

This is a repository copy of Towards a general framework for predicting threat status of data-deficient species from phylogenetic, spatial and environmental information.

White Rose Research Online URL for this paper: http://eprints.whiterose.ac.uk/88375/

Version: Accepted Version

Article:

Jetz, W. and Freckleton, R.P. (2015) Towards a general framework for predicting threat status of data-deficient species from phylogenetic, spatial and environmental information. Philosophical Transactions B: Biological Sciences, 370 (1662). ISSN 0962-8436

https://doi.org/10.1098/rstb.2014.0016

Reuse

Unless indicated otherwise, fulltext items are protected by copyright with all rights reserved. The copyright exception in section 29 of the Copyright, Designs and Patents Act 1988 allows the making of a single copy solely for the purpose of non-commercial research or private study within the limits of fair dealing. The publisher or other rights-holder may allow further reproduction and re-use of this version - refer to the White Rose Research Online record for this item. Where records identify the publisher as the copyright holder, users can verify any specific terms of use on the publisher's website.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



1	Toward a general framework for predicting threat status of data-
2	deficient species from phylogenetic, spatial and environmental
3	information
4	
5	
6	Walter Jetz ¹ * & Robert P. Freckleton ² *
7	
8	¹ Department of Ecology and Evolutionary Biology, Yale University, PO Box 208106,
9	New Haven CT 06520-8106, USA.
10	
11	² Department of Animal & Plant Sciences, University of Sheffield, Sheffield S10 2TN,
12	United Kingdom
13	
14	Running head: Predicting species threat status
15	
16	
17	
18	
19	
20	* Authors contributed equally

21 Abstract

In taxon-wide assessments of threat status many species remain not included due to lack of data. Here we present a novel spatial-phylogenetic statistical framework that uses a small set of readily available or derivable characteristics, including phylogenetically imputed body mass and remotely-sensed human encroachment, to provide initial baseline predictions of threat status for data-deficient species. Applied to assessed mammal species worldwide the approach effectively identifies threatened species and predicts the geographic variation in threat. For the 483 data-deficient species the models predict highly elevated threat, with 69% 'at-risk' species in this set, compared to 22% among assessed species. This results in 331 additional potentially threatened mammals, with elevated conservation importance in rodents, bats and shrews, and countries like Colombia, Sulawesi, and the Philippines. These findings demonstrate the future potential for combining phylogenies and remotely sensed data with species distributions to identify species and regions of conservation concern.

Introduction

Human activities continue to cause the loss of many species together with the function and services they provide [1]. In the face of these mounting threats and limited resources to conserve species [2], tools are required to identify those of greatest conservation concern. Global IUCN Red List assessments [3] have provided important knowledge about the state of biodiversity and have helped to identify priority species and regions for conservation [4-7]. On this basis approximately 20% of mammal, bird and amphibian species are currently identified as threatened [3]. In order to minimize potential biases in perceived patterns of biodiversity threat, species should be assessed comprehensively or at least representatively. In addition to undiscovered species [4, 8], species with too little information for threat categorization ('data-deficient species') are thus a major concern. The IUCN assessment process relies on available field-based knowledge of e.g. population size, rate of decline, and range size of each species to assigned threat status [9-12]. Paucity of data, e.g. due to financial or logistical limitations for field studies, makes complete assessments impossible for some species, with little prospects for change in the near future. The number of species lacking data may be substantial, with e.g. 2,436 of 11,806 recognized mammal and amphibians species classified as 'data-deficient' in 2011 [3], including 834 extant mammals. The potential for data-deficient species to change absolute threat levels of taxa has been acknowledged [4, 7]. In the absence of better knowledge, a risk-averse approach may be to simply assume that all data-deficient species are threatened. But given the sheer number of data-deficient species the implications for conservation prioritization may be substantial and carry a high cost if large numbers are in fact not threatened. At the other extreme data-deficient species may

> 59 be no more threatened than assessed species, for example as appears to be the case with 60 data-deficient birds [13]. Model-based initial baseline threat predictions for data-deficient 61 species and a general framework to provide them for groups assessed in the future would 62 thus hold multiple benefits for conservation practice.

