

# A robust new metric of phenotypic distance to estimate and compare multiple trait differences among populations

Rebecca J. SAFRAN<sup>1\*§</sup>, Samuel M. FLAXMAN<sup>1§</sup>, Michael KOPP<sup>2§</sup>, Darren E. IRWIN<sup>3</sup>, Derek BRIGGS<sup>4</sup>, Matthew R. EVANS<sup>5</sup>, W. Chris FUNK<sup>6</sup>, David A. GRAY<sup>7</sup>, Eileen A. HEBETS<sup>8</sup>, Nathalie SEDDON<sup>9</sup>, Elizabeth SCORDATO<sup>10</sup>, Laurel B. SYMES<sup>11</sup>, Joseph A. TOBIAS<sup>9</sup>, David P. L. TOEWS<sup>3</sup>, J. Albert C. UY<sup>12</sup>

<sup>1</sup> Department of Ecology and Evolutionary Biology, University of Colorado, Boulder CO 80309, U.S.A.

<sup>2</sup> Mathematics and Biosciences Group, Faculty of Mathematics, University of Vienna, A-1090 Vienna, Austria; Current address: Evolutionary Biology and Modeling Group, Faculty of Sciences, Aix-Marseille University, 13331 Marseille, France.

<sup>3</sup> Department of Zoology, University of British Columbia, 6270 University Boulevard, Vancouver, BC, Canada, V6T 1Z4.

<sup>4</sup> Research & Evaluation Methodology, School of Education, University of Colorado, Boulder, CO 80309, U.S.A.

<sup>5</sup> School of Biological and Chemical Sciences, Queen Mary University of London, London, E1 4NS

<sup>6</sup> Department of Biology, Colorado State University, Fort Collins, CO 80523, U.S.A.

<sup>7</sup> Department of Biology, California State University, Northridge, CA 91330, U.S.A.

<sup>8</sup> School of Biological Sciences, University of Nebraska, Lincoln, NE 68588, U.S.A.

<sup>9</sup> Edward Grey Institute, Department of Zoology, University of Oxford, Oxford, OX1 3PS, UK

<sup>10</sup> Committee on Evolutionary Biology, The University of Chicago, Chicago, IL, 60637, U.S.A.

<sup>11</sup> Department of Biology, Dartmouth College, Hanover, NH 03755, U.S.A.

<sup>12</sup> Department of Biology, University of Miami, Coral Gables, FL 33146, U.S.A.

**Abstract** Whereas a rich literature exists for estimating population genetic divergence, metrics of phenotypic trait divergence are lacking, particularly for comparing multiple traits among three or more populations. Here, we review and analyze via simulation Hedges'  $g$ , a widely used parametric estimate of effect size. Our analyses indicate that  $g$  is sensitive to a combination of unequal trait variances and unequal sample sizes among populations and to changes in the scale of measurement. We then go on to derive and explain a new, non-parametric distance measure, " $\Delta p$ ", which is calculated based upon a joint cumulative distribution function (CDF) from all populations under study. More precisely, *distances are measured in terms of the percentiles in this CDF at which each population's median lies*.  $\Delta p$  combines many desirable features of other distance metrics into a single metric; namely, compared to other metrics,  $p$  is relatively insensitive to unequal variances and sample sizes among the populations sampled. Furthermore, a key feature of  $\Delta p$ —and our main motivation for developing it—is that it easily accommodates *simultaneous comparisons of any number of traits across any number of populations*. To exemplify its utility, we employ  $\Delta p$  to address a question related to the role of sexual selection in speciation: are sexual signals more divergent than ecological traits in closely related taxa? Using traits of known function in closely related populations, we show that traits predictive of reproductive performance are, indeed, more divergent and more sexually dimorphic than traits related to ecological adaptation [*Current Zoology* 58 (3): 426–439, 2012].

**Keywords** Effect size, Phenotype divergence, Sexual dimorphism, Sexual selection, Speciation

Inferences about the role of adaptation in population differentiation and speciation are often made by comparing phenotypic divergence and population genetic divergence. An active area of research and debate concerns the role of sexual selection in the process of speciation (e.g., Lande, 1981; West-Eberhard, 1983; Price, 1998; Panhuis et al., 2001; Boul et al., 2007; Ritchie, 2008; van Doorn et al., 2009; Kraaijeveld et al.,

2011). Whereas divergence in sexual traits is a common form of phenotypic differentiation among populations and sister taxa (e.g., Endler and Houde, 1995; Seehausen and van Alphen, 1999; Gray and Cade, 2000; Uy and Borgia, 2000; Irwin et al., 2001, 2008; Safran and McGraw, 2004; Rodríguez et al., 2004; Mendelson et al., 2005; Johnsen et al., 2006; Svensson et al., 2006; Boul et al., 2007; Uy et al., 2008; Seddon et al., 2008; Free-

Received Oct. 22, 2011; accepted Feb. 9, 2012.

\* Corresponding author. E-mail: rebecca.safran@colorado.edu. §Equal contribution.

© 2012 *Current Zoology*

man-Gallant et al., 2009), many questions still remain about how sexual signal divergence is related to speciation. In particular, researchers are interested in estimating differences in the extent of trait divergence that is underlain by ecological adaptation or sexual selection as a way to examine mechanisms that maintain modern population differences. In turn, such analyses can be used to infer a role of either natural or sexual selection in the process of divergence (Mayr, 1947; Maan and Seehausen, 2011).

An issue underlying all research concerning divergence among closely related populations concerns the metrics employed to examine estimates of both phenotypic and genetic distance. Whereas a rich and controversial literature exists for estimates of genetic distance (e.g. Wright, 1943, 1951, 1965, 1973, 1978; Slatkin 1987; Charlesworth, 1998; Excoffier, 2001; Charlesworth et al., 2003; Hedrick, 2005) there are relatively few resources for metrics of phenotypic distance, yet such metrics are fundamental for comparing trait differences among populations. In particular, the literature on metrics of phenotype distance, a sub-set of effect size metrics, (reviewed by Grissom and Kim, 2001; Nagawa and Cuthill, 2007) is focused on comparisons between two populations and on cases where traits follow the assumptions of parametric statistics (e.g., Grissom and Kim, 2001). Yet, limitations in the widespread utility of these metrics exist, particularly when trait distributions deviate from assumptions underlying parametric methods, when traits under study are measured in different units (e.g., size vs. color), when there are unequal sample sizes among groups under comparison, or when simultaneous analysis of more than one trait and/or more than two populations is desired.

Here, we offer a non-parametric and potentially powerful new metric, " $\Delta p$ ", that overcomes the aforementioned limitations. We analyze the performance of  $\Delta p$  by comparing its behavior to that of a very commonly employed effect size metric, Hedges'  $g$  (Hedges, 1981), which is a variant of Cohen's  $d$  (Cohen, 1969) and belongs to a class of parametric effect size measures that essentially calculate a difference in means, scaled (divided) by some measure of the standard deviation in one or both groups being compared (see Grissom and Kim, 2001; Nagawa and Cuthill, 2007). (In fact,  $g$  and  $d$  are practically identical metrics, with the only difference being that  $d$  does not utilize the "-2" correction seen below in the denominator of equation (2).) As we illustrate, the behavior of such metrics can be sensitive to unequal variances and sample sizes among groups being

compared. We show that  $\Delta p$  does not have this sensitivity, that its behavior is at least as reliable as that of  $g$  for both normally and non-normally distributed data sets, and that it offers the crucial, additional advantage of being amenable to comparisons involving more than two populations and/or two or more traits simultaneously.

We note here that, like other distance measures, the measures we present below are descriptive. Hence, we suggest the following protocol: standard statistical methods are first used to establish the significance of differences between populations. Then, the method we present below can be used to generate quantitative descriptions of differences between two or more populations for examining questions about (1) the degree of phenotypic divergence of a single trait, (2) the overall degree of phenotypic differentiation across all traits being considered, (3) the degree of sexual dimorphism in a trait within populations (if applicable) relative to the phenotypic differences between populations, and (4) the ranking of traits, regardless of their units of measure, in order of which traits are most phenotypically divergent.

After explaining the newly derived phenotype distance metric, we illustrate its utility by applying it to a number of empirical data sets where the function of traits in either a sexual signaling or ecological adaptation context has been previously explored, such that we can compare trait distance between populations for both sexual and ecological traits. We also use this metric to make comparisons of males to females within closely related populations to test a prevailing yet largely untested assumption about using sexual dimorphism as a proxy of sexual selection (e.g., Kraaijeveld et al., 2011). Here, we can address whether known sexual signals - when they are present in both males and females - are more sexually dimorphic than ecological traits. The overall goal of this contribution is to present methodological recommendations, so that appropriate metrics of trait distance are available for making comparisons among closely related populations.

*Limitations of Hedges'  $g$*  In the following, we outline some of the limitations of Hedges'  $g$  that have motivated the development of our new metric  $\Delta p$ . Many of our points have also been made by Grissom and Kim (2001). We wish to emphasize that our paper is not meant as a general critique of Hedges'  $g$ , which has many useful properties: one only needs to know means, standard deviations, and sample sizes to calculate it, it is unit-less, and it has achieved widespread usage in a va-

riety of disciplines. Indeed, it is because of the widespread acceptance of  $g$  that we choose to use it for comparisons here: the behavior of  $g$  sets a standard that a newly proposed metric should meet and exceed.

Hedge's  $g$  is computed as

$$g_{y,zj} = \frac{\bar{x}_{yj} - \bar{x}_{zj}}{s_{y,zj}^*} \quad (1)$$

where the subscripts on  $g$  in equation (1) indicate that the calculation of  $g$  between two populations, denoted  $y$  and  $z$ , was for the  $j$ th trait that was measured in these populations.  $\bar{x}_{yj}$  and  $\bar{x}_{zj}$  are sample means for the  $j$ th trait in the two sampled populations ( $y$  and  $z$ , respectively) and  $s_{y,zj}^*$  is a measure of pooled sample standard deviation. This measure is weighted by sample size and is defined as

$$s_{y,zj}^* \equiv \sqrt{\frac{(n_{yj} - 1)S_{yj}^2 + (n_{zj} - 1)S_{zj}^2}{n_{yj} + n_{zj} - 2}}, \quad (2)$$

where  $n_{yj}$  and  $n_{zj}$  are the sample sizes of observations of the  $j$ th trait in the two populations and  $S_{yj}^2$  and  $S_{zj}^2$  are the sample variances (see Table 1 for definitions of all symbols). Note that Hedges and Olkin (1985) give an additional correction factor for  $g$ , which should be applied if the overall sample size is small.

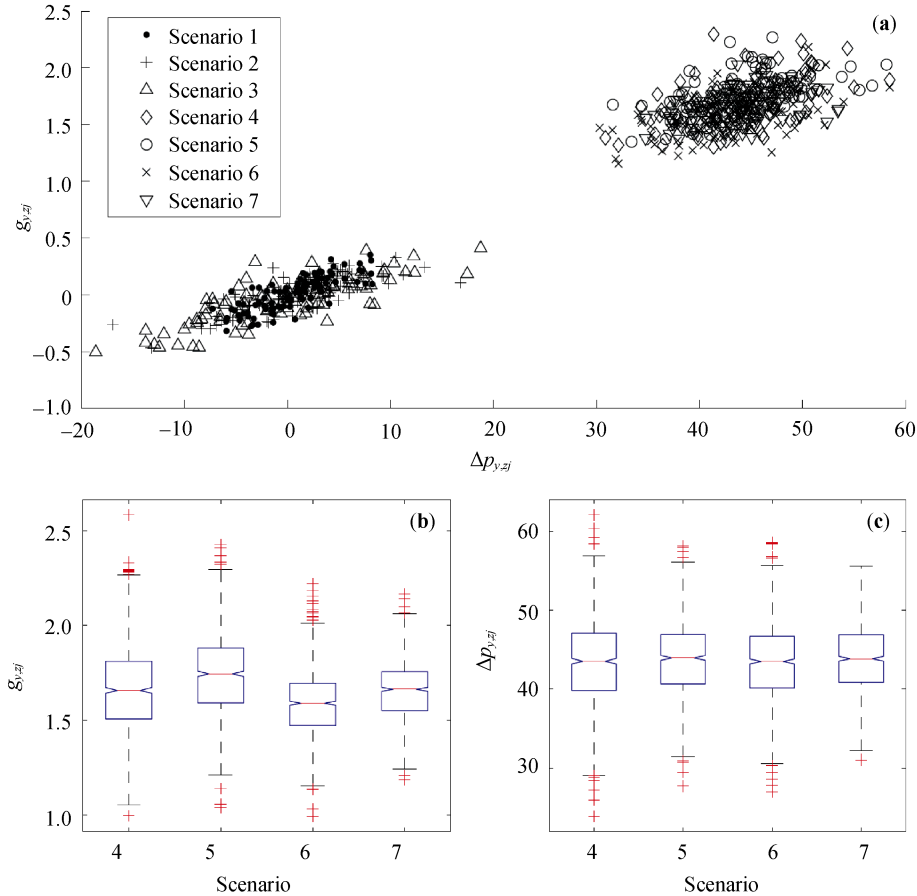
As mentioned above,  $g_{y,zj}$  has properties that limit its utility in certain situations.

