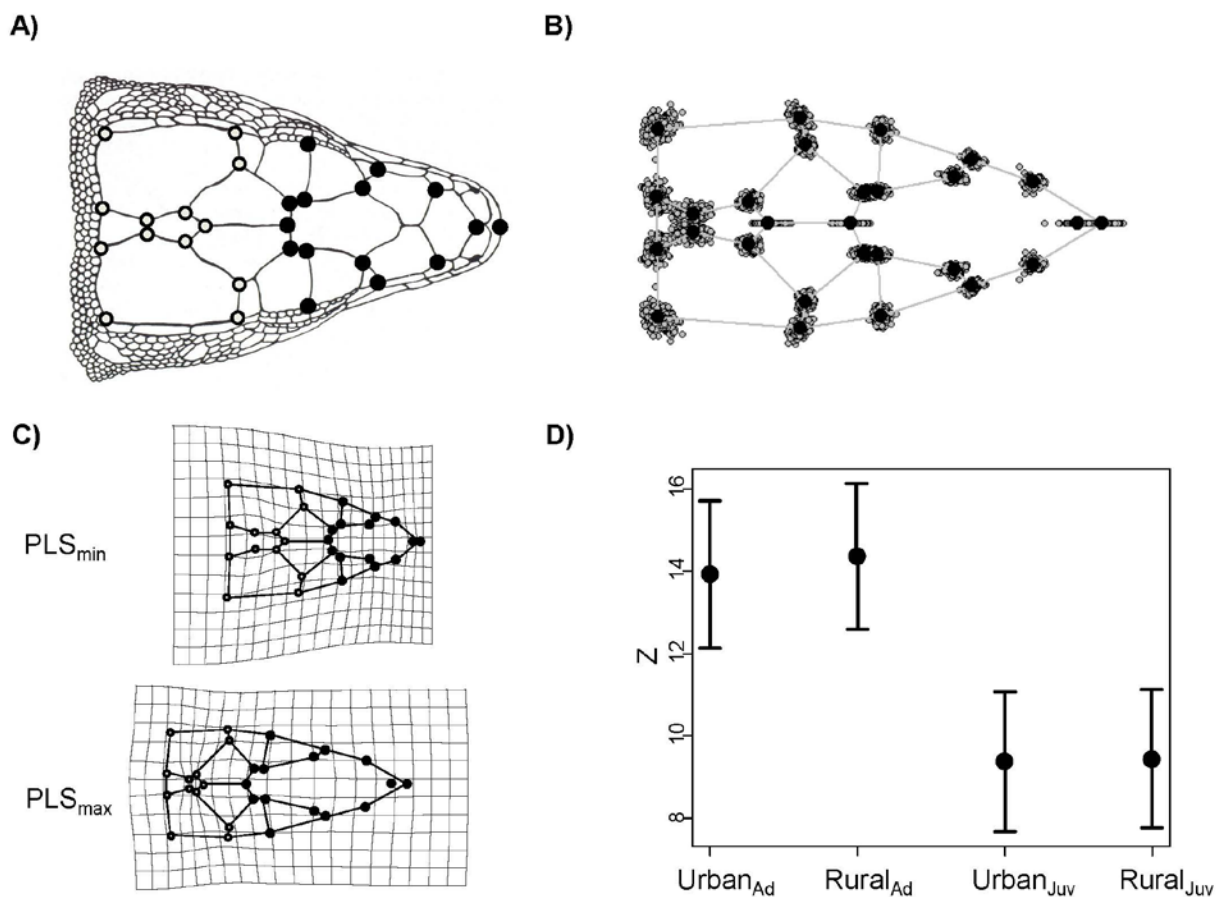
468



469

470 **Fig. 2.** Proportion of times the null hypothesis was rejected in 100 runs of 200 comparisons of
 471 morphological integration for (A) multivariate data and (B) Procrustes residuals. Samples of size, $n = 30$
 472 were obtained from populations ($N = 10,000$) with no integration (no effect) and substantial integration
 473 (effect). Means with 95% confidence limits are shown, unless the proportion was rather invariant, in
 474 which case, error bars are not included. Dotted lines are shown for an expected type I error rate of 0.05 or
 475 a statistical power of 0.80; within-population comparisons simulate type I errors (first two) and between-
 476 population comparisons simulate statistical power (last).

477



478
 479 **Fig. 3.** Graphical summary of results from the empirical example. (A) Locations of 28 landmarks on the
 480 dorsal view of a lizard head. Landmarks from the anterior module are designated by open circles while
 481 those from the posterior module are designated by closed circles (from Lazic et al. 2015). (B) Procrustes
 482 superimposition of all specimens (gray) with the mean specimen shown in black. (C) Thin-plate spline
 483 deformation grids representing specimens at the extremes of the PLS axis for the adult rural population.
 484 Deformation grids are accentuated by a factor of two to facilitate visual interpretation. (D) Levels of
 485 integration in juveniles and adults from both rural and urban populations shown as z-scores and their 95%
 486 confidence intervals obtained from the standard error of the permutation sampling distributions for each
 487 group.
 488
 489