The relative paucity of data on ecology and threats of many species stands in stark contrast to our rapidly growing detailed knowledge about species' phylogenetic relationships and geographic distributions. Technology now allows cost-effective and rapid generation of phylogenies for thousands of species. While often remaining coarse in grain [14, 15] or even limited to type specimen and thus for some species a main reason for data-deficient status, more facetted geographic distribution information is increasingly becoming available for many species [15; see also mol.org]. Distribution data permits two types of inference about potential threat status. First, statistical models can quantitatively capture the association between range size and threat status in assessed species [16] and then can be applied to data-deficient species [17-19] to account for this risk component. Second, geographic range information can be intersected with environmental layers that inform about broad-scale environmental niches and associated life history signals (e.g. on fecundity, generation time) related to threat status [16, 20]. And, more directly, remotely-sensed layers of land-cover can provide coarse estimates of potential habitat loss due to human encroachment. Information of this kind has recently been shown to successfully predict threat status in birds [20] and mammals [21]. Modern statistical tools allow the development of models of correlates of current threat levels that incorporate both phylogenetic and spatial data [17, 22-27].

In previous work, modeled threat predictions for data-deficient species have been made without environmental or phylogenetic information [19], or without habitat encroachment information and using eigenvectors [17, 18] which are highly constrained in their ability to appropriately represent both phylogenetic and spatial signals [28, 29]. A general framework that readily capitalizes on the ever increasing availability of species distribution and remote sensing data, and rigorously incorporates phylogenetic and geographic information is thus still missing. In this study we build on our earlier work linking spatial and phylogenetic models [27] and predicting threat data with GIS-derived habitat information [20] to develop such a framework. We demonstrate the approach applied to mammals by parameterizing models of threat status based on readily available variables capturing key aspects of life history, rarity and range loss (body mass, geographic range size, human encroachment on species' ranges) together with spatial and phylogenetic dependency for 3,703 mammal species across 16 orders with sufficient information to be assessed by the Red List. We then apply these models to 483 species classified as data-deficient species. We show, that the presented framework may offer a cost-effective way for initial baseline threat evaluation of many understudied (and potentially at-risk) species.

101 Methods

Data. We analyzed data on 4,186 terrestrial mammal species from 16 orders in the IUCN
Red List [30] that could also be placed in the mammalian super tree phylogeny [31] (with
recent updates). Of these, 3703 species had been assessed (with 812 deemed threatened,

105	i.e. categories "Vulnerable", "Endangered" and "Critically Endangered") and 483
106	recognized but not assessed (category "Data-Deficient") by IUCN. We gathered
107	information on mammal body masses from [32] and from F. Smith (pers. comm.). One
108	order (Perissodactyla) contained no data-deficient species. We selected native and
109	reintroduced resident and breeding ranges that were extant or probably extant from the
110	IUCN expert range maps [30] which we extracted over a 110x110km grid in Equal Area
111	Cylindrical projection. We overlaid each species range map with information on
112	transformed habitats owing to anthropogenic activities. Specifically, we estimated
113	'Encroachment' as the proportion of expert range transformed by past human activities
114	(i.e. cultivated or managed, mosaics, including cropland and urban areas) according to the
115	Global LandCover 2000 land-cover classification [33]. At 1km native resolution this
116	information is collected at much finer scale than expert range maps and analysis grid
117	[14], but used as range summary measure it offers a concrete first-order estimate of
118	overall range encroachment, and has recently been shown to be a strong correlate of
119	expert-assessed IUCN threat status in birds [20]. We note that other high-resolution
120	global land cover classifications exist and that all suffer from remaining classification
121	errors [34]. As additional metric we also calculated the average Human Influence Index
	errors [5 1]. The additional mourie we also carealated the average righting minuence match