First, it assumes that the trait has the same "true" variance in both populations. Indeed, the term under the square-root sign in equation (2) is an unbiased estimator for this variance (as it is in the two-sample  $t$ -test with equal variance, e.g., Sokal and Rohlf 1995). If, in contrast, the (true) variances in the two populations are different, equation (1) cannot be applied, since the denominator (eq. 2) has no useful interpretation and its expected value will depend on sample sizes. The latter point can be seen by replacing the empirical variances in equation (2) with their "true" counterparts. Then, increasing the sample size for the population with the larger (smaller) true variance will increase (decrease)  $s_{y,zj}^*$  and decrease (increase)  $g_{y,zj}$ . This effect is also illustrated in Fig. 1.

A possible solution to the problem of variance heterogeneity is to define an alternative distance measure in which equation (2) is replaced by the square root of an unweighted average of the sample variances (in analogy to the  $t$ -test with unequal variances). However, the interpretation of such a measure poses conceptual difficulties, since the difference in sample means is scaled by a "virtual" standard deviation that does not apply to any real population (Grissom and Kim, 2001).

**Table 1** Definitions and explanations of notation employed

Symbol	Meaning	Value(s) or range assigned (if applicable)
$N$	Number of populations or groups being compared	Integer, $\geq 2$
$t$	Number of traits measured in each population or group	Integer, $\geq 1$
$i$	Index variable for populations	$i = 1, 2, \dots, N$
$j$	Index variable for traits	$j = 1, 2, \dots, t$
$n_{ij}$	Number of observations of $j$ th trait in $i$ th population	Integer, $> 0$
$k$	Index variable for observations	$k = 1, 2, \dots, n_{ij}$
$x_{ijk}$	$k$ th observation of $j$ th trait in $i$ th population	Empirically determined
$\bar{x}_{ij}$	Sample mean of $j$ th trait in $i$ th population	Empirically determined
$\hat{x}_{ij}$	Sample median of $j$ th trait in $i$ th population	Empirically determined
$S_{ij}^2$	Sample variance of $j$ th trait in $i$ th population	Empirically determined
$g_{y,zj}$	Hedges' $g$ statistic computed for the $j$ th trait measured in populations $y$ and $z$ ( $y, z \in \{1, 2, \dots, N\}$ )	See equation (1)
$s_{y,zj}^*$	Pooled standard deviation used in calculation of Hedges' $g$ statistic	See equation (2)
$d$	Cohen's $d$ statistic	See text for description
$p_j(u)$	cumulative distribution function for trait $j$ , expressed as a percentage	See equation (3)
$\Delta p_{y,z}$	Distance between populations $y$ and $z$ , calculated over all traits	See equation (5)
$\Delta p_{y,zj}$	Distance between populations $y$ and $z$ , calculated for trait $j$	See equation (4)



**Fig. 1** Comparisons between the behavior of  $g_{y,zj}$  and  $\Delta p_{y,zj}$  using simulated data following normal distributions

1000 data sets were generated with the parameters listed for each “Scenario” (Table 2). (a):  $g_{y,zj}$  and  $\Delta p_{y,zj}$  are strongly correlated with each other (Pearson’s  $\rho = 0.98$ ). In Scenarios 1-3, there was no difference between the populations; in Scenarios 4-7 the means and standard deviations of the populations truly differed. For visual clarity, only 100 randomly selected points from each scenario are plotted here. (b) and (c): Box-and-whisker-plots of  $g_{y,zj}$  and  $\Delta p_{y,zj}$  (respectively) in the four scenarios in which the two mock populations truly differed in their underlying means (see Table 2). In these plots, the centerline is the median value of the metric, the box shows the interquartile range (IQR), the whiskers extend to up to  $1.5 \times \text{IQR}$  beyond the box, and the “+” symbols show points outside the latter range.

Furthermore, if one wishes to compare more than two populations (e.g. the three possible pairwise comparisons between three populations), each pairwise difference will be scaled by a different pooled standard deviation, rendering it awkward if not impossible to make meaningful quantitative comparisons between them.

A second limitation of  $g_{y,zj}$  is that its value depends on the scale of measurement. For example, researchers will often apply non-linear transformations (such as log or arcsin transforms) to make their data meet the assumptions of parametric statistical methods (including normality and variance homogeneity). However, such transformations will also alter the calculated values of  $g_{y,zj}$ . This may be problematic if one wishes to compare distance measures for different traits and only some of the traits have been transformed or different traits have

different natural scales of measurement (e.g., additive vs. multiplicative). Indeed, problems of this kind may often occur in sexual selection research when researchers aim to compare the divergence of naturally- versus sexually-selected traits (e.g., size vs. color).

In some cases, instead of comparing the divergence of different traits, one might want to have a single divergence measure involving multiple traits. Such a measure is given by the Mahalanobis distance (Mahalanobis, 1936; Arnegard et al., 2010), which may be seen as a multivariate generalization of  $g_{y,zj}$  that also takes into account correlations between traits. However, the Mahalanobis distance faces the same restrictions as  $g_{y,zj}$ , that is, it requires a single estimate of the variance-covariance matrix for all populations, and its exact value will depend on scale(s) of measurement.

In summary, the limitations of  $g_{y,zj}$  regarding unequal variances and scales of measurement not only affect pairwise comparisons, but also limit its applicability for comparisons involving multiple traits and/or populations.

A more useful distance metric would work for simple pairwise comparisons involving two populations and a single trait, but would also work for considering more than two populations and more than one trait simultaneously, so that (1) all pairwise effect sizes, even for traits measured in different units, would all be on the same scale, and (2) measures of overall distance (involving all traits at once) could be computed.

In the following sections, we introduce the derivation of  $\Delta p$  and via simulation and examples using empirical data, present its utility in a number of contexts in which a flexible measure of trait distance is required.

## 1 Methods

### 1.1 A nonparametric distance measure for arbitrary numbers of traits and population

Our new metric  $\Delta p$  does not make any assumptions about trait distributions or variances (i.e., it is nonparametric). Instead, it is based upon the joint (average) cumulative distribution function (cdf) across all populations for a given trait. Suppose there are  $N$  populations and  $t$  traits being considered. Let  $x_{ijk}$  denote the  $k$ th observation of the  $j$ th trait in the  $i$ th population, and  $n_{ij}$  the number of samples for trait  $j$  taken from population  $i$ . The joint empirical cdf for trait  $j$  (expressed in percentiles) is defined as

$$p_j(u) = \frac{100}{N} \sum_{i=1}^N \frac{1}{n_{ij}} \sum_{k=1}^{n_{ij}} \mathbf{1}\{x_{ijk} \leq u\} \quad (3)$$

where  $u$  is any given value of trait  $j$ , and  $\mathbf{1}\{\bullet\}$  is an indicator function that returns 1 if its argument is true and 0 otherwise. For illustration, imagine that the data for trait  $j$  from all populations have been pooled and sorted in increasing order. Let  $x_{\min,j}$  denote the global minimum and  $x_{\max,j}$  the global maximum.  $p_j(u)$  is a step function which starts out at zero (for  $u < x_{\min,j}$ ), jumps up by  $100/(N n_{ij})$  at each  $x_{ijk}$ , and reaches 100 at  $u = x_{\max,j}$ . Importantly, by making the height of the jumps inversely proportional to the size of the sample a given data point stems from, we make sure that each sampled population contributes equally to  $p_j(u)$ , independent of sample or population size (e.g., for  $N = 2$ , each population is responsible for 50% of the total increase in  $p_j(u)$ ).

Returning to the individual populations, we then ask: Into what percentile in the overall CDF does the median

of each population fall? In other words, we calculate the value of  $p_j(\hat{x}_{ij})$ , where  $\hat{x}_{ij}$  is the median value of trait  $j$  in population  $i$ , and we repeat this for all  $N$  populations. Our measure of phenotypic distance between populations  $y$  and  $z$  with respect to trait  $j$  is then defined as

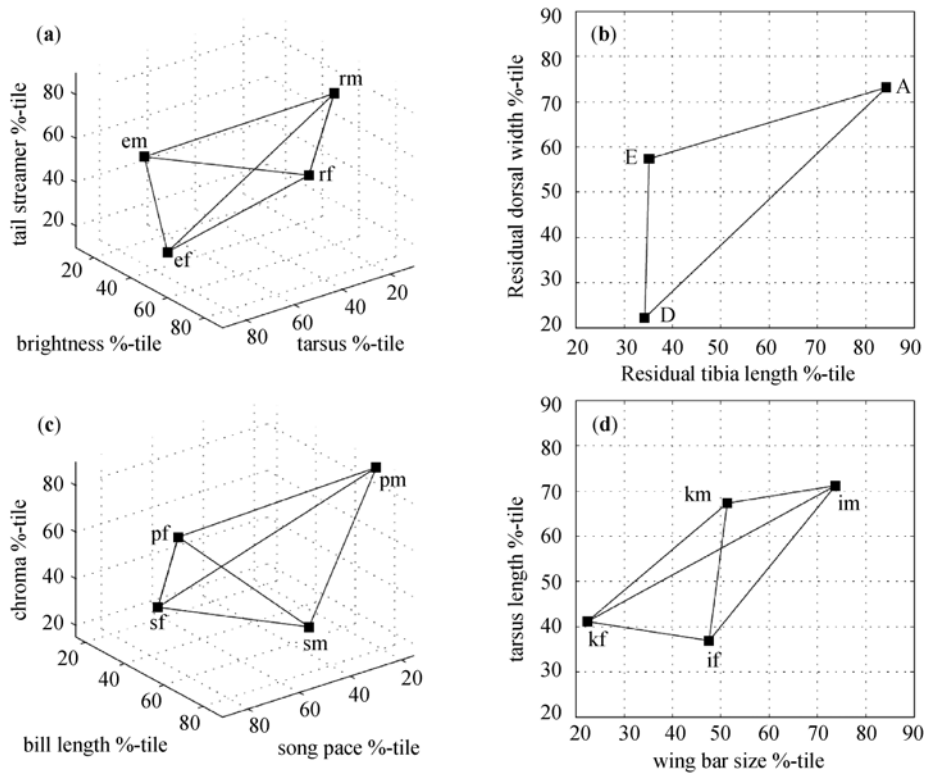
$$\Delta p_{y,zj} \equiv p_j(\hat{x}_{yj}) - p_j(\hat{x}_{zj}). \quad (4)$$

As with  $g_{y,zj}$ ,  $\Delta p_{y,zj}$  can be positive or negative, in this case depending on whether population  $y$  or population  $z$  has the larger median. Importantly—and in contrast to  $g_{y,zj}$  and other phenotypic distance measures—if there are more than two populations, all pairwise  $\Delta p_{y,zj}$  values will be based on the same overall CDF, and hence, will be directly comparable. A numerical example outlining the above calculations is included in the online supplementary materials (Appendix 1) as a Microsoft Excel spreadsheet.

If data are available for more than one trait, the above analysis can be repeated for each trait separately. Again, the results will be comparable, because each phenotypic distance is measured at the appropriate scale (i.e., with respect to the overall CDF for that trait). In addition, we can also define an overall phenotypic distance for a pair of populations considering all traits simultaneously. The idea is to view the  $p_j(\hat{x}_{ij})$  values (with  $j = 1, \dots, t$ ) of a single population as a set of “coordinates” for that population in a  $t$ -dimensional trait-percentile space (illustrated in Fig. 2). The coordinates for different populations naturally lend themselves to a notion of distance: for any number of traits being considered, we can calculate a Euclidean distance between the two populations using their percentile coordinates. We denote this distance between two groups or populations as

$$\Delta p_{y,z\bullet} \equiv \sqrt{\sum_{j=1}^t (p_j(\hat{x}_{yj}) - p_j(\hat{x}_{zj}))^2}, \quad (5)$$

where  $y, z \in \{1, 2, \dots, N\}$  refer to two of the populations that were measured, and the subscript “ $\bullet$ ” denotes that this distance is calculated over all traits. Note that, unlike the single trait distance  $\Delta p_{y,zj}$ ,  $\Delta p_{y,z\bullet}$  is always positive; however, in the limiting case of just a single trait (i.e.,  $j = t = 1$ ), from equation (5) we obtain  $\Delta p_{y,z\bullet} = |\Delta p_{y,z1}|$ . Note also that  $\Delta p_{y,z\bullet}$  is expected to increase as more traits are added to the analysis (as is true of any Euclidean distance as more dimensions are added).



