

1 Measuring individual identity information in animal signals: Overview
2 and performance of available identity metrics
3

4 Pavel Linhart¹, Tomasz Osiejuk¹, Michal Budka¹, Martin Šálek^{2,3}, Marek Špinka^{4,5}, Richard Policht^{4,6},
5 Michaela Syrová^{4,7}, Daniel T. Blumstein^{8,9}

6

7 Affiliations:

8 1 Department of Behavioural Ecology, Adam Mickiewicz University, Umultowska 89, 61-614, Poznań,
9 Poland

10 2 The Czech Academy of Sciences, Institute of Vertebrate Biology, Květná 8, 603 65 Brno, Czech
11 Republic

12 3 Faculty of Environmental Sciences, Czech University of Life Sciences Prague, Kamýcká 1176,
13 Suchbátka, 16521 Prague, Czech Republic

14 4 Department of Ethology, Institute of Animal Science, Přátelství 815, Prague, Uhřetěves, 104 00,
15 Czech Republic

16 5 Department of Ethology and Companion Animal Science, Faculty of Agrobiological Sciences, Food and Natural
17 Resources, Czech University of Life Sciences Prague, Kamýcká 129, 165 21 Prague 6, Czech Republic

18 6 Department of Game Management and Wildlife Biology, Faculty of Forestry and Wood Sciences,
19 Czech University of Life Sciences Prague, Kamýcká 129, 165 21 Prague 6, Czech Republic

20 7 Department of Zoology, Faculty of Sciences, University of South Bohemia, Branišovská 31, České
21 Budějovice, 370 05, Czech Republic

22 8 Department of Ecology and Evolutionary Biology, University of California, 621 Young Drive South,
23 Los Angeles, CA 90095-1606, USA

24 9 Rocky Mountain Biological Laboratory, Box 516, Crested Butte, CO 81224, USA

25

26 Corresponding author: Pavel Linhart

27

28 Abstract

29 1. Identity signals have been studied for over 50 years but, and somewhat remarkably, there is
30 no consensus as to how to quantify individuality in animal signals. While there are a variety
31 of different metrics to quantify individuality, these methods remain un-validated and the
32 relationships between them unclear.

33 2. We contrasted three univariate and four multivariate identity metrics (and their different
34 computational variants) and evaluated their performance on simulated and empirical
35 datasets.

36 3. Of the metrics examined, Beecher's information statistic (H_S) performed closest to
37 theoretical expectations and requirements for an ideal identity metric. It could be also easily
38 and reliably converted into the commonly used discrimination score (and vice versa).

39 Although Beecher's information statistic is not entirely independent of study sampling, this
40 problem can be considerably lessened by reducing the number of parameters or by
41 increasing the number of individuals in the analysis.

42 4. Because it is easily calculated, has superior performance, can be used to quantify identity
43 information in single variable or in a complete signal and because it indicates the number of
44 individuals that can be discriminated given a set of measurements, we recommend that
45 individuality should be quantified using Beecher's information statistic in future studies.

46 Consistent use of Beecher's information statistic could enable meaningful comparisons and
47 integration of results across different studies of individual identity signals.

48 **Keywords:** Individual recognition, Social behavior, Identity signal, Beecher's Information Statistic,
49 Acoustic identification, Acoustic discrimination, Vocal individuality, Discriminant analysis

50

51 Introduction

52 The fact that individuals differ in consistent ways is a both a central attribute of life and one that
53 underlies a number of theoretically important questions such as explaining cooperative behavior or
54 understanding the evolution of sociality (Crowley et al., 1996; Bradbury & Vehrencamp, 1998;
55 Tibbetts, 2004). Such individuality can also be important in wildlife conservation as well when it is
56 used to help census or monitor individuals based on individually-distinctive traits (Terry & McGregor,
57 2002; Blumstein et al., 2011). And, because, animals may base their decisions on the identity of the
58 individual with whom they interact or respond to (Wilkinson, 1984; Godard, 1991), there may be
59 selection to both produce individually-distinctive signals, and selection to discriminate among them
60 (Tibbetts & Dale, 2007; Wiley, 2013).

61 Quantification of individual identity (individuality) requires the assessment of variation in one or
62 more traits between at least two individuals. For identity signals to function properly, they should
63 maximize the between-individual variation and minimize the within-individual variation (Beecher,
64 1982, 1989). A variety of identity metrics have proliferated because of recognized biases (e.g., it is
65 more likely to find similar individuals in larger populations and, hence, it will be more difficult to
66 discriminate individuals in large populations or studies involving more individuals). These biases
67 make the comparison of results among studies unreliable (Beecher, 1989; Mathevon, Koralek,
68 Weldele, Glickman, & Theunissen, 2010). Additionally, some existing metrics were considered
69 unsuitable for a particular signal type (Searby & Jouventin, 2004). Nevertheless, new alternatives
70 were not always thoroughly tested and were not shown to be superior to the metrics they attempted
71 to replace. Furthermore, there are methodological problems that result from the calculation of
72 particular identity metrics, and some studies have used different equations to calculate the same
73 identity metric. Thus, somewhat remarkably given its importance, there is no consensus about how
74 to properly measure identity. As a result, researchers have generally avoided quantitative
75 comparisons between studies (Insley, Phillips, & Charrier, 2003). In a few cases, researchers tried to
76 overcome problems with identity metrics in comparative analyses by using exactly the same methods

77 across involved species (Beecher, Medvin, Stoddard, & Loesche, 1986; Lengagne, Lauga, & Jouventin,
78 1997; Pollard & Blumstein, 2011). Thus, hundreds of isolated studies have been published on
79 individuality in animal signals but because they used different metrics there is limited prospect that
80 we can benefit from the cumulative evidence of these studies. The lack of a commonly used identity
81 metric is a major impediment toward understanding the evolution of identity signaling and indeed,
82 the evolution of individuality.

83 Here, we review previously developed univariate (quantifying individuality within a single trait)
84 and multivariate metrics (quantifying individuality across multiple traits) that have been used to
85 quantify individual identity information in signals and we test their performance on simulated and
86 empirical datasets. In particular, we examine the following metrics: F-value, Potential of individual
87 coding PIC, Beecher's information statistic H_s , Information capacity H_M , and Mutual information MI.
88 We further evaluate different computational variants found in the literature in case of PIC and H_s
89 (see Table 1 and Supplement 1 for a detailed overview of metrics and their variants).

Table 1. Overview of the identity metrics and their variants

Metric variant and equation	description	reference	IDmeasurer function
$F = \frac{MS_b}{MS_w}$	F from one-way ANOVA where the individual is treated as independent variable and trait as dependent variable; MS_b = between group mean squares; MS_w = within group mean squares; calcF	e.g., Miller, 1978	calcF
$PIC_{between\ tot} = \frac{CV_{between\ tot}}{CV_w}$	$CV_{between\ tot}$ = between individual coefficient of variation calculated from all data points; CV_w = within individual coefficient of variation	e.g., Robisson, Aubin, & Bremond, 1993	calcPICbetweentot
$PIC_{between\ means} = \frac{CV_{between\ means}}{CV_w}$	$CV_{between\ means}$ = between individual coefficient of variation calculated with means from each individual; CV_w = within individual coefficient of variation	e.g., Lein, 2008	calcPICbetweenmeans
$H_{Sntot} = \log_2 \sqrt{\frac{F + n_{tot} - 1}{n_{tot}}}$	F = ANOVA F-value; n_{tot} = total sample size	possible variant from Beecher, 1989	calcHSntot
$H_{Sngroups} = \log_2 \sqrt{\frac{F + n_{groups} - 1}{n_{groups}}}$	F = ANOVA F-value; n_{groups} = number of groups (individuals)	possible variant from Beecher, 1989; e.g., Pollard, Blumstein, & Griffin, 2010	calcHSngroups
$H_{Snpergroup} = \log_2 \sqrt{\frac{F + n_{pergroup} - 1}{n_{pergroup}}}$	F = ANOVA F-value; $n_{pergroup}$ = number of samples in each group (observations per individual)	possible variant from Beecher, 1989	calcHSnpergroup
$H_{Svarcomp} = \log_2 \frac{\sigma_T}{\sigma_W}$	σ_T = total variance in mixed model; σ_W = residual variance associated with random factor in mixed model	Beecher, 1989; Carter, Logsdon, Arnold, Menchaca, & Medellin, 2012	calcHSvarcomp
$H_{Snpergroup} = \log_2 \sqrt{\frac{F + n_{pergroup} - 1}{n_{pergroup}}}$	F = ANOVA F-value; $n_{pergroup}$ = number of samples in each group (observations per individual); original variables are subjected to PCA to get uncorrelated components and H_s is calculated and summed over each independent component	Beecher, 1989	calcHSnpergroup
$H_M = \log_2 \sqrt{\frac{F_M + n - 1}{n}}$ $F_M = \frac{n - 1}{g - 1} * \frac{dist_t - g * dist_w}{dist_w}$	$dist_t$ = sum of distances of all samples from their centroid; $dist_w$ = sum of distances of samples within individual to its centroid; n = number of observations; g = number of groups;	Searby & Jouventin, 2004	calcHM
$DS = \frac{c}{N}$	C = samples correctly classified by Discriminant analysis; N = total number of samples	e.g., Hafner et al., 1979	calcDS
$MI = \sum_{ij} \log_2 \frac{p(i,j)}{p(i) * p(j)}$	$p(i)$ = probability of predicted individual; $p(j)$ = probability of actual individual; $p(i,j)$ = probability of match between predicted and actual individual	Mathevon et al., 2010	calcMI