124 Summary of approach

To summarize our approach, we first imputed the body masses of species for which data are missing and then used Generalized Linear Models that include phylogenetic and spatial dependence to predict IUCN status. We account for statistical uncertainty in our

2
0
ა
4
2 3 4 5 6 7 8 9 10 1 12 3 4 15 16 7 8 9 20 1 2 2 2 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 3
5
6
0
7
8
a
10
4.4
11
12
10
13
1/
14
15
16
10
17
18
10
19
20
01
21
22
~~
23
24
<u> </u>
25
06
20
27
28
20
29
30
~
31
30
02
33
04
34
35
~~
36
37
07
38
20
29
40
41
42
43
44
45
46
47
11
48
49
49
50
51
эı
52
50
53
54
55
56
57
58
hu
59 60

estimates of body mass by using Multiple Imputation. In order to incorporate uncertainty
in our overall predictions, we express the model outputs as threat probabilities; i.e. given
the predictions of the model and the statistical uncertainty in these, what is the probability
that each species is threatened (i.e. IUCN categories Vulnerable, Endangered or Critically
Endangered) or not?

133

134 Statistical modelling framework

135 The starting point for our analyses is a linear statistical model relating the values of a trait 136 of interest to a set of predictors [24, 26]. The errors are assumed to have a multivariate 137 normal distribution with mean 0 and a variance-covariance matrix that is defined by the phylogeny [23, 24, 26] and spatial distances [27]. Predictions from our models were 138 139 generated by using the fitted parameter values together with the degree of phylogenetic 140 and spatial similarity of species using the approach described in [26]. Our predictions 141 therefore account for the phylogenetic /spatial structure in the data, i.e. they have the 142 property that closely related, or species that live in the same place, should be similar to 143 each other. We calculated variances for predicted values using the formulae in [24]. These variances are used to calculate the variance in estimates of body mass and IUCN 144 145 status (below).

146

147 **Phylogenetic and spatial models for trait covariances**

148 We use the generalized least squares (GLS) approach described in Freckleton & Jetz [27]

149 to account for both spatial and phylogenetic effects. A parameter ϕ is included in the

2
3 4 5
4
5
6
7
7 8 9 10 11 12 13 14 15 16 17
8
9
10
11
12
13
1/
15
10
16
1/
18
19
20
21
22
<u>, , , , , , , , , , , , , , , , , , , </u>
17 18 19 20 21 22 23 24 25 27 28 29 31 23 34 35 37 38 90
24
25
26
27
28
29
30
21
20
3Z
33
34
35
36
37
38
39
40
40 41
42
43
44
45
46
47
48
49
49 50
52
53
54
55
56
57
58
59
60

1 S

150	model to account for the influence of space. According to this model, of the total
151	variance, a proportion ϕ is attributed to spatial variance, $(1 - \phi)$ is due to the non-spatial
152	component. We also used the λ transformation suggested by Pagel [22, 37]. In the
153	context of modelling spatial and phylogenetic effects simultaneously, the λ statistic
154	allows us to include trait variation independent of both phylogeny and space in our
155	analysis: a proportion $\gamma = (1 - \phi) (1 - \lambda)$ of the trait variation is independent of phylogeny
156	or space [27]. This approach is akin to including a 'nugget' in a spatial model [38]. We
157	estimated ϕ and λ by maximum likelihood [39].
158	
159	The spatial matrix was calculated and tested using the approach described in Freckleton
160	& Jetz [27]. The spatial matrix reduces the spatial configuration of the data to a series of
161	pairwise distances that measure the distance between each species. Following Freckleton