**Fig. 2 Representations of population locations in trait-percentile spaces**

(a) is based upon data on three traits from males and females (with sexes considered as separate “groups”) for each of two subspecies of barn swallows *Hirundo rustica*. “em” and “ef” represent *H. r. erythrogaster* males and females (respectively; sample sizes = 71-98 for the three traits); “rm” and “rf” represent *H. r. rustica* males and females (sample sizes = 32-76 for the three traits). (b) Data on two traits from females in three species of painted forest toadlets (“A”, “D”, and “E” refer, respectively, to *Engystomops petersi* sp. A and sp. D and *E. freibergeri*). Sample sizes for both traits for A, D, and E are 9, 9, and 36, respectively. (c) Data on *Hypocnemis peruviana* males (“pm”) and females (“pf”) and *Hypocnemis subflava* males (“sm”) and females (“sf”). Sample sizes vary from 20-27 for both sexes for bill length and chroma; sample sizes for song pace are 5 (sf), 9 (pf), 17 (pm), and 21 (sm). (d) Data on populations of Hume’s Warblers from Kyrgyzstan (males = “km”; females = “kf”; sample sizes = 39 for both sexes for both traits) and India (males = “im”, sample size = 56; females = “if”, sample size = 45). Information about the data contained in these figures can be found in online Appendix 3.

We note that for  $t > 1$ , the interpretation of  $\Delta p_{y,z\bullet}$  as a Euclidean distance neglects correlations between traits within samples (e.g., correlations between traits  $j$  and  $j+1$  within population  $y$ ). Our justification is that, while the amount and direction of divergence relative to within-population correlations poses some extremely important questions (e.g. Schluter 1996), it does not seem possible to derive a single metric that captures all aspects of this problem, especially if the orientation of principal components differs between populations. (If the orientation of principle components is similar across populations, then if significant correlations between traits exist,  $\Delta p_{y,z\bullet}$  can be calculated using results from principle components analysis.)

We have developed a MATLAB script to import data, perform all of the above calculations, and automatically generate a .csv file of results on Hedges’  $g$  values,  $\Delta p_{y,z\bullet}$ ,  $\Delta p_{y,zj}$ , and other useful descriptive statistics.

The commented source code (.m files), along with plain text explanations, metadata, and example input and output files are all freely available from the second author (SMF) upon request and have also been archived at SourceForge.net (<http://sourceforge.net/projects/deltap/files/>). Details of the derivation of confidence intervals are located in online Appendix 4.

## 1.2 Evaluating the performance of $\Delta p$

In order to explore and illustrate the behavior of  $\Delta p_{y,z\bullet}$ , we used numerical simulations to generate pseudo-random data sets on hypothetical traits from mock populations. We then applied the above methods to these data sets to explore realistic scenarios involving equal and unequal means for traits among populations, unequal sample sizes among populations, and unequal variances in a trait among populations (Table 2). We compared the behavior of  $\Delta p_{y,z\bullet}$  in these “Scenarios” (Table 2) to the behavior of  $g_{y,zj}$ . By necessity, for the

purposes of directly comparing  $g_{y,zj}$  and  $\Delta p_{y,zj}$ , we considered only two mock populations and a single hypothetical trait (since that is all that can be used in a single calculation of  $g_{y,zj}$ ); that is, we compare  $g_{y,zj}$  and  $\Delta p_{y,zj}$ .

For these comparisons, we generated pseudo-random data following normal and exponential distributions. We show results below for several scenarios involving normally distributed data; additional scenarios with different sample sizes, means, standard deviations, and non-normal data are given in the supplemental materials (see online Appendix 2). For each “Scenario” in Table 2, we generated 1000 pairs of random samples from the two populations, each with the specified sample size and following the specified distribution. For each Scenario, we could thus calculate  $g_{y,zj}$  and  $\Delta p_{y,zj}$  1000 times, independently. Analysis of variance (ANOVA) was performed to compare whether  $g_{y,zj}$  systematically differed in different Scenarios, and likewise for  $\Delta p_{y,zj}$ .

### 1.3 Applying $\Delta p$ to empirical data sets involving traits of known function

To demonstrate the application of  $\Delta p$  we solicited data from researchers working on systems where traits related to sexual signaling and ecological adaptation are well characterized. Criteria for inclusion of data were as follows: 1) the underlying mechanisms generating trait variation are fairly well-understood, such that one trait can be assumed to be predominantly underlain by natural selection via ecological adaptation and another to be predominantly underlain by sexual selection via variation in reproductive performance. 2) Data are from

closely related taxa, ranging from sister taxa to sub-species to geographically isolated populations. See online Appendix 3 for details about the individual study systems, and the field and lab methods used to generate the unpublished data given in Tables 3 and 4; references for published data are given when available in Tables 3 and 4.

## 2 Results

### 2.1 Evaluating the performance of $\Delta p$

In simulations comparing the performance of  $\Delta p_{y,zj}$  with  $g_{y,zj}$ , two important categories of results emerged (Fig. 1). First, if  $\Delta p_{y,zj}$  is a valid distance metric, it should reproduce some aspects of the behavior of the well-established metric,  $g_{y,zj}$ . This was indeed the case: (1)  $\Delta p_{y,zj}$  and  $g_{y,zj}$  were very tightly correlated, (2) they both were centered on zero for cases when populations did not truly differ (Scenarios 1–3 in Figure 1a), and (3) they were both larger than zero for cases when the means of the populations truly differed (Scenarios 4–7 in all panels of Fig. 1). Second, however, we also found that  $g_{y,zj}$  is much more sensitive to the combination of unequal variances and sample sizes (e.g., Grissom and Kim 2001; Fig. 1b) than  $\Delta p_{y,zj}$  (Fig. 1c). In Scenarios 3–7, the true difference between populations was constant, but the simulated values of  $g_{y,zj}$  varied systematically depending upon which population had the larger sample size and which had the larger variance (ANOVA on data in Fig. 1b:  $F_{3,996}^3 = 104.75$ ,  $P < 10^{-64}$ ). When the population with the smaller variance is sampled the most, then  $g_{y,zj}$  will tend to overestimate the distance between populations; when the population with the larger variance is sampled the

**Table 2** Parameters<sup>a</sup> used in simulated “Scenarios” used to compare  $g_{y,zj}$  and  $\Delta p_{y,zj}$

Scenario	Population y			Population z		
	$n_{yj}$	$\mu_{yj}^b$	$\sigma_{yj}^c$	$n_{yz}$	$\mu_{yz}^b$	$\sigma_{yz}^c$
1: unequal sample sizes only	56	103.36	8.89	120	103.36	8.89
2: unequal sample sizes only; means are equivalent but differ from scenario 1	56	90.27	6.77	120	90.27	6.77
3: unequal variances only	56	90.27	6.77	56	90.27	8.89
4: unequal means and variances only	56	103.36	8.89	56	90.27	6.77
5: unequal means, variances, and sample sizes	56	103.36	8.89	120	90.27	6.77
6: as in 5, but with reversed sample sizes	120	103.36	8.89	56	90.27	6.77
7: as in 4, but with larger sample sizes	120	103.36	8.89	120	90.27	6.77

<sup>a</sup> Sample sizes, means, and standard deviations were all inspired by a real data set on tail streamer lengths (in mm) for two subspecies of barn swallows (Safran and Evans, unpublished data).

<sup>b</sup> Assumed true population mean used for generating pseudorandom data.

<sup>c</sup> Assumed true population standard deviation used for generating pseudorandom data.

Table 3 Side by side comparisons of  $g_{y,f}$  and  $\Delta p_{y,f}$  for sexual vs. ecological traits in two closely related populations

Species	Comparison <sup>1</sup>	Sex Trait	$g_{y,f}$	$\Delta p_{y,f}$	Ecological Trait	$g_{y,f}$	$\Delta p_{y,f}$	source
Barn swallows <i>Hirundo rustica</i>	<i>H. r. erythrogaster</i> vs. <i>H. r. rustica</i> $n_{1s} = 71, n_{2s} = 53, n_{1e} = 82, n_{2e} = 76$	Tail length	-1.80	-46.36 (-49.18, -40.49)	Tarsus length	1.53	39.25 (29.67, 48.44)	Safran and Evans <i>unpubl</i>
Hume's warbler <i>Phylloscopus humei</i>	Kyrgyztan vs India $n_{1s} = n_{1e} = 39; n_{2s} = n_{2e} = 56$	Wing bar size	-0.71	-28.50 (-38.10, -13.28)	Tarsus length	-0.25	-1.75 (-29.95, 18.57)	Scordato <i>unpubl</i>
Greenish warblers <i>Phylloscopus trochiloides</i>	<i>P. viridanus</i> vs. <i>P. plumbeitarsus</i> $n_{1s} = 5, n_{2s} = 9, n_{1e} = 12, n_{2e} = 15$	Song units	-3.96	-47.78 (-57.89, -27.78)	Tarsus length	0.24	11.67 (-29.17, 40.83)	Irwin et al., 2001; Irwin et al., 2009
Pacific / Winter wrens <i>Troglodytes pacificus / hiemalis</i>	<i>T. hiemalis</i> vs. <i>T. pacificus</i> $n_{1s} = 13, n_{2s} = 19, n_{1e} = 9, n_{2e} = 34$	Song freq	-3.62	-49.39 (-54.66, -37.85)	Tarsus length	0.36	12.91 (-24.35, 44.28)	Toews and Irwin, 2008
Tree crickets <i>Oecanthus forbesi</i>	Wisconsin vs. Ohio $n_{1s} = 8, n_{2s} = 14, n_{1e} = 8, n_{2e} = 13$	File tooth #	-1.03	-40.18 (-64.29, -16.96)	Underwing length	-0.24	-3.87 (-40.38, 35.10)	Symes <i>unpubl</i>
Field Crickets <i>Gryllus</i>	<i>G. texensis</i> vs <i>G. rubens</i> <sup>2</sup> $n_{1s} = 164, n_{2s} = 122, n_{1e} = 119, n_{2e} = 102$	Pulse rate	4.79	50.91 (49.38, 54.27)	Ovipositor length	-2.40	-47.90 (-49.37, -44.75)	Izzo and Gray, 2004; Gray et al., 2001; Gray <i>unpubl</i>
Painted forest toadlet <i>Engystomops petersi</i>	<i>E. petersi</i> sp. A vs. sp. D. $n_{1s} = 4, n_{2s} = 12, n_{1e} = 14, n_{2e} = 19$	Call dom freq	-2.43	-50.00 (-66.67, -29.17)	Tibia length	1.85	41.54 (10.53, 48.68)	Boul et al., 2007; Funk et al., 2008; Funk <i>unpubl</i> .
Warbling antbirds <i>Hypocnemis</i> spp.	<i>H. peruviana</i> vs. <i>H. subflava</i> $n_{1s} = 17, n_{2s} = 21, n_{1e} = 20, n_{2e} = 27$	Song pace	-1.16	-33.75 (-45.67, -4.76)	Bill length	0.34	16.11 (-13.61, 41.02)	Tobias and Seddon, 2009 (song); Seddon and Tobias <i>unpubl</i> (morphology)
Manakins <i>Manacus</i>	<i>M. vitellinus</i> vs. <i>M. candei</i> $n_{1s} = n_{1e} = 17, n_{2s} = n_{2e} = 15$	Plumage brightness	-3.85	-50.20 (-56.86, -44.81)	Bill length	-0.24	0 (-25.10, 40.59)	Stein and Uy, 2006; Uy <i>unpubl</i>
Wolf spiders <i>Schizocosa</i>	<i>S. bilineata</i> vs <i>S. crassipalpata</i> $n_{1s} = 87, n_{2s} = 51, n_{1e} = 86, n_{2e} = 52$	Leg length	-1.25	-41.68 (-47.50, -33.40)	Cephalothorax width	-1.38	-43.43 (-47.29, -34.35)	Hebets <i>unpubl</i>

<sup>1</sup> In the "Comparison" column,  $n_{1s}$  denotes the sample size for the sexual trait in the first named population in the comparison,  $n_{1e}$  the sample size for the ecological trait in that population,  $n_{2s}$  the sample size for the sexual trait in the second named population, and  $n_{2e}$  the sample size for the ecological trait in the second population.

<sup>2</sup> For this comparison only, the sex trait occurs in males only and the ecological trait in females only.

Comparisons are males vs. males with the exception of *Gryllus* crickets in which case only males possess the sexual signal and females the ecologically relevant trait. A negative value of  $g_{y,f}$  and  $\Delta p_{y,f}$  indicates that the first population named has a smaller mean trait value compared to the second. For example, in the case of barn swallows, the length of tail streamers in the subspecies *erythrogaster* is shorter compared to *rustica*. Numbers in parentheses following  $\Delta p_{y,f}$  are 95% confidence intervals (see methods).