92 We compare the performance of metrics to hypothetical ideal identity information metric. The
93 main principle of measuring individual identity in continuous traits is to quantify the ratio of between
94 and within individual variation (Beecher, 1982, 1989; Robisson, Aubin, & Bremond, 1993; Searby &
95 Jouventin, 2004). Thus, an ideal individual identity metric should be expressed on a ratio scale with a
96 meaningful zero value, equivalent to the situation when there is no between individual variation.
97 Further, there is no expected upper limit for individuality. High between to within individual variation
98 ratio indicates easy discrimination of individuals.

99 The datasets for the assessment of individual identity in different species vary in properties such
100 as the number of individuals, the number of samples per individual, the number of variables
101 measured (i.e., number of individualistic traits) and the covariance between the multiple variables
102 measured. Hence, we further propose that an ideal identity metric should be robust or respond
103 predictably to these dataset parameters to allow meaningful comparisons between studies.
104 Therefore, an ideal identity metric: 1) should not be systematically biased by the sampling effort, i.e.,
105 there should be no systematic effects of number of individuals and number of calls per individual in a
106 study on individuality estimate, and the sampling should ideally only impact on precision of
107 individuality estimate; and 2) in the multivariate case, it should well capture the intrinsic
108 multidimensionality of identity signals. In particular, it should rise with number of meaningful
109 variables because each of the uncorrelated variables can encode another level of individual variation.
110 In addition, it should also decrease with covariance between the variables because increasing
111 covariance between variables essentially decreases the number of independent variables. For our
112 comparison we gave the same weight to all criteria because these are very basic requirements and
113 an ideal metric should fulfill all of them. In addition, we will list other potential pros and cons of each
114 metric to provide a comprehensive evaluation of existing metrics.

115 We also wished to see if each of two commonly used metrics (Beecher's information statistic H_s ,
116 and discrimination score DS) could be converted to the other metric. We focused on H_s and DS

117 metrics only. DS has been used in the vast majority of past studies and DS has been found to
118 correlate well with potentially unbiased H_5 in a previous study (Beecher, 1989). However, the
119 previous study only tested the relationship between H_5 and DS on datasets with equal number of
120 individuals and observations per individual, thus, ignoring the known biases associated with DS.
121 Reliable conversion of DS into potentially unbiased H_5 could facilitate comparative analyses of results
122 reported in past and future studies.

123 Material and methods

124 We used R for simulations and statistical analysis (R Core Team, 2012). Functions to calculate identity
125 metrics, associated functions and datasets are available within an IDmeasurer package. This package
126 is available on CRAN (<https://cran.r-project.org/web/packages/IDmeasurer/index.html>) and GitHub
127 (<https://github.com/pygmy83/IDmeasurer>).

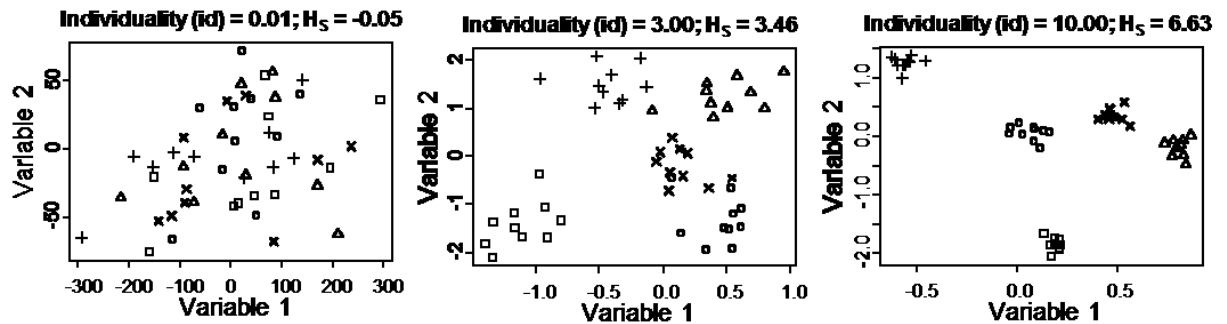
128 Datasets

129 **Simulated datasets.** Datasets were constructed to mimic typical data on individuality. Parameters of
130 datasets vary among studies. In particular, there are different numbers of individuals, observations
131 (samples) per individual, variables, and different covariances among variables. Effect of these
132 parameters was simulated along with individuality within datasets. The level of individual identity in
133 data was modified by changing the ratio of between and within individual variance in accordance
134 with theoretical assumptions of individual identity signals and previous studies (Beecher, 1989;
135 Searby & Jouventin, 2004). We developed R scripts involving “rnorm” and MASS package (Venables &
136 Ripley, 2002) “mvrnorm” function to generate the datasets. These functions generate random values
137 with a given standard deviation around pre-specified mean and, in “mvrnorm”, with pre-specified
138 covariance.

139 We constructed datasets with univariate and multivariate normal distributions with parameters
140 covering a wide range of values, specifically, five values for individuality ($id = 0.01, 1, 2.5, 5, 10$), five
141 values for number of observations per individual ($o = 4, 8, 12, 16, 20$), eight values for number of
142 individuals ($i = 5, 10, 15, 20, 25, 30, 35, 40$). Additionally, for multivariate datasets, five values for

143 covariance among variables ($cov = 0, 0.25, 0.5, 0.75, 1$) and five values for number of variables ($p = 2,$
144 $4, 6, 8, 10$). Thus, 200 and 5000 unique parameter combinations were possible in case of univariate
145 and multivariate datasets respectively. Individuality (id) represents the ratio of standard deviations
146 between and within individuals ($id = SD_{between} / SD_{within}$; $SD_{between}$ was calculated from means for each
147 individual and SD_{within} was set to be $SD_{between} / id$) (Fig. 1). A single covariance (cov) value was used in
148 the variance-covariance matrix to define covariances between all pairs of variables. For univariate
149 datasets, we first generated individual means for a predefined number of individuals “ i ” (normal
150 distribution, “ $rnorm$ ” function, mean = 1000, $SD_{between} = 1$) and then we generated a predefined
151 number of random observations “ o ” around each individual mean (normal distribution, “ $rnorm$ ”
152 function, mean = individual mean, $SD_{within} = SD_{between} / individuality$ “ id ”). In the multivariate case, we
153 first created a matrix representing mean individual values of variables for each of the individuals
154 (multivariate normal distribution, “ $mvrnorm$ ” function, mean for each variable = 0, variance-
155 covariance matrix). Variances on the diagonal of the covariance matrix were set equal to 1 (hence
156 $SD_{between} = 1$) and all covariances between variable pairs were set equal to the predefined covariance
157 “ cov ”. Then, we generated a predefined number of random observations “ o ” around each individual
158 and a variable mean (“ $rnorm$ ” function, mean = individual mean, $SD_{within} = SD_{between} / individuality$
159 “ id ”).