162 & Jetz [27] we did this by calculating the distances between the centroids of the ranges of 163 each pair of species. The assumption is, therefore, that the variance between species' 164 traits grows linearly with spatial distances. As we showed before, this assumption can be 165 tested graphically and in the analyses reported here, as well as in Freckleton & Jetz [27], 166 this assumption was found to be adequate. Following Freckleton & Jetz [27], in order to 167 aid interpretation of the model we define λ' as the relative contribution of phylogeny (λ' $=\lambda (1 - \phi)$) once the effects of space have been accounted for . This parameterisation 168 169 allows a simple interpretation of the joint estimates of ϕ and λ because, as shown in

170 Freckleton & Jetz [27], the sum of γ , λ' , and ϕ is always 1. These parameters can be

171 interpreted as the individual proportional contributions to variance of the different

 172 variance components.

174 Imputation of mammalian body mass

We used estimates of mammalian body masses of 3462 species in the 16 analyzed orders to predict the values for the 723 species without body mass data. For each order we used the GLS approach described above to predict body mass based on the species with body mass data along with phylogenetic and spatial information. We conducted this analysis at the level of orders as previous analysis has shown that the Brownian model, modified to allow for varying degrees of phylogenetic dependence, provides an adequate description of body mass variation within orders [40]. Body mass was log-transformed prior to analysis.

For species missing body mass we used the predicted values predicted as estimates of log mass in the modeling of IUCN threat status. A problem with using single imputation of this sort is that although parameter estimates should be unbiased [41], there is a possibility of under-estimation of variances for parameters using this method. We therefore conducted significance tests for our models using multiple imputation. For this we calculated for all species lacking body mass data predicted values using the above GLS model, along with a variance for each prediction (using the method in [24], see above). These estimates formed the basis for the multiple imputations [for further background on the method see 42, 43, 44; for specific implementation here see also Nakagawa & Freckleton 2008]. We used 10 imputations, and the statistical tests reported

in Table S1 are the outcome of this analysis. We found that in practice the variance
across the imputations was very small indeed so that this step was not vital in this case,
although this need not always be true.

> In order to evaluate the accuracy of the predictions of body mass we used a simple randomization. Estimates of λ and ϕ from the best fitting model for each order were used to construct a variance-covariance matrix. This variance matrix formed the basis for generating randomized multivariate normally distributed data (using the rmvnorm in the R mythorm package). Species originally missing data were then removed and their values imputed. The correlation of these imputed values with the true values was then calculated. Note that because this analysis is conducted on randomly generated data, this is different from a cross-validation which is based on removal of data from the original data and would not normally be conducted using single-species removals. This was repeated 1000 times per missing species per order. The results of this analysis are summarized in Table S2.

210 Application to IUCN categories

The IUCN categories were treated as a five point ordinal scale ranging from "Least Concern", 1, to "Critically Endangered", 5. Although the response variable is a discrete ordinal variable, the models described observed threat levels well, offering explanatory power equal to, or better than that found in previous studies (Table S1, Fig. S1). This same approach has been taken in other recent analyses of threat status [45]. We compared

1	
2	
3	
4	
5	
6	
7	
7	
8 9	
9	
10	
11 12	
13	
14	
14 15	
15	
16	
1/	
18	
19	
20	
21	
 19 20 21 22 23 24 25 26 27 28 29 30 	
22	
23	
24	
25	
26	
27	
28	
20	
23	
30	
31	
32	
33	
34	
35	
34 35 36 37 38 39	
27	
37	
38	
40	
41	
42	
43	
44	
45	
46	
47	
48	
49	
50	
51	
52	
54	
55	
56	
57	
58	
59	
23	

60

216	our results with those of generalized linear models in which responses are treated as
217	multinomial or ordered logistic responses, which yielded very similar results, but are
218	unable to address the spatial and phylogenetic covariance (see Figs S2, S3 and below).
219	The main problem in generating an output from the model is that a fitted / predicted value
220	is a point estimate and does not account for the statistical uncertainty in our estimates. To
221	incorporate uncertainty, after the analysis we converted our predictions of IUCN status
222	into probabilities of threat. Previous analyses have taken a similar approach in the
223	analysis of threat status, but instead converting the threat to a binary variable before the
224	regression analysis [17]. This has the disadvantage that information on the ordinal nature
225	of the IUCN scale is ignored. Our analysis, however, retained the continuous information
226	in the model fitting: for example we account for the fact that a species classified as
227	category 5 (Critically Endangered) is more at risk than a species in category 3
228	(Vulnerable).