Table 4 Using  $\Delta P_{y,f}$  to compare sexual and ecological trait dimorphism

Species	Comparison	Dimorphism Sex Trait	$\Delta P_{y,f}$		Dimorphism Ecological Trait	$\Delta P_{y,f}$		source
			Pop 1	Pop 2		Pop 1	Pop 2	
Barn swallows <i>Hirundo rustica</i>	<i>H. r. rustica</i> vs. <i>H. r. erythrogaster</i>	Tail length	35.40 (29.22, 42.92)	38.24 (34.01, 44.70)	Tarsus length	-8.79 (-26.10, 0.00)	0.56 (-8.06, 8.25)	Safran and Evans, <i>unpubl.</i>
Hume's warbler <i>Phylloscopus humei</i>	Kyrgyzstan vs. India	Wing bar size	28.96 (15.74, 41.83)	26.13 (10.42, 41.28)	Tarsus length	26.19 (8.86, 41.13)	34.24 (13.67, 48.06)	Scordato, <i>unpubl.</i>
Warbling antbirds <i>Hypocnemis</i> spp.	<i>H. peruviana</i> vs. <i>H. subflava</i>	Song pace	-62.08 (-77.24, -39.94)	-42.43 (-60.92, -30.24)	Bill length	20.23 (-13.70, 42.69)	21.81 (-1.04, 44.77)	Tobias and Seddon, 2009 (song); Seddon and Tobias, <i>unpubl.</i> (morphology).
Mankins <i>Manacus</i>	<i>M. vitellinus</i> vs. <i>M. candei</i>	Plumage brightness	25.74 (20.74, 28.68)	74.70 (70.15, 78.03)	Bill length	0 (-16.08, 47.04)	-31.41 (-56.92, 0.00)	Stein and Uy, 2006; Uy, <i>unpubl.</i>
Wolf spiders <i>Schizocosa</i>	<i>S. bilineata</i> vs <i>S. crassipalpata</i>	Leg length	7.79 (-16.38, 25.41)	33.11 (21.49, 41.42)	Cephalothorax width	-47.69 (-54.13, -37.04)	-35.95 (-47.31, -24.01)	Hebets, <i>unpubl.</i>

In each analysis shown below, four distinct groups were included simultaneously (2 sexes  $\times$  2 populations). A positive value indicates that males have the larger mean value for a trait; a negative value indicates that females have the larger mean value.

most,  $g_{y,zj}$  will tend to underestimate distance between populations. With large numbers of repeated simulations, slight differences could also be detected for  $\Delta p_{y,zj}$  (ANOVA on data in Fig. 1c:  $F_{3,996}^3 = 3.36$ ,  $P < 0.02$ ). However, the degree of this sensitivity was an order of magnitude less for  $\Delta p_{y,zj}$  than for  $g_{y,zj}$ : mean values of  $g_{y,zj}$  ranged from 1.589 in Scenario 6 to 1.744 in Scenario 5, a difference of 10%. By contrast mean values of  $\Delta p_{y,zj}$  ranged from 43.34 in Scenario 5 to 43.83 in Scenario 5, a difference of only 1%.

## 2.2 Applying $\Delta p$ to empirical data sets involving traits of known function

The results summarized in Fig. 2 and Table 3 exemplify the utility of  $\Delta p_{y,zj}$  and  $\Delta p_{y,z\bullet}$ . Namely, multiple traits measured in different units, from across multiple populations can be compared simultaneously. In Figures 2a, 2c, and 2d, simultaneous comparisons of males from different populations, females from different populations, and dimorphism within populations can all be made. For example, Fig. 2a shows the extent to which barn swallow tail streamers are (i) sexually dimorphic in both populations (compare “ef” to “em” and “rf” to “rm” on the z-axis), (ii) divergent across populations (compare “em” to “rm” and “ef” to “rf”), and (iii) similar between *H. r. rustica* females and *H. r. erythrogaster* males (compare “rf” and “em” on the z-axis). Fig. 2b, data from two traits in females from three closely related populations of painted forest toadlets indicates the different axes of phenotype distance among these three closely related populations.

Although no formal conclusions about the relative significance of sexual selection and ecological adaptation in the process of population divergence can be drawn from Table 3 (as these require phylogenetic correction and time-since-divergence analyses), our comparisons strongly indicate greater distances between sexual traits compared to ecological traits, leading to the inference that sexual traits are more divergent in closely related taxa compared to those traits related to ecological adaptation. This conclusion is supported by two aspects of the results shown in Table 3. First, the point estimates of  $\Delta p_{y,zj}$  are greater (in magnitude) for the sexual trait than the ecological trait in 9 of 10 cases. Secondly, the 95% confidence intervals around  $\Delta p_{y,zj}$  do not include zero for any of the sexual traits, yet they do include zero for six of the 10 ecological traits. Moreover, Table 4

indicates that sexual trait dimorphism may generally be greater than ecological trait dimorphism, where the function of each phenotypic trait has been addressed through empirical field study. In 7 of 10 comparisons (using results within each population in Table 4), traits with known sexual signaling function are more dimorphic compared to traits related to ecological adaptation.

## 3 Discussion

Testing predictions of hypotheses about the role of sexual selection in speciation - and many other investigations related to trait divergence - requires researchers to compare the relative degree of inter-population divergence for very different types of traits (e.g. size and color). Here, we have emphasized that commonly used parametric distance metrics, such as Hedge’s  $g$  ( $g_{y,zj}$ ), have several drawbacks, which limit their usefulness in such studies. First, the definitions of many of these metrics assume that the trait distributions in divergent populations have equal variances (reviewed above). If the variances are unequal (which will not always be known or apparent with empirical data), the expected value obtained from equation (1) depends on differences in sample size (Fig. 1). Second, the numerical value of  $g_{y,zj}$  depends on the scale of measurement, and this metric will be affected if the data are subjected to a nonlinear transformation. This makes it difficult to compare the degree of divergence of different traits that may have been measured in very different ways (i.e., the problem of comparing “apples with oranges”).

With these problems in mind, we developed a novel, non-parametric distance measure,  $\Delta p$ , which does not depend on equality of variances, is independent of the scale of measurement (because it is non-parametric), and facilitates comparisons of several traits across several populations.  $\Delta p$  is based on comparing the location of population medians in the joint (trait-wise) cumulative distribution function (CDF) across all populations. Viewed differently,  $\Delta p$  compares population medians after transforming the data into percentiles of the joint cdf (this view of percentiles as an alternative scale of measurement is illustrated in all panels of Fig. 2). The percentile scale serves as a common frame of reference for all comparisons involving a given trait. In addition, percentiles provide a natural normalization (since they always range from 0 to 100), and they are independent of the original scale of measurement (because they only depend on the ranking of the raw data). These properties, in turn, allow for meaningful comparisons of divergence

measures for different traits. In sum, measuring divergence at the percentile scale makes it possible to really compare “apples to apples”.

$\Delta p$  may also be interpreted as a measure of overlap between two distributions (see Huberty and Lowman, 2000). This is most clearly seen in the case of two populations and a single trait. If we assume, for simplicity, that both sample distributions are symmetric, then the maximal possible value of  $\Delta p_{y,zj}$  is 50 (because the median of the smaller distribution is at least at the 25 percentile of the joint CDF, and the median of the larger distribution is at most at the 75 percentile). The difference between the actual value of  $\Delta p_{y,zj}$  and the maximal value (50) is determined by how much the lower tail of the larger distribution overlaps with the median of the smaller distribution, and vice versa.

### 3.1 Empirical comparisons

Recent hypotheses about speciation propose that sexual signal divergence is accompanied by ecological trait divergence, predicting that sexual selection plays a role in speciation – in cases with and without gene flow – when ecological contexts differ (e.g., van Doorn et al., 2009). According to this model, sexual trait divergence in closely related populations should coincide with ecological trait divergence, but this is not the case in the various systems explored to demonstrate the utility of  $\Delta p$  (Table 3). Table 3 presents data on the divergence shaped predominantly by sexual or natural selection. Although not a formal quantitative comparison in which phylogenetic relationships or a metric of time since divergence would need to be accounted for, a striking pattern *when comparing closely related species only* is that sexual signals are more strongly divergent than ecological traits among disparate taxonomic groups. Moreover, the values of  $\Delta p_{y,zj}$  are estimated on the same scale although these various acoustic signals, color variation, and morphological traits are measured in fundamentally different units. Thus, although in most cases  $g_{y,zj}$  and  $\Delta p_{y,zj}$  provide similar information about which traits are more divergent,  $\Delta p_{y,zj}$  provides the advantage that ecological and sexual trait differentiation are *directly* comparable. An interesting exception is the wolf spider case, which suggests that the ecological trait is slightly more divergent compared to the sexual trait.

For those taxa in which sexual signals are present in both males and females, we derived dimorphism estimates using  $\Delta p_{y,zj}$  to compute the differences in eco-

logical and sexual traits between males and females. Similar to the case in Table 3, Table 4 is not a formal analysis of whether sexual traits are more dimorphic than ecological traits, though among the five taxa examined, support for greater dimorphism in sexual signals is evident. An interesting exception, again, is the wolf spiders which suggest that leg length (a putative sexual trait in these species) is either hardly dimorphic (*S. bilineata*) or very dimorphic (*S. crassipalpata*) and that in both taxa, the ecological trait (cephalothorax width) is equally dimorphic but the direction of dimorphism differs (in *S. crassipalpata* females are larger than males). It is important to note that *S. bilineata* males develop brushes upon their tibial forelegs upon maturation – a secondary sexual trait that makes them distinctly dimorphic (Stratton 2005), potentially relieving foreleg length from sexual selection in this species. Additionally, due to the potential for sexual cannibalism in spiders, selection from ecological selection versus sexual selection is often intertwined, making predictions less apparent. Whereas the data from Table 4 are not conclusive evidence to support the use of sexual dimorphism as a proxy of sexual selection on phenotypic traits (e.g., Kraaijeveld et al., 2011), they do indicate that – in traits of known function – sexual traits may tend to be more dimorphic compared to those underlain predominantly by natural selection in the study systems described in Table 4.

Finally, as illustrated in Fig. 2, whereas an overall metric of distance can be obtained across multiple traits from multiple populations, the advantage of  $\Delta p_{y,z\bullet}$  is that the effect of one trait on overall distance among taxa can be quantified. For example, in Hume’s Warblers (Fig. 2d), it is clear that the sexual signal wing bar size rather than tarsus length is a major contributor to overall phenotype distance between these closely related taxa (compare females, “if” and “kf”, on the two axes; compare males, “im” and “km” on the two axes). Fig. 2d also shows that sexual dimorphism within populations is at least as pronounced as phenotypic divergence (within a sex) among populations.

### 3.2 Caveats and cautions in using $\Delta p_{y,zj}$ and $\Delta p_{y,z\bullet}$

One important consideration to keep in mind with  $\Delta p_{y,zj}$  is that while the ability to use more than two populations simultaneously is a strength of this distance measure, the magnitude of  $\Delta p_{y,zj}$  will change if a new population is added in the construction of  $p_j(u)$ . For example, suppose that  $p_j(u)$  is constructed for two

populations ( $y$  and  $z$ ) and  $\Delta p_{y,zj}$  is calculated. Now suppose that observations from trait  $j$  in a third population ( $w$ ) are added, and  $p_j(u)$  is recalculated to reflect the observations on all three populations.  $\Delta p_{y,zj}$  may now be reduced in magnitude if population  $w$  had more extreme trait values than the other populations; alternatively, if  $w$  was intermediate between  $y$  and  $z$ , then  $\Delta p_{y,zj}$  would be increased in magnitude. This property of  $\Delta p_{y,zj}$  is a direct consequence of the fact that percentiles given by  $p_j(u)$  are always bounded on the interval  $[0,100]$ , regardless of how many populations are being considered. The important consideration here is that if one wishes to compare the magnitudes of different  $\Delta p_{y,zj}$  values that were calculated *independently* from one another—for example, as might be done in a meta-analysis—then it is important that the calculations (1) involve the same set of populations (or at least, comparable sets of populations, e.g. 2 sympatric and one allopatric population) and (2) do not involve “saturation” of the metric (see below). In order to facilitate ease of conducting meta-analyses, we suggest it might be useful for any researcher reporting  $\Delta p_{y,zj}$  to report the pairwise distances (calculated from just two populations) along with the distances calculated for  $>2$  populations. However, whenever possible, this issue should be avoided by using one of the strengths that  $\Delta p_{y,zj}$  offers: all the populations should be put into the same analysis (rather than calculating  $\Delta p_{y,zj}$  values independently in different analyses). When all populations are compared in a single analysis, all comparisons of  $\Delta p_{y,zj}$  values will truly be “apples to apples”. The more general point here is that one of our main motivations for developing  $\Delta p_{y,z\bullet}$  was the need for a way to fairly compare distances among arbitrary numbers of populations and traits *simultaneously* and all on the same scale. When analyses are performed that way (i.e., one analysis using all appropriate data simultaneously), comparisons of magnitudes of  $\Delta p_{y,z\bullet}$  values (and  $\Delta p_{y,zj}$  values calculated as part of  $\Delta p_{y,z\bullet}$ ) will be valid.