160 We asked how dataset parameters (i, o, p, cov, id) influenced the value of each identity metric.
161 To explore this, 20 randomization cycles were run for each unique combination of parameter values.
162 For example, in the multivariate case, $20 * 5000 = 100\ 000$ independent datasets were generated
163 (datasets 1-20: $i = 5, o = 4, p = 2, cov = 0, id = 0.01$; datasets 21-40: $i = 10, o = 4, p = 2, cov = 0, id =$
164 0.01 ; ... ; datasets 99 981-100 000: $i = 40, o = 20, p = 10, cov = 1, id = 10$). Identity metrics were
165 calculated for each dataset.

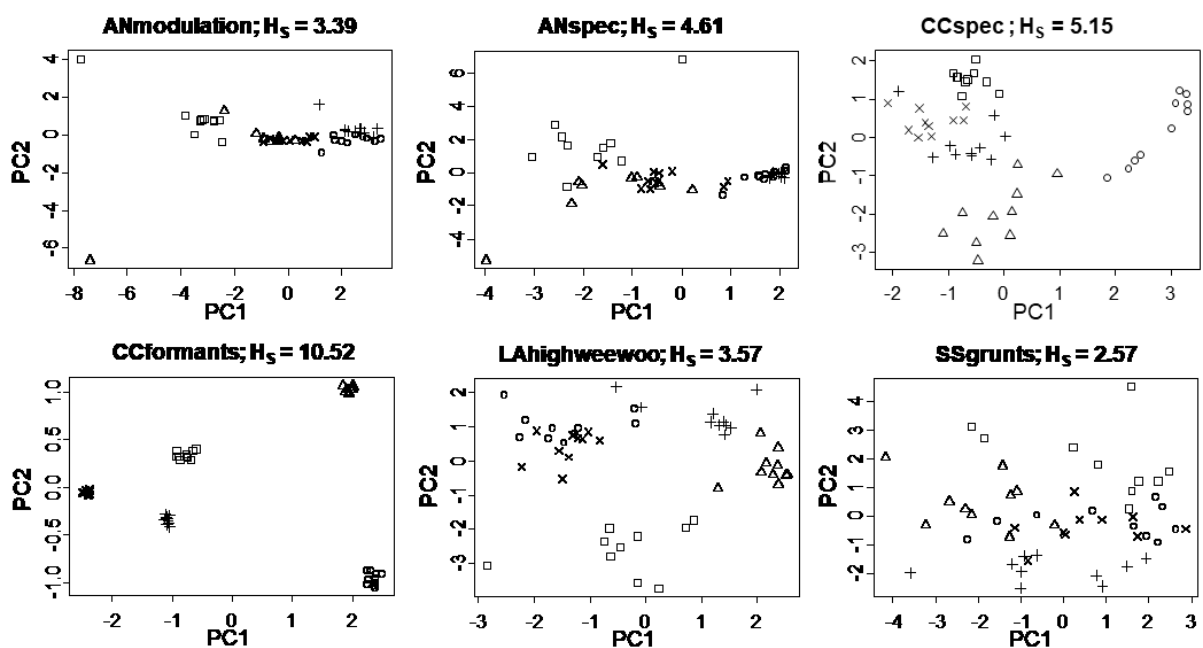


166

167 **Figure 1.** Illustration of three artificial multivariate datasets that differ only in the individuality used
 168 to generate datasets. Settings for the function generating these datasets: $i = 5$, $o = 10$, $p = 2$, $cov = 0$,
 169 $id = 0.01, 3, \text{ and } 10$.

170 **Empirical datasets.** While the general performance of identity metrics was evaluated on simulated
 171 datasets, empirical datasets were used to evaluate the consistency of DS and H_5 metrics and
 172 reliability of H_5 and DS conversion on real data. We used six empirical datasets from four different
 173 species: little owls *Athene noctua* (ANmodulation, ANspec) (Linhart & Šálek, 2017), corncrake *Crex*
 174 *crex* (CCformants, CCspec) (Budka & Osiejuk, 2013), yellow-breasted boubous *Laniarius atrofasciatus*
 175 (LAhighweewoo) (Osiejuk et al. unpublished data), and domestic pigs *Sus scrofa* (SSgrunts) (Syrová,
 176 Policht, Linhart, & Špinková, 2017) (Figure 2). In two species – corncrakes and little owls – calls were
 177 described by two different sets of variables. In little owls, we described calls by frequency
 178 modulation by measuring fundamental frequency at 10 measuring points evenly spread through the
 179 duration of the call (ANmodulation) or parameters describing the distribution of the frequency
 180 spectrum such as peak frequency, minimum and maximum frequencies and frequencies dividing
 181 spectrum by energy content (ANspec). In corncrakes, we used formants (CCformants) and
 182 parameters describing the distribution of the frequency spectrum (CCspec) (see the Supplement 2 for
 183 detail description of empirical datasets). Because datasets varied with respect to the number of
 184 individuals (33 – 100) and the number of calls per individual available (10 – 20), we scaled all datasets
 185 down to lowest common denominator by randomly selecting individuals and calls from bigger
 186 datasets. Eventually, each dataset had 33 individuals and 10 calls per individual. Each dataset also

187 used different numbers of variables to describe the calls' acoustic structure (ANmodulation = 11,
 188 ANspec = 7, CCformants = 4, CCspec = 7, LAhighweewoo = 7, SSgrunts = 10). In all these empirical
 189 datasets, assumptions of multivariate normality were tested (Korkmaz, Goksuluk, & Zararsiz, 2014),
 190 but not met. We found various issues on the level of univariate variables and the whole dataset. For
 191 instance, there were issues with outliers, skewness, kurtosis, and multimodal distributions (see
 192 Supplement 2 for univariate histograms and multivariate Chi-square Q-Q plot). Normality issues are
 193 common for research studies on acoustic individual identity. Authors deal with normality issues by
 194 eliminating problematic variables (e.g., Couchoux & Dabelsteen, 2015), using non-parametric
 195 classification methods (e.g., Mielke & Zuberbuehler, 2013), or by relying on robustness of cross-
 196 validated DFA and PCA towards relaxed assumptions (e.g., Mathevon et al., 2010). We used the last
 197 approach. If the assumptions of discriminant analysis are not met the results should be less stable
 198 when using different sampling and hence our results should be viewed as conservative.



199

200 **Figure 2.** Illustration of empirical datasets. Five individuals were randomly sampled from each
 201 dataset of 33 individuals and all 10 calls per individual were selected. H_5 for a full dataset is shown.
 202 Data were centered and scaled and subjected to PCA. The first two Principal Components are
 203 plotted.

204 Statistical analysis
205 The relationship between a given identity metric and each of the parameters was assessed
206 graphically by plotting the mean value and the 95% confidence intervals of an identity metric against
207 all of the modelled data parameters separately. We then used a one-way ANOVA to test whether an
208 identity metric was constant across all levels of a parameter. One-way ANOVA along with graphical
209 evaluation of relationships between metrics and model parameters was preferred over multivariate
210 regression because it simply, but adequately, addresses our main question (i.e., does the metric
211 change in response to model parameter?) without the need to specify and compare many different
212 multivariate regression models. If we found significant differences, we followed up these with post-
213 hoc Tukey tests to identify which parameter levels differed. Due to the large number of comparisons,
214 we only reported comparisons of neighboring parameter levels. We used linear and non-parametric
215 loess regression to convert H₅ to DS and vice versa. Loess regression identifies a function that best
216 describes complex data by fitting simple models to sequential subsets of data. Its main advantage is
217 that it does not require specifications of the function and, hence, it is suitable for modeling of
218 complex relationships. Loess regression included the number of individuals and the number of calls
219 per individual as additional predictors. We used Spearman correlation coefficients to quantify
220 between-metric consistency of ranking individuality in datasets. Pearson correlations were used to
221 assess consistency within identity metrics in full and partial datasets. We then used Friedman tests,
222 followed by a series of Wilcoxon tests (for post-hoc comparison of differences between levels), to
223 compare correlation coefficients obtained for each pair of the metrics.