A second consideration is that as differences between groups being compared become large,  $\Delta p_{y,zj}$  will eventually “saturate.” For example, in a pairwise comparison of body mass of hummingbirds and cheetahs,

$\Delta p_{y,zj}$  will be at its expected maximum magnitude (approximately 50 for a two-population analysis: see note below). The same would be true of an independent pairwise comparison of body mass between hummingbirds and elephants. The solution here is once again to use the features that  $\Delta p_{y,zj}$  offers: the data on hummingbirds, cheetahs, and elephants should all be included in a single analysis, in which case  $\Delta p_{y,zj}$  will resolve distances between populations appropriately. Another point to consider here is that we expect that most applications of  $\Delta p_{y,zj}$  and  $\Delta p_{y,z\bullet}$  will involve closely related groups, in which case “saturation” of the metric is unlikely to diminish its utility. For example, across the wide range of taxa and types of traits shown in Tables 3 and 4, it would have been problematic if a comparison involved two values of  $\Delta p_{y,zj}$  that were both near saturation values. There was only one case in which this occurred: in Table 3, in the row for field crickets in the genus *Gryllus*,  $|\Delta p_{y,zj}|$  was near 50 for *both* the sexual trait and the ecological trait. However, in this case, the narrow, non-overlapping confidence intervals around each estimate of  $\Delta p_{y,zj}$  still permit a meaningful comparison showing that it is highly likely that the sexual trait is more phenotypically divergent than the ecological trait. In other cases where comparisons did not produce unequivocal differences, it is sample size (and associated wide confidence intervals) rather than saturation that is the limiting factor. We note that when only two populations are considered, the theoretical expected maximum value of  $\Delta p_{y,zj}$  with infinite sample sizes is 50. However, values slightly larger than this can be realized for real data sets and for the confidence intervals around  $\Delta p_{y,zj}$ , especially when sample sizes are small, as is seen occasionally in Table 3. This is because—with finite sample sizes—there is no reason that the medians of two non-overlapping trait distributions must fall exactly at the 25<sup>th</sup> and 75<sup>th</sup> percentiles in the joint CDF,  $p_j(u)$ . In particular, deviations can occur when the medians coincide exactly with one or more trait values. We also note that Table 4 has values larger than 50 for a different reason: there are 4 populations included simultaneously in the calculations of  $\Delta p_{y,zj}$ .

A third and practical consideration is that calculating  $\Delta p_{y,zj}$ , and thus  $\Delta p_{y,z\bullet}$ , requires raw data. Calculation

of  $g_{y,zj}$  requires only having means, standard deviations, and sample sizes, which are often easy to obtain from published works; by contrast,  $\Delta p_{y,zj}$  utilizes a distribution of data. While the latter contributes to its desirable properties, it also means that one cannot calculate  $\Delta p_{y,zj}$  without access to original data sets (or at least, a random subsample of data from an original data set). In the current academic climate of free, electronic access to original data sets—which indeed, is now required upon publication by a number of journals in ecology and evolutionary biology (Fairburn, 2011)—we expect that the need for original data will be much less of an impediment than it might have been even just a decade ago. In light of this transition and because of issues related to the number of populations in a study and saturation, we recommend publishing both  $\Delta p_{y,zj}$  and  $g_{y,zj}$  side by side in studies related to phenotype distances, noting the advantages and disadvantages associated with each effect size metric.

**Acknowledgements** We thank Matthew Arnegard, Carlos Botero, Tamra Mendelson, Rafael Rodriquéz and Sander van Doorn for excellent discussions about the need for a new phenotypic distance metric and Maria Servedio for the invitation and encouragement to formalize our ideas. This research was supported as part of the Sexual Selection and Speciation working group by the National Evolutionary Synthesis Center (NESCent), NSF #EF-0905606. RJS and SMF were supported by the University of Colorado and National Science Foundation grant IOS-0717421 to RJS. MK was supported by a grant from the Vienna Science and Technology Fund (WWTF) to the Mathematics and Biosciences Group at the University of Vienna. EAH thanks Mitch Bern for use of his Master's thesis data and was supported by the National Science Foundation grant IOS - 0643179. DEI and DPLT were supported by the Natural Sciences and Engineering Research Council of Canada (Discovery Grants 311931-2005 and 311931-2010 to DEI, CGS-D to DPLT). NS and JAT were supported by the Royal Society, British Ecological Society and John Fell Fund (Oxford University). ES supported by NSF-DDIG, the American Ornithologists Union, the University of Chicago, and the American Philosophical Society Lewis and Clark award. JACU was funded by National Science Foundation grant IOS 0306175.

## References

- Arnegard ME, McIntyre PB, Harmon LB, Zelditch ML, Crampton WGR et al., 2010. Sexual signal evolution outpaces ecological divergence during electric fish species radiation. *Am. Nat.* 176: 335–356.
- Boul KE, Funk WC, Darst CR, Cannatella DC, Ryan MJ, 2007. Sexual Selection Drives Speciation in an Amazonian Frog. *Proc. Roy. Soc. B* 274: 399–406.
- Charlesworth B, 1998. Measures of divergence between populations and the effect of forces that reduce variability. *Mol. Biol. Evol.* 15: 538–543.
- Charlesworth B, Charlesworth D, Barton NH, 2003. The effects of genetic and geographic structure on neutral variation. *Annu. Rev. Ecol. Evol. Syst.* 34: 99–125.
- Cohen J, 1969. *Statistical power analysis for the behavioral sciences*. 1<sup>st</sup> edn. New York: Academic Press.
- Endler J, Houde AE, 1995. Geographic variation in female preferences for male traits in *Poecilia reticulata*. *Evolution* 49: 456–468.
- Excoffier L, 2001. Analysis of population subdivision. In: Balding DJ, Bishop M, Cannings C ed. *Handbook of Statistical Genetics*. New York: John Wiley, 271–307.
- Fairburn DJ, 2011. The advent of mandatory data archiving. *Evolution* 65: 1–2.
- Freeman-Gallant CR, Taff CC, Morin DF, Dunn PO, Whittingham LA et al., 2009. Sexual selection, multiple ornaments, and age- and condition-dependent signaling in the common yellowthroat. *Evolution* 64: 1007–1017.
- Funk WC, Angulo A, Caldwell JP, Ryan MJ, Cannatella DC, 2008. Comparison of morphology and calls of two cryptic species of *Physalaemus* (Anura: Leiuperidae). *Herpetologica* 64: 290–304.
- Gray DA, Cade WH, 2000. Sexual Selection and Speciation. *Proc. Natl. Acad. Sci. USA* 97: 14449–14454.
- Gray DA, Walker TJ, Conley BE, Cade WH, 2001. A morphological means of distinguishing females of the cryptic field cricket species *Gryllus rubens* and *G. texensis* (Orthoptera: Gryllidae). *Florida Entomologist* 84: 314–315.
- Grissom RJ, Kim JJ, 2001. Review of assumptions and problems in the appropriate conceptualization of effect size. *Psych. Methods* 6: 135–146.
- Hedges LV, 1981. Distribution theory for Glass's estimator of effect size and related estimators. *J. Educ. Stat.* 6: 107–128.
- Hedrick PW, 2005. A standardized genetic differentiation measure. *Evolution* 59: 1633–1638.
- Huberty CJ, Lowman LL, 2000. Group overlap as a basis for effect size. *Educ. Psych. Measurement* 60: 543–563.
- Irwin DE, Bensch S, Price TD, 2001. Speciation in a ring. *Nature* 409: 333–337.
- Irwin DE, Thimman MP, Irwin JH, 2008. Call divergence is correlated with geographic and genetic distance in greenish warblers *Phylloscopus trochiloides*: A strong role for stochasticity in signal evolution? *Journal of Evolutionary Biology* 21: 435–448.
- Izzo AS, Gray DA, 2004. Cricket song in sympatry: Examining reproductive character displacement and species specificity of song in *Gryllus rubens*. *Annals of the Entomological Society of America* 97: 831–837.
- Johnsen A, Andersson S, Fernandez JG, Kempnaers B, Pavel V et

- al., 2006. Molecular and phenotypic divergence in the blue-throat subspecies complex. *Molec. Ecol.* 15: 4033–4047.
- Kraaijeveld K, Femmie Kraaijeveld-Smit JL, Maan M, 2011. Sexual selection and speciation: The comparative evidence revisited. *Biol. Rev.* 86: 367–377.
- Lande R, 1981. Models of speciation by sexual selection on polygenic traits. *Proc. Natl. Acad. Sci USA* 78: 3721–3725.
- Maan ME, Seehausen O, 2011. Ecology, sexual selection and speciation. *Ecol. Letters*. 14: 591–602.
- Mahalanobis PC, 1936. On the generalised distance in statistics. *Proceedings of the National Institute of Sciences of India* 2: 49–55.
- Mayr E, 1947. Ecological factors in speciation. *Evolution* 1: 263–288.
- Mendelson TC, Shaw KL, 2005. Sexual behaviour: Rapid speciation in an arthropod. *Nature* 433: 375–376.
- Nagawa S, Cuthill IC, 2005. Effect size, confidence interval and statistical significance: A practical guide for biologists. *Biol. Rev.* 82: 591–605.
- Panhuis TM, Butlin R, Zuk M, Tregenza T, 2001. Sexual selection and speciation. *Trends Ecol. Evol.* 6: 364–371.
- Price TD, 1998. Sexual selection and natural selection in bird speciation. *Philosophical Transactions of the Royal Society of London B* 353: 251–260.
- Price TD, 2008. *Speciation in Birds*. Greenwood Village: Roberts and Company.
- Ritchie MG, 2007. Sexual selection and speciation. *Annual Review of Ecology, Evolution and Systematics* 38: 79–102.
- Rodríguez RL, Sullivan LE, Cocroft RB, 2004. Vibrational communication and reproductive isolation in the *Enchenopa binotata* species complex of treehoppers (Hemiptera: Membracidae). *Evolution* 58, 571–578.
- Safran RJ, McGraw KJ, 2004. Plumage coloration, not length or symmetry of tail-streamers, is a sexually selected trait in North American barn swallows. *Behavioral Ecology* 15: 455–461.
- Seehausen O, Van Alphen JM, 1999. Can sympatric speciation by disruptive sexual selection explain rapid evolution of cichlid diversity in Lake Victoria? *Ecology Letters* 2: 262–271.
- Seddon N, Merrill RM, Tobias JA, 2008. Sexually selected traits predict patterns of species richness in a diverse clade of suboscine birds. *American Naturalist* 171: 620–631.
- Schluter D, 1996. Adaptive radiation along genetic lines of least resistance. *Evolution* 50: 1766–1774.
- Slatkin M, 1987. Gene flow and the geographic structure of natural populations. *Science* 236: 787–792
- Sokal RR, Rohlf FJ, 1995. *Biometry: The Principles and Practice of Statistics in Biological Research*. 3<sup>rd</sup> edn. Freeman: San Francisco.
- Stein AC, Uy JAC, 2006. Plumage brightness predicts male mating success in the lekking golden-collared manakin. *Behavioral Ecology* 17: 41–47.
- Stratton GE, 2005. Evolution of ornamentation and courtship behavior in Schizocosa: Insights from a phylogeny based on morphology (Araneae, Lycosidea). *Journal of Arachnology* 33: 347–376
- Svensson EI, Eroukhmanoff F, Friberg M, 2006. Effects of natural and sexual selection on adaptive population divergence and premating isolation in a damselfly. *Evolution* 60: 1242–1253.
- Tobias JA, Seddon N, 2009. Signal design and perception in *Hypocnemis antbirds*: Evidence for convergent evolution via social selection. *Evolution* 63: 3169–3189.
- Toews DPL, Irwin DE, 2008. Cryptic speciation in a Holarctic passerine revealed by genetic and bioacoustic analyses. *Molecular Ecology* 17: 2691–2705.
- Uy JAC, Borgia G, 2000. Sexual selection drives rapid divergence in bowerbird display traits. *Evolution* 54: 273–278.
- Uy JAC, Moyle RG, Filardi CE, 2008. Plumage color and song differences mediate species recognition between incipient flycatcher species of the solomon islands. *Evolution* 63: 153–164
- van Doorn S, Edelaar P, Weissing FJ, 2009. On the origin of species by natural and sexual selection. *Science* 326: 1704–1707.
- West- Eberhard MJ, 1983. Sexual selection, social competition, and speciation. *Quarterly Review of Biology* 58: 155–183
- Wright S, 1943. Isolation by distance. *Genetics* 28: 114–128.
- Wright S, 1951. The genetical structure of populations. *Ann. Eugen.* 15: 323–354.
- Wright S, 1965. The interpretation of population structure by *F*-statistics with special regard to systems of mating. *Evolution* 19: 395–420.
- Wright S, 1973. Analysis of gene diversity in subdivided populations. *Proc. Natl. Acad. Sci.* 70: 3321–3323
- Wright S, 1978. *Evolution and the Genetics of Populations: Variability within and among Natural Populations*. Chicago: Chicago Univ. Press.

## Appendix 1 The following is a hypothetical example of how $\Delta p$ is calculated for a single trait

For ease of display, sample sizes here are small. Suppose we have the following observations on three populations:

Population	Observation values, sorted in ascending order														
A	0.48	2.30	3.00	4.13	5.64	5.69	6.08	6.72	8.67	10.54	11.07	12.16			
B	3.98	6.49	6.69	6.84	8.22	8.68	8.78	8.79	8.81	8.82	9.59	10.52	10.54	10.72	11.08
C	0.17	3.00	5.56	5.79	6.57	6.64	6.74	8.09	8.69	9.88	9.98	11.67	13.11	13.32	

Each of these observations accounts for some percent of the overall distribution of the data; for populations to contribute equally to this distribution (i.e., WITHOUT weighting by sample size), we need assign a weight to each observation that is inversely proportional to the sample size from the population of origin.

Population	Sample Size	Weight Assigned per observation
A	12	$(1/12) / 3 = 0.0278$
B	15	$(1/15) / 3 = 0.0222$
C	14	$(1/14) / 3 = 0.0238$

Note: the weight is divided by 3 because there are three populations

Using these weights, we can construct an empirical probability mass function (PMF) and cumulative distribution function (CDF); the CDF is used to assign

percentiles to the observed values, which in turn will be used to determine the percentile corresponding to a given median trait value.

Population	Observation Value	Weight (PMF)	CDF	Percentile (CDF * 100)	Population	Observation Value	Weight (PMF)	CDF	Percentile (CDF * 100)
C	0.17	0.0238	0.0238	2.38	A	8.67	0.0278	0.5516	55.16
A	0.48	0.0278	0.0516	5.16	B	8.68	0.0222	0.5738	57.38
A	2.30	0.0278	0.0794	7.94	C	8.69	0.0238	0.5976	59.76
A	3.00	0.0278	0.1071	10.71	B	8.78	0.0222	0.6198	61.98
C	3.00	0.0238	0.1310	13.10	B	8.79	0.0222	0.6421	64.21
B	3.98	0.0222	0.1532	15.32	B	8.81	0.0222	0.6643	66.43
A	4.13	0.0278	0.1810	18.10	B	8.82	0.0222	0.6865	68.65
C	5.56	0.0238	0.2048	20.48	B	9.59	0.0222	0.7087	70.87
A	5.64	0.0278	0.2325	23.25	C	9.88	0.0238	0.7325	73.25
A	5.69	0.0278	0.2603	26.03	C	9.98	0.0238	0.7563	75.63
C	5.79	0.0238	0.2841	28.41*	B	10.52	0.0222	0.7786	77.86
A	6.08	0.0278	0.3119	31.19	A	10.54	0.0278	0.8063	80.63
B	6.49	0.0222	0.3341	33.41	B	10.54	0.0222	0.8286	82.86
C	6.57	0.0238	0.3579	35.79	B	10.72	0.0222	0.8508	85.08
C	6.64	0.0238	0.3817	38.17	A	11.07	0.0278	0.8786	87.86
B	6.69	0.0222	0.4040	40.40	B	11.08	0.0222	0.9008	90.08
A	6.72	0.0278	0.4317	43.17	C	11.67	0.0238	0.9246	92.46
C	6.74	0.0238	0.4556	45.56	A	12.16	0.0278	0.9524	95.24
B	6.84	0.0222	0.4778	47.78	C	13.11	0.0238	0.9762	97.62
C	8.09	0.0238	0.5016	50.16	C	13.32	0.0238	1.0000	100.00
B	8.22	0.0222	0.5238	52.38					

two identical values

two identical values

\* observation values highlighted in red text are those equal to or immediately below each population's medians.

Next, we locate the medians of the trait values in the three populations and the associated percentiles (highlighted in red text)

Population	Median	Percentile
A	5.89	28.41
B	8.79	64.21
C	7.42	47.78

And from these we calculate  $\Delta p$  values pairwise:

Population 1	Population 2	$\Delta p$
A	B	-35.8
A	C	-19.37
B	C	16.43

Running our MATLAB code on the provided example data file ("DataUsedInExampleSpreadsheet.csv") will generate these same results

## Appendix 2 More numerical studies of the behavior of $\Delta p$

### I. A comparison of the performance of Hedge's $g$ and $\Delta p$ for non-normally distributed data.

For the results shown in Fig. S1, Table S1 gives the parameters used in constructing several "Scenarios" that were explored with simulated data. For all of these scenarios, pseudorandom data following an exponential distribution were generated from MATLAB's `random()` function. Methods otherwise follow those used to produce Fig. 2 in the main text. The patterns shown here (Fig. S1) very closely resemble those shown in the main text (Fig. 2): both  $g_{y,zj}$  and  $\Delta p_{y,zj}$  perform similarly with exponentially distributed data as they did with normally distributed data, and  $\Delta p_{y,zj}$  shows less sensitivity to differences in sample size and sample variance than does  $g_{y,zj}$ .

### II. A comparison of the performance of Hedge's $g$ and $\Delta p$ when sample sizes are small

For the results shown in Figure S2, Table S2 gives the parameters used in constructing several "Scenarios" that were explored with simulated data. For all of these scenarios, pseudorandom data following a normal dis-

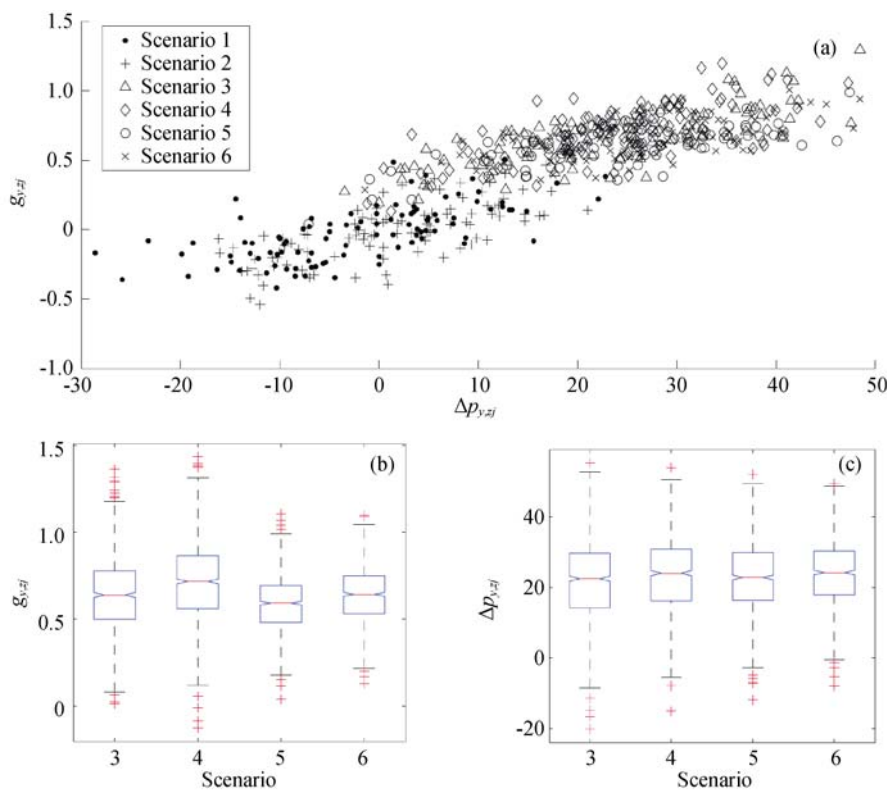
tribution were generated from MATLAB's `randn()` function. Methods otherwise follow those shown in Table 2

**Table S1 Parameters<sup>a</sup> used to generate pseudorandom data following exponential distributions which was used to generate results shown in Fig. S1**

Scenario	Population $y$		Population $z$	
	$n_{yj}$	$\mu_{yj}$	$n_{zj}$	$\mu_{zj}$
1: unequal sample sizes only	31	30	59	30
2: unequal sample sizes only	31	15	59	15
3: unequal means and variances <sup>b</sup> only	31	30	31	15
4: unequal means, variances, and sample sizes	31	30	59	15
5: as in 4, but with reversed sample sizes	59	30	31	15
6: as in 3, but with larger sample sizes	59	30	59	15

<sup>a</sup> The parameters in this table are not meant to reflect any specific population; rather, the purpose here is simply to compare the performance of  $g_{y,zj}$  and  $\Delta p_{y,zj}$  for non-normally distributed data.

<sup>b</sup> We note for clarity that no variance parameter is given in the table because an exponential distribution has just a single parameter (the mean); the expected variance of data from an exponential distribution is  $\mu_{ij}^2$ .



**Fig. S1 Comparisons between the behavior of  $g_{y,zj}$  and  $\Delta p_{y,zj}$  using simulated data following exponential distributions**

1000 data sets were generated with the parameters listed for each "Scenario" in Table S1. Scenarios 1 and 2 involve no true difference between populations; Scenarios 3-6 involve a true difference between the populations. Comparing (b) and (c), we see again (as in Figure 2 in the main text), that  $g_{y,zj}$  shows more sensitivity to variation in sample size and variance than does  $\Delta p_{y,zj}$ .



and Fig. 2 in the main text. The main differences between these examples and those shown in the main text (Table 2 and Figure 2 in the main text) are that here the sample sizes are smaller (Table S2), as are the calcu-

lated distances (compare Fig. S2b,c with Fig. 2b,c).

The patterns shown here (Fig. S2) again resemble those shown in the main text (Fig. 2).

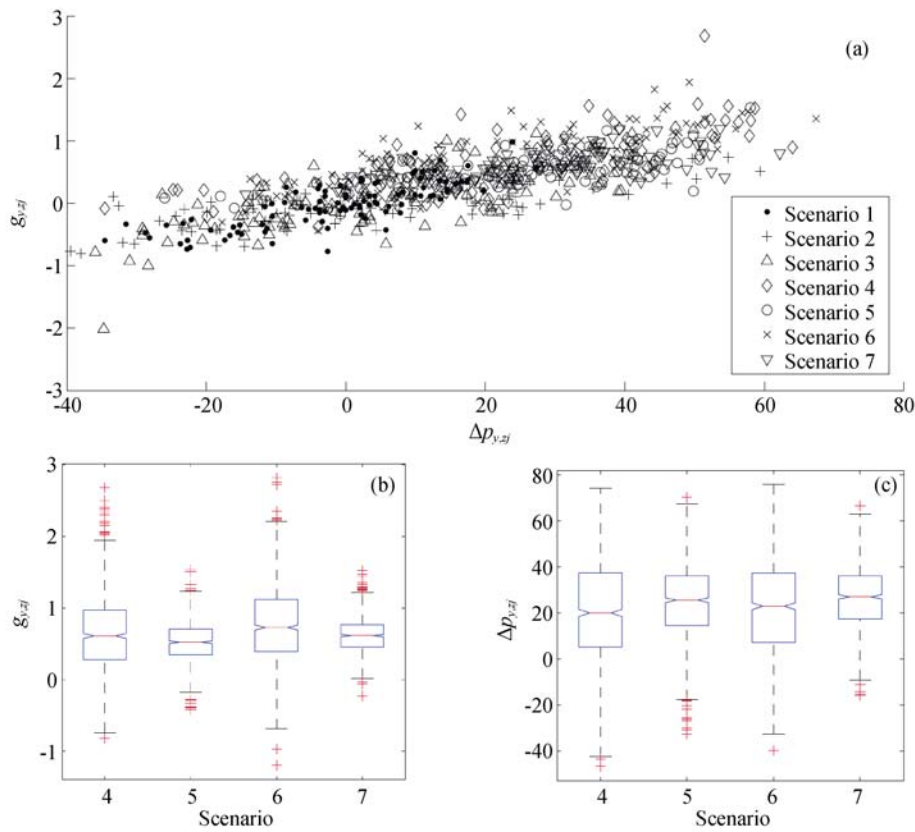
**Table S2** Parameters<sup>a</sup> used in simulated “Scenarios” used to compare  $g_{y,zj}$  and  $\Delta p_{y,zj}$  in Fig. S2

Scenario	Population y			Population z		
	$n_{yj}$	$\mu_{yj}^b$	$\sigma_{yj}^c$	$n_{zj}$	$\mu_{zj}^b$	$\sigma_{zj}^c$
1: unequal sample sizes only	9	0.40	0.50	36	0.40	0.50
2: unequal sample sizes only; means equivalent (but differ from Scenario 1)	9	-0.02	0.84	36	-0.02	0.84
3: unequal variances only	9	-0.02	0.84	9	-0.02	0.5
4: unequal means and variances only	9	0.40	0.50	9	-0.02	0.84
5: unequal means, variances, and sample sizes	9	0.40	0.50	36	-0.02	0.84
6: as in 5, but with reversed sample sizes	36	0.40	0.50	9	-0.02	0.84
7: as in 4, but with larger sample sizes	36	0.40	0.50	36	-0.02	0.84

<sup>a</sup> Sample sizes, means, and standard deviations were all inspired by a real data set on residual dorsal widths from two populations of female painted forest toadlets (the trait on the y-axis of Figure 1b in the main text; populations used were those labeled A and E in that Figure).