224 Results

225 The comparison of available univariate and multivariate metrics to an ideal metric is shown in Table
 226 2.

	zero	limit	id	cov	p	o	i	points	pros	cons
Univariate Metrics										
ideal	y	n	+			ns	ns	5/5		
F	y	n	+			+	ns	4/5		sample dependent
PIC _{between_{tot}}	y	n	+			ns	ns	5/5	intuitive and straightforward calculation; allows separate assessment of within and between individual variation	not meaningful for variables with positive and negative values; cannot be summed or averaged over different variables = univariate only
PIC _{between_{means}}	n	n	+			ns	ns	4/5		Converges to non-meaningful value for no individuality in data
H _{S_{tot}}	y	n	+			ns	-	4/5		sample dependent; incorrect HS variant
H _{S_{npergroup}}	y	n	+			ns	ns	5/5	standard variant of HS; univariate and multivariate	
H _{S_{ngroups}}	y	n	+			+	-	3/5		sample dependent; incorrect HS variant
H _{S_{varcomp}}	y	n	+			ns	ns	5/5	allows including various covariates in mixed models	values twice as big as in case of standard HS _{npergroup}
Multivariate Metrics										
ideal	y	n	+	-	+	ns	ns	7/7		
DS	n	y	+	-	+	+	-	3/7	population and individual metric; the most commonly used metric	sample dependent; not suitable for high individuality signals because values are limited from the top
H _S	y	n	+	-	+	ns	+	6/7	univariate and multivariate; partial sample dependence is introduced by PCA but can be to large extent eliminated; biologically meaningful - provides number of unique individual signatures within population; good theoretical framework for both discrete and continuous individuality traits	partially sample dependent
H _M	y	n	+	ns	ns	ns	ns	5/7	sample independent; various types of similarity metrics can be potentially used (euclidean distances, Jaccard similarity, string edit distance, dynamic time warping, etc.)	number of independent variables needs to be known to calculate total identity information
MI	n	y	+	-	+	-	+	3/7	could be applied with various classification methods	sample dependent; not suitable for high individuality signals because values are limited from the top

227 'zero' – metric has a meaningful zero; 'limit' – metric is limited from the top by an asymptote; 'id' – change in response to increasing identity information in data; 'cov' – response to increasing covariance between variables; 'p' – response to increasing number of variables; 'o' – response to increasing number of calls per individual; 'i' – response to increasing number of individuals; 'y' - yes; 'n' – no; '+' – increase; '-' – decrease; 'ns' – not significant, does not change with a parameter.

228 **Table 2.** The comparison of available univariate and multivariate metrics to a hypothetical ideal metric and
 229 summary of their pros and cons. We summed the number of matches (points) to compare different metrics to
 230 the ideal metric.

231 Univariate metrics

232 All explored univariate metrics increased with increasing individuality in the data. However, only

233 PIC_{between_{tot}}, PIC_{between_{means}}, H_{S_{npergroup}} and H_{S_{varcomp}} estimates were independent of the number of calls

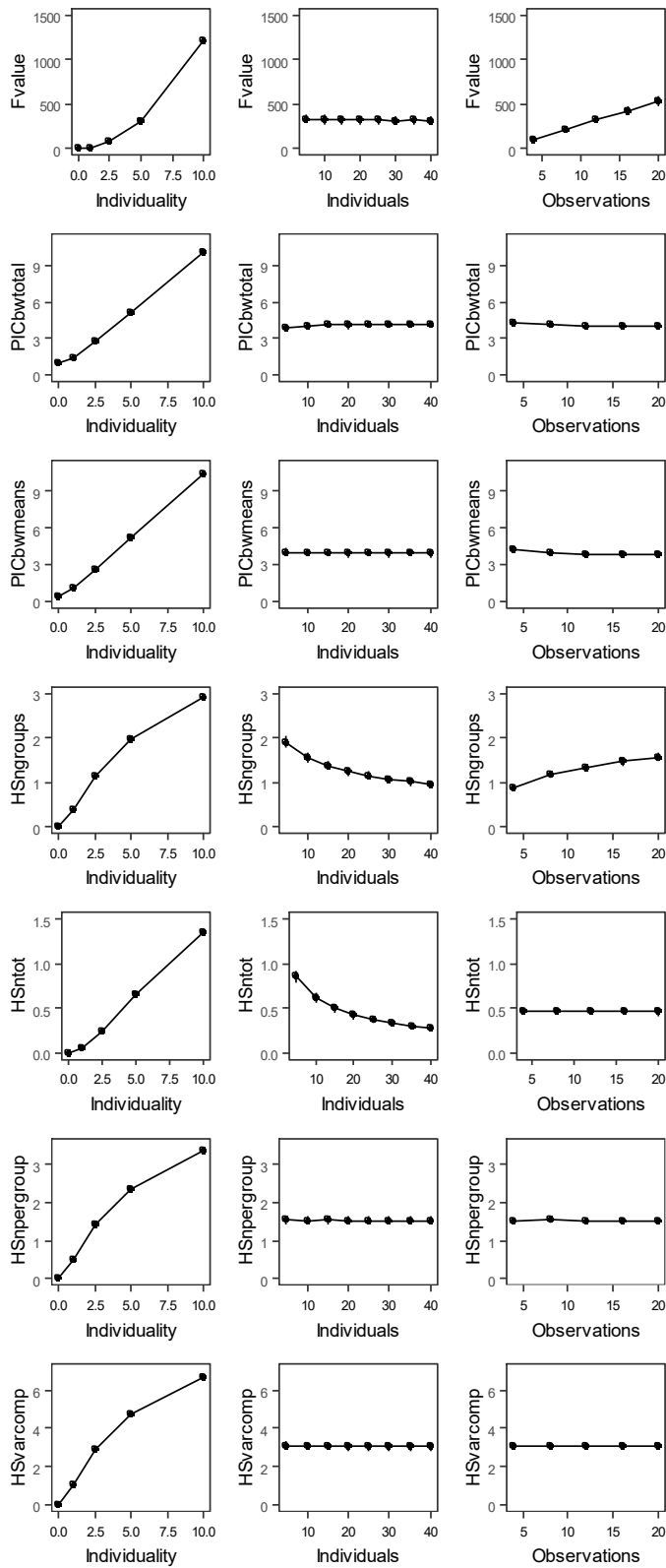
234 and the number of individuals used to calculate the metric (Figure 3). These general patterns were

235 qualitatively identical when all simulated data were pooled or if only one of the parameters (number

236 of calls, number of individuals, individuality) was changed at a time and the others were kept

237 constant at the middle value (see Supplement 3 for detailed results including ANOVA tests).

238 All four sampling-independent metrics ($PIC_{\text{betweentot}}$, $PIC_{\text{betweenmeans}}$, $H_{\text{Snpergroup}}$ and H_{Svarcomp}) were
239 highly correlated (Spearman correlation, all $r > 0.99$). $H_{\text{Snpergroup}}$ and H_{Svarcomp} correctly converged to 0
240 in the case when individuality was set to be negligible ($id = 0.01$), while $PIC_{\text{betweentot}}$ and $PIC_{\text{betweenmeans}}$
241 converged to higher values (1.01 and 0.32 respectively). $PIC_{\text{betweentot}}$ reflects the number of potential
242 individual signatures within a population in same way as 2^{H_s} does (Beecher, 1989), and, both,
243 $PIC_{\text{betweentot}}$ and $2^{H_{\text{Snpergroup}}}$ reflect the ratio of between to within individual variation. Hence,
244 convergence of $PIC_{\text{betweentot}}$ to 1 could be also seen as desirable quality and meaningful value for a
245 signal with no individuality. H_{Svarcomp} was equal to $2 * H_{\text{Snpergroup}}$ (see Supplement 4 for details). We
246 further considered only the $H_{\text{Snpergroup}}$ variant in multivariate analyses.