<sup>b</sup> Assumed true population mean used for generating pseudorandom data.

<sup>c</sup> Assumed true population standard deviation used for generating pseudorandom data.



**Fig. S2** Comparisons between the behavior of  $g_{y,zj}$  and  $\Delta p_{y,zj}$  using simulated data following normal distributions using the parameters given in Table S2

1000 data sets were generated with the parameters listed for each “Scenario” in Table S2. The interpretation of Scenarios and panels is the same as in Figure 2 in the main text. Comparing (b) and (c), we see again (as in Fig. 2 in the main text), that  $g_{y,zj}$  shows more sensitivity to variation in sample size and variance than does  $\Delta p_{y,zj}$ .

## Online Appendix 3 Field Methods for Unpublished Data in Tables 3 and 4

### *Hirundo rustica* / Rebecca J. Safran, Matthew R. Evans

Data were collected at two different field locations. Data from North America are from barn swallows breeding near Ithaca New York whereas data from the United Kingdom are from barn swallows breeding near Cornwall, southern England. In both study areas, individuals were captured at breeding sites using mist nets and individually marked for later identification. The following morphological measurements were collected from each individual: streamer length (outer rectrix; mm), wing length (flattened 9<sup>th</sup> primary; mm), tarsus length (mm), and body mass following methods outlined in Safran and McGraw 2004.

### Hume's Warbler/ Elizabeth Scordato

Hume's warbler *Phylloscopus humei* is a small, insectivorous passerine distributed throughout the Himalayas, central Asia, and western Siberia. Data were collected on wing-bar size (a sexually selected plumage trait; Marchetti 1998) and tarsus length from males and females at two different breeding populations: Keylong, Himachal Pradesh, India (2008, 32°37'54.05"N, 77°10'39.00"E) and Ala Archa Park, Kyrgyzstan (2009, 42°38'33.92"N, 74°29' 13.70"E). Adult birds were caught in mist nets when they arrived on the breeding grounds and ringed for future identification. Tarsus length for each bird was measured with brass calipers as the distance from the ankle joint to the last undivided scute at the base of the foot, to the nearest 0.1 mm. The 4<sup>th</sup> greater covert feather (counting in from distal to the body) was collected from each wing for later measurement of wing-bar size. Adults were sexed behaviorally (only males sing), or by the presence of a brood patch (only females incubate). After completing data collection, each feather was mounted on a microscope slide. Feathers were measured under a dissecting microscope with an illuminated stage at 10x magnification. Wing-bar size was measured as the length of the unmelanized portion of the feather rachis (mm). The average of the right and left feather measurements was used as the wing-bar size for that individual; this has been shown to correlate strongly with wing-bar size based on measurements of the 2<sup>nd</sup>, 4<sup>th</sup>, and 6<sup>th</sup> coverts combined (Price and Pavelka 1996). Measurements were randomized across the two years to minimize bias in measurements over time. Detailed methods in Scordato *et al* (in review).

### Tree Crickets / Laurel Symes

*Oecanthus forbesi* were collected in August 2010 at Bong Recreation Area (42.63451, 88.130378) in Kansasville, WI and Dawes Arboretum (39.98843, 82.41172) in Newark, OH. Males were located visually and acoustically and were hand-captured. Tooth number was determined by photographing the top right tegmina under a dissecting scope and counting the teeth in the file. The underwings were carefully detached from the preserved specimen, affixed to a transparency sheet, scanned on a high resolution flatbed scanner with transparency adaptor, and measured to the nearest thousandth of a centimeter in ImageJ (available at <http://rsbweb.nih.gov/ij/>).

### *Gryllus* Field Crickets / David Gray

Pulse rates of male songs were calculated from the pulse periods measured from digitized recordings of males' songs from throughout the species' ranges (collection localities are detailed in Gray and Cade, 2000; Izzo and Gray, 2004; Gray *et al.*, 2001; Gray, 2011). Pulse rates were corrected for temperature based on empirically determined regression analyses of pulse rates on temperature for these species. Females' ovipositor lengths and pronotal widths were measured to the nearest 0.1 mm; data presented here use the ratio of ovipositor length/pronotal width.

### *Engystomops petersi* Painted Forest Toadlet / W. Chris Funk

Field Methods. A total of 26 *Engystomops petersi* clade A from two sites in Ecuador; 33 *E. petersi* clade D from three sites in Ecuador and Peru; and 87 *E. freibergi* from 7 sites in Peru and Brazil were analyzed. Recent systematic analyses based on mitochondrial and nuclear DNA sequences, morphology, and calls reveal that clades A and D are distinct species (Funk *et al.*, in review).

Morphological Data. Three morphometric variables were measured from fixed specimens that have been shown to vary among *Engystomops* species (Funk *et al.*, 2008): snout-vent length (SVL), dorsum width, and tibia length. Measurements were taken using digital calipers accurate to the nearest 0.01 mm following Funk *et al.* (2008). The residuals of the linear regressions of dorsum width vs. SVL and tibia length vs. SVL were used in the analyses.

Call Data. Call analyses focused on call dominant frequency, as previous work has shown that this parameter is likely important in female preferences (Boul

et al., 2007; Guerra and Ron, 2008). Call recordings had a sampling rate of 44 kHz and were analyzed with program Raven 1.2 (Charif et al., 2004). Fast Fourier transformation (FFT) size was 2048 and the spectral analysis had a frequency resolution of 21.5 Hz.

#### **Antbirds / Nathalie Seddon, Joseph A. Tobias**

*Hypocnemis* antbirds are small (11–12 cm, 10–14 g) passerine birds of the understory and borders of rainforest in the Amazon basin and foothills of the Andes (Zimmer and Isler, 2003). From 2004 to 2009 (Sep–Dec), we studied *H. peruviana peruviana* and *H. subflava collinsi* within their region of overlap in SW Amazonia, at the Centro de Investigación y Conservación de Río Los Amigos (CICRA; 12°34'07"S, 70°05'57"W), Depto. Madre de Dios, Peru (Tobias and Seddon, 2009; Tobias et al., 2011). A total of 20 males and 24 females of *H. peruviana*, and 27 males and 20 females of *H. subflava* were tape-lured into a 12 x 4 m mistnets, and a range of standard biometric data were collected. These included culmen (bill) length measured using dial callipers (to the nearest 0.01 mm) from the anterior end of the nares. We also collected three feathers from the center of the chest and flank plumage patches, and quantified spectral reflectance following standard methodology described in detail in Tobias and Seddon (2009). For the purposes of this study, we used the reflectance data to quantify variation in the intensity of white in the chest feathers (chest chroma), calculating chest chroma as the proportion of total reflectance occurring between 400 nm and 490 nm, i.e. R400-490/R300-700. Meanwhile, songs were recorded from a total of 17 male and 9 female *H. peruviana*, and 21 male and 5 female *H. subflava*. The methods used to record songs, generate spectrograms and measure song pace are given in detail in Tobias and Seddon (2009).

#### **Manacus / J. Albert C. Uy**

Manakin leks were sampled and monitored in Gamboa, Republic of Panama (*M. vitellinus*) from Feb 15 through April 23, 2004 (Stein and Uy, 2006) and La Selva, Costa Rica (*M. candei*) from April 30 to May 18, 2001. Individuals were caught at or near the lek sites with mist nets, and were fitted with uniquely numbered aluminum and colored leg bands to facilitate identification. At the time of capture we took standard measurements of tarsus, wing chord, and tail lengths (mm), bill dimensions (mm) and body mass (g).

An Ocean Optics USB2000 spectroradiometer and a Xenon flash light source (Ocean Optics PX-2) was used to obtain reflectance scans of the bird's plumage (Uy

and Endler, 2004). The tip of a micron fiber optic probe was housed in a hollow, black anodized aluminum sheath with an angled tip that contacted the bird's plumage. This technique ensured that (1) the Xenon flash was the only light source for our scans, (2) the distance between the probe and color patch was standardized at 1 cm, and (3) the angle of measure was standardized at 45 (this angle reduces specular glare). To allow for comparison across different measures, a spectrally flat 97% reflecting spectralon white standard (Labsphere, North Sutton, NH) was used in addition to a dark current reading to standardize each scan. Scans were taken from an 3 mm diameter circle at 0.40 nm intervals across 300–700 nm. Reflectance scans were taken at the time of capture. Scanning was focused on the colorful golden patches of *M. vitellinus* and bright white patches of *M. candei* because these patches are expanded and most visible to females during courtships (Uy and Endler, 2004). Spectral measures of the same patches for individuals with brown plumage (i.e., immature individuals or females) were also taken. Three consecutive measurements were taken for each of these patches by lifting the probe and replacing it back onto the patch. We averaged the three readings taken at each color patch to get a mean estimate for each male's beard and collar.

Brightness or intensity of the plumage patches was calculated by summing its reflectance from 300 to 700 nm, the approximate visible spectrum of most avian species (Hart, 2001). This measure assays the total intensity of light reaching the eye from the color patch. Chroma, or color purity of a signal, was estimated by calculating the Euclidian distance between a color patch and the origin of a tetrahedron space with each vertex representing one of the four avian cones. The perception model developed by Chittka (1992) was implemented in Avicol ver. 5.0 (Gomez, 2011). Ambient light used for the perception model came from irradiance measures of ambient light at *Manacus* leks.

#### **Wolf Spiders / Eileen Hebets**

*Schizocosa crassipalpa* and *S. bilineata* are sister taxa with partially overlapping distributions across Eastern North America (Comstock, 1940; Kaston, 1948; Dondale and Redner, 1990; Sierwald et al., 2005; Vaccaro et al., 2010). *Schizocosa crassipalpa* males lack foreleg ornamentation, yet produce an asymmetric tap of their forelegs during courtship (Emerton and Dondale, 1978; Stratton, 2005; Vaccaro et al., 2010). *Schizocosa bilineata* males develop black pigmentation and brushes

on their foreleg tibiae upon maturation and their courtship display incorporates an incremental raising and lowering of one or both forelegs, culminating in a quick foreleg tap on the substrate (Emerton and Dondale, 1978; Stratton, 2005; Vaccaro et al., 2010). Given the role of the forelegs in the courtship display of both of these species, foreleg length is hypothesized to be under sexual selection.

**Collection and Maintenance** – Subadult *S. crassipalpata* ( $n = 103$ ; 1–2 molts from maturity) were collected in March 2009 from Bath Nature Preserve, Summit County, OH ( $41^{\circ} 10' 35.7414''$ ,  $81^{\circ} 38' 52.7928''$ W). Subadult *S. bilineata* ( $n = 173$ ; 2–3 molts from maturation) were collected in March 2010 from the Ohio State University campus at Newark, Licking County, OH ( $40^{\circ} 4' 30.0144''$ N,  $82^{\circ} 26' 32.8272''$ W). Spiders were housed individually, maintained on a 12:12 light:dark cycle, and provided with a constant source of water.

**Trait Quantification** Cephalothorax width (CW) was quantified as an over-all size measure – our trait hypothesized to reflect mostly natural selection. Mature individuals were sacrificed and photographed using a Leica light microscope (Leica Microsystems, Bannockburn, IL, USA) fitted with a SPOT Flex 15.2 64 Mp camera (SPOT Imaging Solutions, Sterling Heights, MI, USA). Using the photographs, the CW was measured at its widest point three separate times using Image-Pro Discovery software (Media Cybernetics, Bethesda, MD, USA); and the mean of the measurements was calculated.

Leg length was quantified as our trait hypothesized to reflect mostly sexual selection. We hypothesize that leg length may be sexually selected as male forelegs are used during leg waving courtship displays. Forelegs from males and females were removed and photographed using the same microscope set-up as above. The exterior lateral side of each leg was photographed and from the digital images, measurements were made of the entire leg length. Leg length measurements consisted of summing the lengths of the tarsus, metatarsus, patella, and femur.

All data incorporated in these analyses were generated as part of a Master's Thesis by Mitch Bern at the University of Nebraska, Lincoln, USA.

## References

Boul KE, Funk WC, Darst CR, Cannatella DC, Ryan MJ, 2007 Sexual selection drives speciation in an Amazonian frog. Proc. R. Soc. B 274: 399–406.

Charif RA, Clark CW, Fristrup KM, 2004. Raven 1.2 User's

Manual. Ithaca, New York: Cornell Laboratory of Ornithology.

Chittka L, 1992. The color hexagon: A chromaticity diagram based on photoreceptor excitations as a generalized representation of color opponency. J. Comp. Physiol. A. 170: 533–543.

Comstock JH, 1940. The Spider Book: A Manual for the Study of the Spiders and Their Near Relatives, the Scorpions, Pseudoscorpions, Whipscorpions, Harvestmen and Other Members of the Class Arachnida, Found in America North of Mexico, with Analytical Keys for Their Classification and Popular Accounts of Their Habits. New York: Doubleday, Doran & Company, Inc..

Dondale CD, Redner JH, 1978. Revision of the Nearctic wolf spider genus *Schizocosa* (Aranei-da: Lycosidae). Canadian Entomologist 110:143–181

Dondale CD, Redner JH, 1990. The Insects and Arachnids of Canada, Part 17. The Wolf Spiders, Nurseryweb Spiders, and Lynx Spiders of Canada and Alaska. Ottawa: Biosystematics Research Centre.

Funk WC, Angulo A, Caldwell JP, Ryan MJ, Cannatella DC, 2008. Comparison of morphology and calls of two cryptic species of *Physalaemus* (Anura: Leiuperidae). Herpetologica 64: 290–304.

Gomez D, 2011. A program to analyse spectrometric data. Free executable available at <http://sites.google.com/site/avicolprogram/>.

Gray DA, 2011. Speciation, divergence, and the origin of *Gryllus rubens*: Behavior, morphology, and molecules. Insects 2: 195–209.

Gray DA, Cade WH, 2000. Sexual selection and speciation in field crickets. Proceedings of the National Academy of Sciences, USA 97:14449–14454

Gray DA, Walker TJ, Conley BE, Cade WH 2001. A morphological means of distinguishing females of the cryptic field cricket species *Gryllus rubens* and *G. texensis* (Orthoptera: Gryllidae). Florida Entomologist 84: 314–315

Guerra MA, Ron SR, 2008. Mate choice and courtship signal differentiation promotes speciation in an Amazonian frog. Behav. Ecol. 19: 1128–1135.

Hart NS, 2001. The visual ecology of avian photoreceptors. Progress in Retinal and Eye Research 20: 675–703.

Izzo AS, Gray DA, 2004. Cricket song in sympatry: Examining reproductive character displacement and species specificity of song in *Gryllus rubens*. Annals of the Entomological Society of America 97: 831–837

Kaston BJ, 1948. Spiders of Connecticut. State Geological and Natural History Survey of Connecticut Bulletin 70:1–874

Marchetti K, 1998. The evolution of multiple male traits in the yellow-browed leaf warbler. Anim. Behav. 55: 361–376.

Price TD, Pavelka M, 1996. Evolution of a colour pattern: History, development, and selection. J. Evolution. Biol. 9: 451–470.

Price TD, Pavelka M, 1996. Evolution of a colour pattern: History, development, and selection. J. Evolution. Biol. 9: 451–470.

Safran RJ, McGraw KJ. 2004. Plumage coloration, not length or symmetry of tail-streamers, is a sexually selected trait in North American barn swallows. Behavioural Ecology 15: 455–461.

Scordato ES, Bontrager AL, Price TD, 2012. Cross-generational effects of climate change on expression of a sexually selected trait. Current Biology 22: 78–82

Sierwald P, Draney ML, Prentice T, Pascoe F, Sandlin N et al.,

2005. Spider Species of the Great Lakes States. Proceedings of the Indiana Academy of Science 114:111–206
- Stein AC, Uy JAC 2006 Plumage brightness predicts male mating success in the lekking golden-collared manakin. Behavioral Ecology 17: 41–47.
- Stratton GE, 2005. Evolution of ornamentation and courtship behavior in Schizocosa: Insights from a phylogeny based on morphology (Araneae, Lycosidae). Journal of Arachnology 33: 347–376
- Tobias JA, Gamarra-Toledo V, Garcia-Olaechea D, Pulgarin PC, Seddon N. 2011. Year-round resource defence and the evolution of male and female song in suboscine birds: Social armaments are mutual ornaments. Journal of Evolutionary Biology 24: 2118–2138.
- Tobias JA, Seddon N, 2009 Signal design and perception in Hypocnemis antbirds: Evidence for convergent evolution via social selection. Evolution 63: 3169–3189.
- Uy JAC, Endler JA, 2004. Modification of the visual background increases the conspicuousness of golden-collared manakin displays. Behavioral Ecology 15: 1003–1010.
- Vaccaro R, Uetz GW, Roberts AJ, 2010. Courtship and mating behavior of the wolf spider Schizocosa bilineata (Araneae: Lycosidae). Journal of Arachnology 38:452–459
- Zimmer KJ, Isler ML, 2003. Family Thamnophilidae (typical antbirds). In del Hoyo J, Elliott A, Christie D ed. Handbook of Birds of the World, vol. 8. Barcelona: Lynx Edicions, 448–681.

## Online Appendix 4: Confidence intervals for $\Delta p$

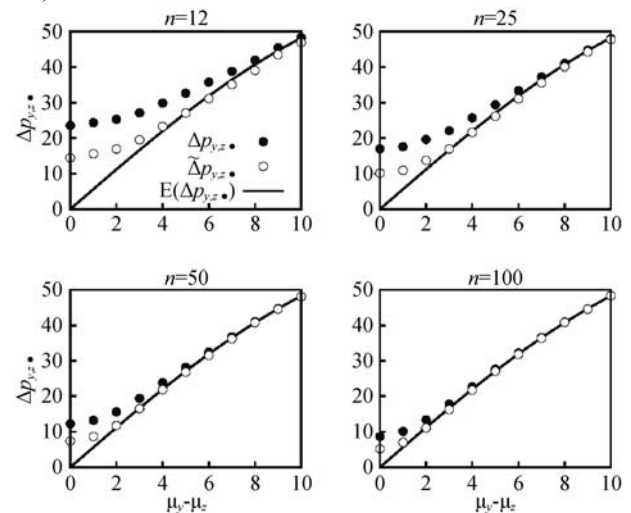
In this Appendix, we describe our methods for calculating confidence intervals (CIs) for the univariate and multivariate versions of  $\Delta p$  ( $\Delta p_{y,zj}$  and  $\Delta p_{y,z*}$ , respectively; see Table 1 in the main text for definitions of all notation used here). We used non-parametric bootstrapping, and largely followed previous work on CIs for parametric effect size measures such as Cohen's  $d$  (Kelley, 2005; Algina et al., 2006) and the Mahalanobis distance (Reiser, 2001; Hess et al., 2007).

The bootstrap is a flexible and widely used tool for statistical inference, including CI construction (Efron, 1987; Carpenter and Bithell, 2000). The basic idea is to estimate the accuracy of an estimator by resampling with replacement from the original data, calculating the estimator from this resample and repeating this process many times. In the context of effect size measures, resampling has to be done independently for each study population. The resamples yield a bootstrap distribution, from which the CI can be deduced. The MATLAB code (see <http://sourceforge.net/projects/deltap/files/>) provided with this paper offers three methods: the percentile method, which uses the 2.5% and 97.5% percentiles of the bootstrap distribution as the lower and upper confidence limits, and two refinements of the former, the bias-corrected (BC) and bias-corrected and accelerated (BCa) methods (Efron 1987). Our code for the latter is adapted from Kelley (2005).

Methods for CI construction can be tested with simulated data that have been drawn from distributions with known parameters. While a full analysis is beyond the scope of this paper, in the following, we briefly summarize the results of preliminary simulations using both normally and non-normally distributed data with equal or unequal variances between populations. The primary measure for CI performance is coverage probability, that is, the fraction of CIs (calculated from independent data sets) that bracket the true parameter value. The nominal coverage for a CI is equal to its confidence level (e.g., 95%), but the actual coverage may differ (e.g., due to violations of model assumptions). Our simulations indicate that, for the univariate  $\Delta p_{y,zj}$ , all three methods of CI construction behave very similarly and, in all cases tested, produce coverage very close to the nominal value (usually between 92.5% and 97.5% for 95% CIs, which is considered good performance, see Bradley 1978; coverage for the percentile method tends to be slightly above 95%, while coverage

for the other two methods tends to be slightly below 95%).

The situation is more complicated for the multivariate effect size  $\Delta p_{y,z*}$  (recall, from the main text, that the latter is a measure of the *overall* distance between populations when data on *two or more traits are being used simultaneously*). The reason is that  $\Delta p_{y,z*}$  (like any multivariate distance metric) cannot be negative. Hence, if the true distance between two populations is small (or zero), sampling variance will almost inevitably lead to an overestimate. Indeed, simulations show that  $\Delta p_{y,z*}$  is positively biased for small differences in population means (Fig. S3). Furthermore, the same problem applies to the bootstrap distribution and, hence, the CI estimates: When the populations being compared have very similar means, all three types of CIs have extremely low coverage (e.g., when the populations are identical, the percentile CI almost never contains the true value  $\Delta p_{y,z*} = 0$ ).



**Fig. S3 Bias and bias-correction in the multivariate estimator**

$\Delta p_{y,z*}$ . Closed symbols show the mean value of  $\Delta p_{y,z*}$  from 5000 simulations, whereas open symbols show the mean of the corrected estimator  $\Delta \tilde{p}_{y,z*}$  (standard errors are too small to be visible). The solid line is the expected value of  $\Delta p_{y,z*}$  in two infinite populations with two identically distributed traits each. The mean value for both traits is  $\mu_y$  in the first population and  $\mu_z$  in the second population. The horizontal axis shows the difference of these two values. Both traits are normally distributed with standard deviation 10 in both populations. The plot shows results for four different sample sizes per population ( $n = nyj = nzj$ ).

To alleviate this problem, the associated MATLAB code offers two approaches. First, it provides a corrected estimate  $\Delta \tilde{p}_{y,z*}$ , which is less strongly biased than

$\Delta p_{y,z}$ , and is calculated as the maximum of 0 and  $2\Delta p_{y,z} - \bar{b}$ , where  $\bar{b}$  is the mean of the bootstrap distribution (e.g., Shao 1988). Second, we provide a p-value from a permutation test of the null hypothesis  $\Delta p_{y,z} = 0$ . If this test is non-significant, it might be reasonable to set the lower confidence limit to 0. (We note two caveats, though: First, permutations tests assume equal variance [or rather, exchangeability of the data under the null hypothesis]. Second, the bootstrap and permutation approaches may yield conflicting results; for example, when applied to the univariate  $\Delta p_{y,z}$ , it is possible that the test rejects the null hypothesis while, simultaneously, the bootstrap CI includes 0, or vice versa.)

For larger inter-population differences, the bias in  $\Delta p_{y,z}$  disappears, and CI coverage is generally good. In contrast to the univariate case, we find that the percentile method performs better than the BC and BCa methods (which behave very similar to each other and tend to be overoptimistic if  $\Delta p_{y,z}$  is small). For this reason, the default output of the MATLAB code will always show the percentile CIs for the multivariate estimator.

## References

- Algina J, Kesselman HJ, Penfield RD, 2006. Confidence interval coverage for Cohen's effect size statistic. *Educational and Psychological Measurement* 66: 945–960.
- Bradley JV, 1978. Robustness? *British Journal of Mathematical and Statistical Psychology* 31: 144–152.
- Efron B, 1987. Better bootstrap confidence intervals. *Journal of the American Statistical Association* 82: 171–185.
- Carpenter J, Bithell J, 2000. Bootstrap confidence intervals: When, which, what? A practical guide for medical statisticians. *Statistics in Medicine* 19: 1141–1164.
- Hess MR, Hogarty KY, Ferron JM, Kromrey JD, 2007. Interval estimates of multivariate effect sizes: Coverage and interval width estimates under variance heterogeneity and nonnormality. *Educational and Psychological Measurement* 67: 21–40.
- Kelley K, 2005. The effects of nonnormal distributions on confidence intervals around the standardized mean difference: Bootstrap and parametric confidence intervals. *Educational and Psychological Measurement* 65: 51–69.
- Reiser B, 2001. Confidence intervals for the Mahalanobis distance. *Communications in Statistics – Simulation and Computation* 30: 37–45.
- Shao J, 1988. Bootstrap variance and bias estimation in linear models. *The Canadian Journal of Statistics* 16: 371–382.