247

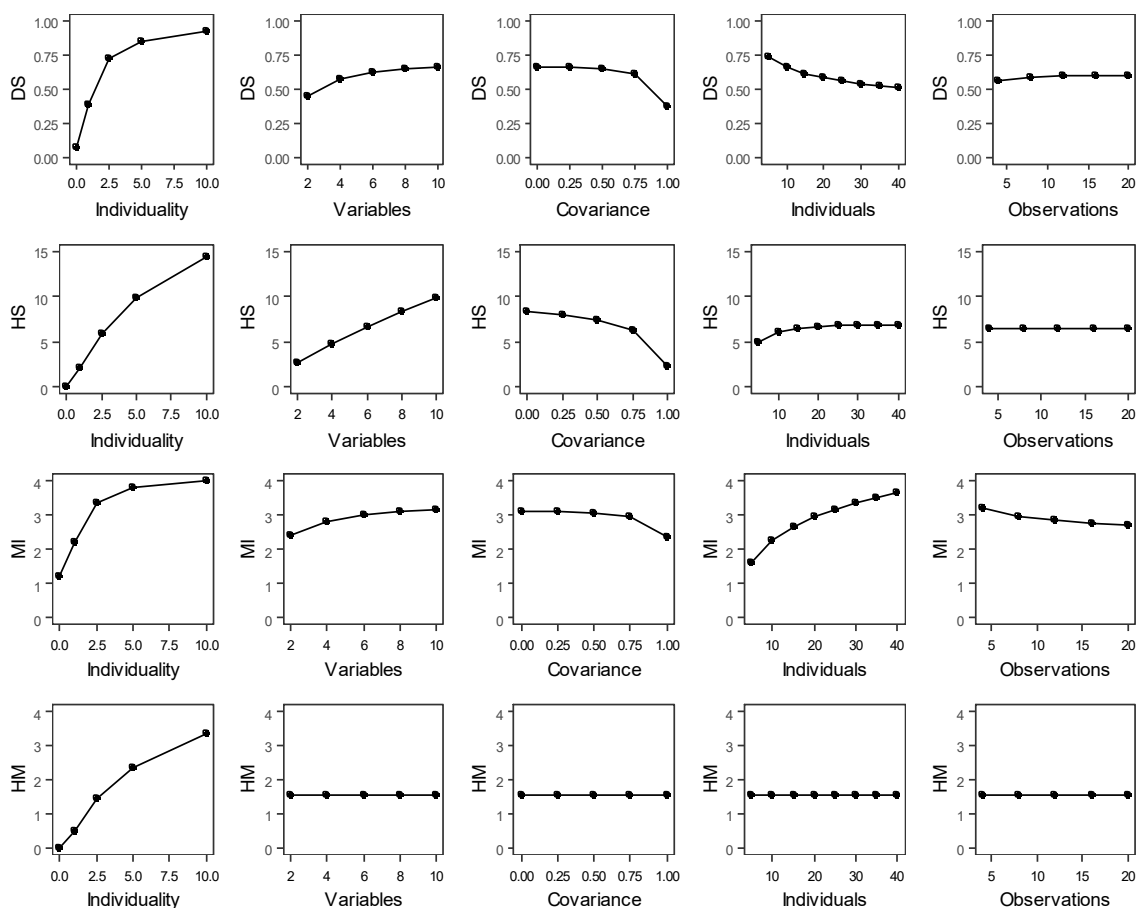
248 **Figure 3.** Variation in univariate identity metrics in response to simulated dataset parameters:

249 individuality, number of calls per individual, and number of individuals. Means and 95% confidence

250 intervals are shown. Graphs were plotted using all simulated univariate data pooled together. For the
 251 graphs with only a single parameter changing at a time see Supplement 3.

252 Multivariate metrics

253 The performance of multivariate identity metrics is illustrated in Figure 4. All metrics increased with
 254 increasing individuality. DS, H_S , and MI increased with increasing number of variables available and
 255 decreased with increasing covariance between variables. Only H_M did not change in response to
 256 increasing the number of individuals. H_S and H_M did not change in response to increasing the number
 257 of calls per individual. These general patterns were qualitatively identical when all simulated data
 258 were pooled or if only one dataset parameter was changed at a time and others were kept constant
 259 at the middle value (see Supplement 5 for detailed results including ANOVA tests).



260

261

262 **Figure 4.** Multivariate identity metrics in response to simulated dataset parameters: individuality,
263 covariance between variables, number of variables, number of calls per individual, and number of
264 individuals. Means and 95% confidence intervals are shown. Graphs were plotted using all simulated
265 multivariate data pooled together. For the graphs with only a single parameter changing at a time
266 see Supplement 4.

267 Despite the different response of metrics to some of the simulated parameters, there was still
268 moderate to high agreement among metrics about identity content in the data (Spearman
269 correlations, mean $r \pm SD = 0.82 \pm 0.07$; minimum $r = 0.71$ for correlation between DS and MI;
270 maximum $r = 0.95$ for correlation between DS and H_5). H_5 had the greatest correlations with other
271 metrics (average $R = 0.88$). We found no advantage to using H_M over H_5 as previously suggested.
272 Instead, H_M was equal to H_5 per variable ($H_M = H_5 / p$) in data with zero covariance between variables.
273 (Supplement 6).

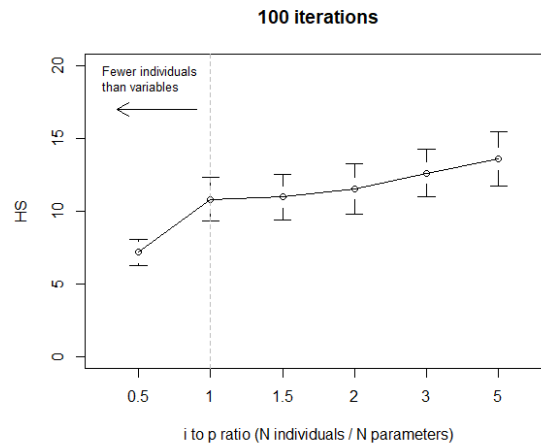
274 Thus, our simulations show that H_5 matched the characteristics of the ideal metric in 6/7 cases,
275 followed by H_M (5/7), DS (4/7), and MI (both 3/7) (Table 1).

276 Potential for removing bias in H_5

277 We observed no significant association between H_5 and the number of individuals in the univariate
278 case so we investigated the origin of the sampling bias in the multivariate case. This bias was only
279 present when data were subjected to Principle Component Analysis (PCA). However, PCA is required
280 to create uncorrelated components for H_5 calculation.

281 It is possible that the more variables measured, the more individuals need to be sampled in order
282 to reduce this bias. We therefore fixed the number of variables to 5, 10, and 20 ($p = 5, 10, 20$) and
283 varied the ratio of the number of individuals to the number of variables 'i to p ratio' from 0.5 to 5 ('i
284 to p ratio' = 0.5, 1, 1.5, 2, 3, 5) by using different numbers of individuals in our simulations ($i = 3, 5, 8,$
285 10, 15, 20, 25, 30, 40, 50, 60, 100 depending on number of variables and "i to p ratio"). The number
286 of calls per individual was set to 10. Individuality and covariance were both chosen randomly in each

287 iteration from predefined intervals used in the earlier simulations (covariance range = [0, 0.25, 0.5,
 288 0.75, 1]; individuality range = [0.01, 1, 2.5, 5, 10]). We used 100 iterations for each 'i to p ratio'. H_S did
 289 not rise significantly after the number of individuals reached at least the number of parameters
 290 (One-way ANOVA $F_{5, 1794} = 7.68, P < 0.001$; no significant differences between levels if 'i to p' ≥ 1 , all p
 291 > 0.132) (Figure 5).



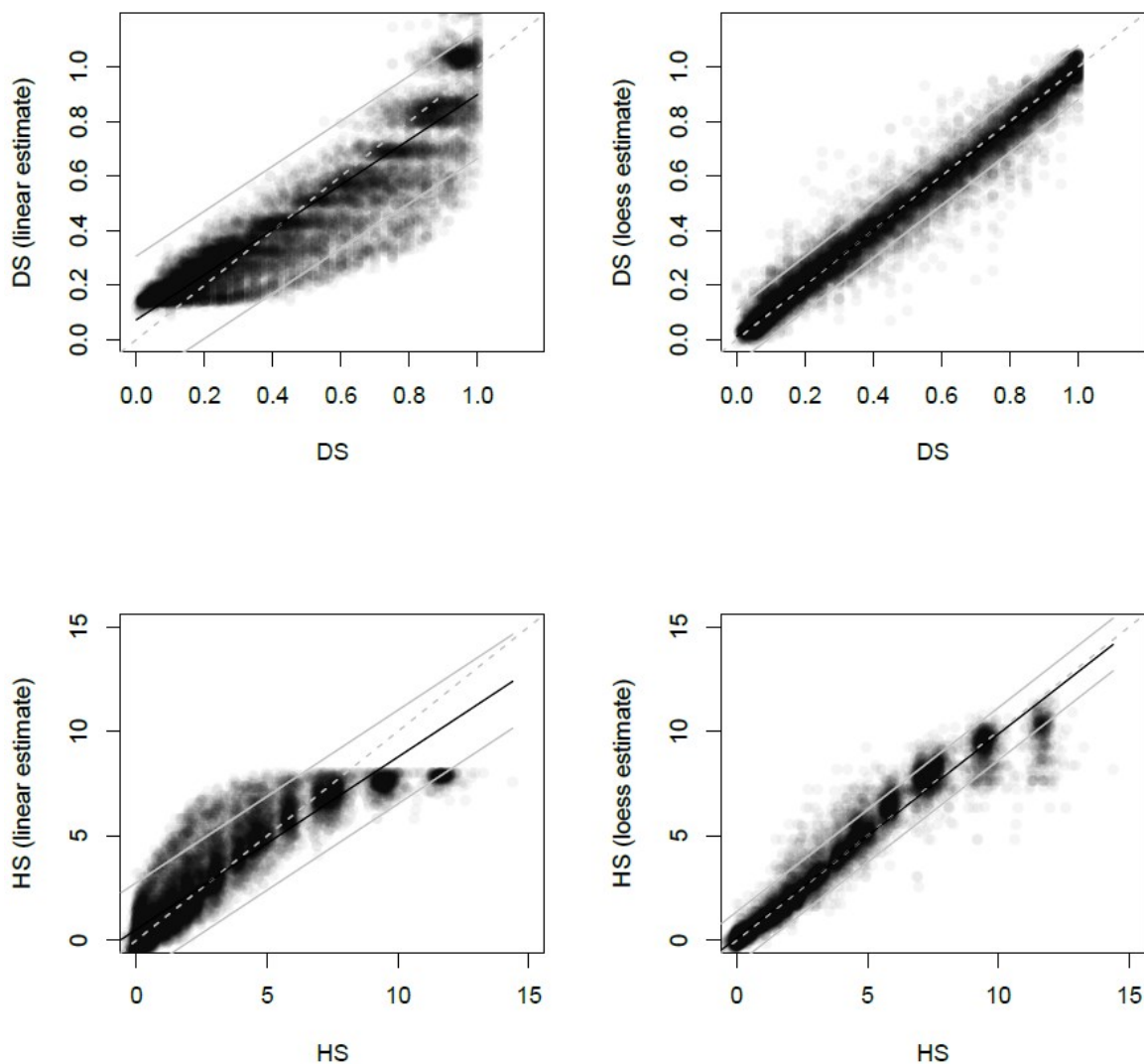
292

293 **Figure 5.** H_S and 'i to p ratio' (number of individuals / number of variables). H_S was under-estimated if
 294 there are fewer individuals than variables. Means and 95% confidence intervals are shown.

295 Converting DS to H_S and vice versa

296 We used simple linear regression and non-parametric loess regression to estimate H_S based on DS
 297 and vice versa. There was a previously suggested linear relationship that had a limit of $H_S = 8$ where
 298 the DS values were 100% correct discrimination (Beecher 1989). Because the H_S values in our original
 299 simulated datasets far exceeded 8, we generated a new set of simulated datasets with individuality
 300 ranging between 0.1 and 2 (id = 0.1, 0.25, 0.5, 0.75, 1, 1.33, 1.66, 2), covariance set to zero (cov = 0),
 301 number of iterations was reduced to 10 (it = 10), and other parameters were set as in previous
 302 models ($p = 2, 4, 6, 8, 10$; $i = 5, 10, 15, 20, 25, 30, 35, 40$; $o = 4, 8, 12, 16, 20$). These settings led to H_S
 303 values up to about 13 for data used for model building, and H_S values up to about 14 in the case of
 304 data used for model testing. These values are much closer to 8 and also much closer to H_S values
 305 reported from nature.

306 Loess models took into account the number of calls per individual and the number of individuals.
 307 We compared the loess conversion and linear conversion models of DS and H_s. In general, loess
 308 estimates were closer to the ideal prediction (intercept = 0, beta = 1) and the loess model reduced
 309 error of both DS and H_s estimates to about a half compared to linear estimates (Figure 6). Both H_s
 310 estimates were underestimated for high values of H_s. The ceiling value is clearly apparent for linear
 311 estimates of H_s. It is still visible in case of loess estimates but loess predictions remain reasonably
 312 good up to about H_s = 10.



313 **Figure 6.** Estimation of H_s and DS based on linear and loess transformation of DS and H_s respectively
 314 for datasets with H_s up to 14.4. **Linear DS estimation:** Intercept = 0.07, Beta = 0.83, R² = 0.83,
 315

316 Standard Error of Estimate (SEE) = 0.12, 95% Prediction interval = predicted value \pm 0.23; **DS loess**
317 **estimation:** Intercept = 0.01, Beta = 0.98, $R^2 = 0.97$, Standard Error of Estimate (SEE) = 0.05, 95%
318 Prediction interval = predicted value \pm 0.10. **Linear H_s estimation:** Intercept = 0.51, Beta = 0.83, $R^2 =$
319 0.83, Standard Error of Estimate (SEE) = 1.14, 95% Prediction interval = predicted value \pm 2.24; **HS**
320 **loess estimation:** Intercept = 0.11, Beta = 0.98, $R^2 = 0.95$, Standard Error of Estimate (SEE) = 0.64,
321 95% Prediction interval = predicted value \pm 1.26.

322 Correlations between calculated and estimated metrics
323 We were further interested in how H_{sest} and DS_{est} might represent H_s and DS of a particular sample of
324 individuals or H_{sfull} and DS_{full} of the whole population. For this purpose, we first generated 50 full
325 datasets with different identity levels representing 50 hypothetical populations of different species.
326 Each dataset comprised 40 individuals, 20 calls per individual, and 10 parameters. For these datasets,
327 individuality was set randomly ranging between 0.2 – 2 (0.1 increments), and the covariance was set
328 randomly ranging between 0.2 – 0.8 (0.1 increments). These settings generated datasets with H_{sfull}
329 values that ranged from 0.22 – 9.89 (mean \pm sd: 4.72 ± 2.95). Then, we repeatedly subsampled full
330 datasets to get partial datasets which simulated different sampling of the population. We
331 subsampled 5-40 individuals and 4-20 calls per individual per dataset in each of total 20 iterations.
332 We also repeatedly subsampled our empirical datasets. We subsampled 5-33 individuals and 4-10
333 calls per individual per dataset in each of total 20 iterations. The number of parameters was not
334 randomized – we always kept the original number of variables.

335 In simulated datasets, H_s and H_{sest} were correlated almost perfectly with each other and with
336 H_{sfull} (all average Pearson $r > 0.97$). There was no difference among the correlation coefficients from
337 correlations between H_{sfull} , H_s , and H_{sest} (Friedman Chi Square = 3.6, $p = 0.165$). In empirical datasets,
338 H_s calculated on partial datasets still reflected the H_{sfull} almost perfectly (average Pearson $r = 0.99$).
339 While H_{sest} reflected H_s of partial dataset (average Pearson $r = 0.90$), and H_{sfull} (average Pearson $r =$
340 0.88) slightly worse, it remained a reasonable fit. However, H_{sest} did not reflect H_{sfull} as precisely as it
341 did H_s (Friedman Chi Square = 33.6, $p < 0.001$, post-hoc test: $H_s - H_{sfull}$ vs. $H_{sest} - H_{sfull}$, $p < 0.001$).

342 DS in simulated datasets was almost perfectly correlated with DS_{est} (average Pearson $r = 0.99$).
343 Although the relationship between DS in full datasets (DS_{full}) and DS and DS_{est} was significantly worse
344 (Friedman Chi Square = 40.0, $p < 0.001$; both post-hoc tests: $p < 0.005$), these associations remained
345 strong (DS_{full} and DS: average Pearson $r = 0.95$; DS_{full} and DS_{est} : average Pearson $r = 0.96$). In empirical
346 datasets, the correlation between DS and DS_{est} was lower than in case of artificial datasets (average
347 Pearson $r = 0.91$). DS and DS_{est} of partial datasets had comparable correlations to DS_{full} (DS_{full} and DS:
348 average Pearson $r = 0.88$; DS_{full} and DS_{est} : average Pearson $r = 0.86$). Thus, the performance of DS and
349 DS_{est} to reflect each other or DS_{full} did not differ (Friedman Chi Square = 0.9, $p = 0.638$).

350 Discussion

351 We provided an overview of the metrics used to quantify individual identity in animal signals in order
352 to identify the best method for reporting individuality in animal signals. Biases associated with some
353 of the commonly used metrics, and the use of different metrics across studies, makes it difficult to
354 compare results and integrate the accumulated knowledge from the numerous published studies on
355 individual identity in animal signals. We show that the assessment of individual identity is relatively
356 straightforward when considering a single trait (univariate case). Both, PIC ($PIC_{betweentot}$) and H_S
357 ($H_{S_{npergroup}}$), performed according to expectations. Multivariate identity metrics based on direct
358 quantification of between to within individual variation ratios (H_S , H_M) performed better than the
359 metrics derived from discrimination of individuals (DS, MI). We confirmed sampling-associated biases
360 where they were reported previously (DS), but we found them even in metrics that had been
361 developed to overcome these biases (H_S , MI). We also described yet unrecognized issues (the need
362 to assess dimensionality for H_M to quantify the total individuality of a signal). We further found that
363 some metrics created values that were so close that they could be viewed as redundant (PIC and H_S ;
364 H_M and H_S) and using them simultaneously brings unnecessary confusion to the field.

365 Based on our review and systematic analysis, we suggest H_S should be routinely reported as the
366 standard individual identity metric because it performed closest to an ideal identity metric in the
367 univariate, as well as in the multivariate case. The partial bias in H_S caused by the number of

368 individuals in a study could be removed by having at least the same number of individuals as the
369 number of variables. H_5 was the most consistent metric and correlated the best with DS and other
370 identity metrics. Further, H_5 could be converted reliably into DS if needed.

371 The robustness of H_M towards sampling bias (number of individuals, number of calls, as well as
372 the number of variables and covariance) is an attractive feature. However, as we show, H_M quantifies
373 identity information per variable and not the identity information of the entire signal. It is necessary
374 to know the effective number of variables to calculate the total identity information of a signal (i.e., if
375 there is perfect covariance between the variables, the effective number of variables is 1 no matter
376 how many variables are used), which may be difficult to assess. On the other hand, H_M uses distances
377 (similarity scores) of samples to calculate individuality and, hence, it could be potentially used not
378 just with Euclidean distances (Searby & Jouventin, 2004, this study) but also together with other
379 various methods assessing similarity (e.g., cross-correlation, dynamic time warping, or string edit
380 distances).

381 Mutual information (MI) is derived from a confusion matrix of discrimination analysis and we
382 show it has similar shortcomings as discrimination scores. Our results that found systematic biases in
383 MI are in line with previous studies that investigated measures of clustering for various machine
384 learning purposes where potentially unbiased variants of MI are constantly searched for (e.g., Amelio
385 & Pizzuti, 2017).

386 **Identity metrics in comparative analyses.** We show that biases associated with DS (the most
387 often used metric) and H_5 (the best metric) are not necessarily fatal for comparisons of different
388 published studies because H_5 and DS values that are based on an entire population or subsamples
389 from a population were well correlated in both simulated and empirical datasets. Additionally, the
390 conversion of sample biased DS values into less biased H_5 values could allow better comparisons
391 between studies. Both H_5 and H_M values were previously found to correlate well with DS (Beecher,
392 1989; Searby & Jouventin, 2004). We extend previous findings for H_5 (Beecher, 1989) to situations

393 with unequal sampling and we show it is possible to convert between H_s and DS with an acceptable
394 amount of error even when datasets differ in the number of individuals and calls per individual, and
395 have important issues associated with multivariate normality (Supplement 2). Discriminant analysis
396 (DA) and Principal component analysis (PCA) used for DS and H_s calculations both assume
397 multivariate normality for optimal results. While using these methods with non-normal data cannot
398 be, in general, recommended, relatively high correlations between our metrics in empirical datasets
399 suggest that DA and PCA scores were quite robust to these normality issues. Discrimination and
400 dimensionality reduction analytical techniques that are able to handle normal and non-normal data
401 definitely need to be considered in future individual identity studies.

402 **Future individual identity metrics.** We hope that our study will stimulate further discussions
403 about how individual identity should be properly measured. Although we suggest that H_s should be
404 generally used to quantify individuality, different metrics or more complex approaches might be
405 required for particular interesting questions. For example, H_s can only provide a population estimate
406 of individual identity. Researchers might be interested in whether distinctiveness of individuals
407 increases during ontogeny (Syrová et al., 2017). In this case, discrimination scores can be reported for
408 each individual, thus making statistical evaluation possible. Furthermore, separate assessments of
409 within- and between individual variations when calculating PIC might be useful to test hypotheses
410 about which of the two has been selected for. Within-individual variation could be reduced by, for
411 example, ritualized behavior while between-individual variation could be increased through, for
412 example, morphological variation in structures producing or carrying the signal (e.g., Sheehan &
413 Nachman, 2014). The dimensionality of identity signals might be an important factor for recognition
414 processes (Trunk, 1979) and evolution could favor low dimensional signals. Paralleling the
415 distribution of individuals in space (territoriality, living in colonies), individual signatures within a
416 population, too, could have random, clumped, or regular distributions depending on the mechanisms
417 behind individual distinctiveness and the degree of plasticity of identity signals.

418 We evaluated the efficacy of all metrics within the acoustic modality only. It is increasingly
419 recognized that signals may employ multiple modalities (Partan & Marler, 1999; Partan, 2013). All of
420 the identity metrics discussed here could be, in principal, used in visual or chemical domains as well.
421 H_5 has an advantage that it could be used both for discrete traits, such as color variants, presence of
422 particular alleles or chemicals, and for continuous traits such as size of visual patterns, duration of
423 calls, etc. (Beecher, 1982, 1989). However, identity information outside the acoustic domain has
424 been rarely quantified and meaningful comparison of individual identity across modalities remains a
425 challenge for the future.

426 It is likely that automatic data collection and analysis techniques will be increasingly applied for
427 various recognition tasks, including individual recognition (Elie & Theunissen, 2018; Stowell,
428 Petrusková, Šálek, & Linhart, 2019). While these methods will allow studying individual identity
429 signalling on unprecedented scales and sample sizes, the resulting classification accuracy scores will
430 be analogous to the discrimination score, with similar positives and drawbacks.. However, many
431 different feature sets, pre-defined or automatically derived from data, as well as many different
432 classification methods could be combined to test for the robustness of identity signals and/or to
433 mimic and test for different alternatives of possible real recognition processes (Elie & Theunissen,
434 2018).

435 **Conclusion.** We suggest that, at the current state of knowledge and methodology development,
436 H_5 should be generally reported as the the “golden standard” individual identity metric to allow the
437 best comparison of individuality in signals across different studies. Given that H_5 may not be
438 sufficient in all cases, we encourage further research to develop new metrics to quantify identity
439 information in signals. However, new metrics should always be appropriately assessed and their
440 performance directly compared to the best existing metrics. We provide datasets and scripts that
441 should help to assess individual identity information in animal signals and benchmark the future
442 metrics.

443 Acknowledgements

444 PL received funding from the European Union’s Horizon 2020 research and innovation programme
445 under the Marie Skłodowska-Curie grant agreement No. 665778 administered by the National
446 Science Centre, Poland (UMO-2015/19/P/NZ8/02507). DTB is supported by the NSF. MŠp, MS, and RP
447 were supported by Czech Science Foundation (GA14-27925S) and Czech Ministry of Agriculture (MZE-
448 RO0718). MŠá work was supported by the research aim of the Czech Academy of Sciences (RVO
449 68081766).

450 Authors’ contributions

451 PL and DTB conceived the ideas and designed methodology; PL, TO, MB, MŠá, MŠp, MS, and RP
452 collected the data; PL analysed the data; PL and DTB led the writing of the manuscript. All authors
453 contributed critically to the drafts and gave final approval for publication.

454 Data Accessibility statement

455 Data and code used for this article are available within IDmeasurer R package currently available on
456 CRAN (<https://cran.r-project.org/web/packages/IDmeasurer/index.html>) and GitHub
457 (<https://github.com/pygmy83/IDmeasurer>).

458 References

- 459 Amelio, A., & Pizzuti, C. (2017). Correction for closeness: Adjusting normalized mutual
460 information measure for clustering comparison. *Computational Intelligence*, 33(3),
461 579–601. doi:10.1111/coin.12100
- 462 Beecher, M. D., Medvin, M. B., Stoddard, P. K., & Loesche, P. (1986). Acoustic adaptations
463 for parent-offspring recognition in swallows. *Experimental Biology*, 45, 179–193.
- 464 Beecher, Michael D. (1982). Signature systems and kin recognition. *American Zoologist*,
465 22(3), 477–490.

466 Beecher, Michael D. (1989). Signaling systems for individual recognition - an information-
467 theory approach. *Animal Behaviour*, 38, 248–261. doi:10.1016/S0003-3472(89)80087-
468 9

469 Blumstein, D. T., Mennill, D. J., Clemins, P., Girod, L., Yao, K., Patricelli, G., ... Kirschel,
470 A. N. G. (2011). Acoustic monitoring in terrestrial environments using microphone
471 arrays: applications, technological considerations and prospectus. *Journal of Applied*
472 *Ecology*, 48(3), 758–767. doi:10.1111/j.1365-2664.2011.01993.x

473 Bradbury, J. W., & Vehrencamp, S. L. (1998). *Principles of animal communication* (1st ed.).
474 Sunderland, MA: Sinauer Associates.

475 Budka, M., & Osiejuk, T. S. (2013). Formant frequencies are acoustic cues to caller
476 discrimination and are a weak indicator of the body size of corncrake males. *Ethology*,
477 119(11), 960–969. doi:10.1111/eth.12141

478 Couchoux, C., & Dabelsteen, T. (2015). Acoustic cues to individual identity in the rattle calls
479 of common blackbirds: a potential for individual recognition through multi-syllabic
480 vocalisations emitted in both territorial and alarm contexts. *Behaviour*, 152(1), 57–82.
481 doi:10.1163/1568539X-00003232

482 Crowley, P. H., Provencher, L., Sloane, S., Dugatkin, L. A., Spohn, B., Rogers, L., & Alfieri,
483 M. (1996). Evolving cooperation: the role of individual recognition. *Biosystems*, 37(1),
484 49–66. doi:10.1016/0303-2647(95)01546-9

485 Elie, J. E., & Theunissen, F. E. (2018). Zebra finches identify individuals using vocal
486 signatures unique to each call type. *Nature Communications*, 9(1), 4026.
487 doi:10.1038/s41467-018-06394-9

488 Godard, R. (1991). Long-term memory of individual neighbors in a migratory songbird.
489 *Nature*, 350(6315), 228–229.

- 490 Insley, S. J., Phillips, A., & Charrier, I. (2003). A review of social recognition in pinnipeds.
491 *Aquatic Mammals*, 29, 181–201.
- 492 Korkmaz, S., Goksuluk, D., & Zararsiz, G. (2014). MVN: An R package for assessing
493 multivariate normality. *The R Journal*, 6(2), 151–162.
- 494 Lengagne, T., Lauga, J., & Jouventin, P. (1997). A method of independent time and frequency
495 decomposition of bioacoustic signals: inter-individual recognition in four species of
496 penguins. *Comptes Rendus De L Academie Des Sciences Serie Iii-Sciences De La Vie-
497 Life Sciences*, 320, 885–891. doi:10.1016/s0764-4469(97)80873-6
- 498 Linhart, P., & Šálek, M. (2017). The assessment of biases in the acoustic discrimination of
499 individuals. *PLOS ONE*, 12(5), e0177206. doi:10.1371/journal.pone.0177206
- 500 Mathevon, N., Koralek, A., Weldele, M., Glickman, S. E., & Theunissen, F. E. (2010). What
501 the hyena's laugh tells: Sex, age, dominance and individual signature in the giggling
502 call of *Crocuta crocuta*. *BMC Ecology*, 10, 9-Article No.: 9. doi:10.1186/1472-6785-
503 10-9
- 504 Mielke, A., & Zuberbuehler, K. (2013). A method for automated individual, species and call
505 type recognition in free-ranging animals. *Animal Behaviour*, 86(2), 475–482.
506 doi:10.1016/j.anbehav.2013.04.017
- 507 Partan, S., & Marler, P. (1999). Communication goes multimodal. *Science*, 283(5406), 1272–
508 1273. doi:10.1126/science.283.5406.1272
- 509 Partan, S. R. (2013). Ten unanswered questions in multimodal communication. *Behavioral
510 Ecology and Sociobiology*, 67(9), 1523–1539. doi:10.1007/s00265-013-1565-y
- 511 Pollard, K. A., & Blumstein, D. T. (2011). Social group size predicts the evolution of
512 individuality. *Current Biology*, 21(5), 413–417. doi:10.1016/j.cub.2011.01.051

513 R Core Team. (2012). *R: A Language and environment for statistical computing*. Vienna,
514 Austria: R Foundation for Statistical Computing. Retrieved from [http://www.R-](http://www.R-project.org/)
515 [project.org/](http://www.R-project.org/)

516 Robisson, P., Aubin, T., & Bremond, J. (1993). Individuality in the voice of the emperor
517 penguin *Aptenodytes-Forsteri* - Adaptation to a noisy environment. *Ethology*, *94*(4),
518 279–290.

519 Searby, A., & Jouventin, P. (2004). How to measure information carried by a modulated vocal
520 signature? *Journal of the Acoustical Society of America*, *116*, 3192–3198.
521 doi:10.1121/1.1775271

522 Sheehan, M. J., & Nachman, M. W. (2014). Morphological and population genomic evidence
523 that human faces have evolved to signal individual identity. *Nature Communications*,
524 *5*, 4800. doi:10.1038/ncomms5800

525 Stowell, D., Petrusková, T., Šálek, M., & Linhart, P. (2019). Automatic acoustic identification
526 of individuals in multiple species: improving identification across recording
527 conditions. *Journal of The Royal Society Interface*, *16*(153), 20180940.
528 doi:10.1098/rsif.2018.0940

529 Syrová, M., Policht, R., Linhart, P., & Špinko, M. (2017). Ontogeny of individual and litter
530 identity signaling in grunts of piglets. *The Journal of the Acoustical Society of*
531 *America*, *142*(5), 3116–3121. doi:10.1121/1.5010330

532 Terry, A. M. R., & McGregor, P. K. (2002). Census and monitoring based on individually
533 identifiable vocalizations: the role of neural networks. *Animal Conservation*, *5*, 103–
534 111. doi:10.1017/s1367943002002147

535 Tibbetts, E. A. (2004). Complex social behaviour can select for variability in visual features: a
536 case study in *Polistes* wasps. *Proceedings of the Royal Society of London B:*
537 *Biological Sciences*, *271*(1551), 1955–1960. doi:10.1098/rspb.2004.2784

538 Tibbetts, E. A., & Dale, J. (2007). Individual recognition: it is good to be different. *Trends in*
539 *Ecology & Evolution*, 22(10), 529–537. doi:10.1016/j.tree.2007.09.001

540 Trunk, G. V. (1979). A problem of dimensionality: a simple example. *IEEE Transactions on*
541 *Pattern Analysis and Machine Intelligence*, 1(3), 306–307.

542 Venables, W. N., & Ripley, B. D. (2002). *Modern Applied Statistics with S* (Fourth). New
543 York: Springer. Retrieved from <http://www.stats.ox.ac.uk/pub/MASS4>

544 Wiley, R. H. (2013). Specificity and multiplicity in the recognition of individuals:
545 implications for the evolution of social behaviour. *Biological Reviews*, 88(1), 179–
546 195. doi:10.1111/j.1469-185X.2012.00246.x

547 Wilkinson, G. S. (1984). Reciprocal food sharing in the vampire bat. *Nature*, 308(5955), 181–
548 184. doi:10.1038/308181a0

549