229

To produce these threat probabilities we calculated the probability that each species was
threatened or not from the predictions of IUCN status. This was simply done by
calculating:

233
$$P_i^{threatened} = Z\left(\frac{y_i^{pred}-2.5}{\sigma_i}\right)$$
 (1)

where Z() is the cumulative z (standardized normal) distribution, y^{pred} is the predicted value and σ is its standard deviation. This is the probability that the predicted value of species *i* is *greater* than 2.5 (see also [17, 20]. The choice of threshold in equation (1) is dependent on the interpretation of the categories and how these relate to continuous

model predictions. With eqn. (1), a species with an IUCN status predicted to be 2.5 (i.e. in between "near threatened" and "vulnerable") will have a threat probability of 0.5. We repeated the analysis using a threshold of 2 which yielded a visually clearer discrimination between the higher IUCN categories, but did essentially not affect the results of Fig. 1(Fig. S4), because the probabilities are simply rescaled such that the mean probability is 0.5 at a predicted value of 2 rather than 2.5. The results in Fig 2 are also extremely similar (Fig. S5), because the estimates of the proportions of species to be threatened or not are set by a threshold estimated from the data by receiver operator characteristic analysis (below). Thus, results were broadly invariant to the choice of threshold in equation (1).

We used the full models in Table S1 for making predictions and did not attempt model reduction. There were several reasons for this. First, model reduction by elimination of variables (e.g. based on statistical significance) has undesirable consequences, such as degenerate sampling distributions and model selection bias [46]. Second, examination of the coefficients for the predictors indicated that, independent of statistical significance, the directions of effects were usually quite consistent between orders. For example 15 out of 16 coefficients for the effect of body mass are positive even if all are not statistically significant (Table S1); 12 of 16 coefficients for the encroachment variable are positive (Table S1). Finally, we checked predictions with and without the least significant variables and confirmed that the R² values were not unduly inflated and giving a false impression of good fit. In order to test the predictive ability of the threat probabilities we assessed how well the fitted threat probabilities predicted for assessed species were able

2
3
1
4 E
5
6
7
8
9
10
10
11
12
13
14
15
16
17
10
18
19
20
21
2 3 4 5 6 7 8 9 10 11 2 3 14 15 16 7 8 19 20 12 23 24 25 6 7 8 9 10 11 23 14 15 16 17 8 19 20 12 23 24 25 6 7 8 29 30 13 23 34 5 6 37 8 39 40
23
20
24
25
26
27
28
20
20
30
31
32
33
34
35
36
07
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

261	to distinguish threatened from non-threatened species using the Area Under the Curve
262	(AUC) in the Receiver Operating Characteristic (ROC) curve [47]. AUC varies between
263	0.5, which indicates that the predictions are no better than random, and 1, which is
264	perfect agreement between observed and predicted. As threshold for assigning
265	probabilities into binary categories of threatened and non-threatened, we used the value at
266	which sensitivity equaled specificity in a given order.

267

268 Model Approach and Limitations

269 The methodology we have used is based on currently available tools and will be 270 improved by future developments that include techniques such as logistic and 271 multinomial generalized linear mixed models that could account for phylogenetic and 272 spatial dependence and would enable to better model the discrete ordinal state variable 273 [48, 49]. However such tools require very large datasets: logistic regression requires large 274 amounts of data because binary observations contain relatively little information. 275 Multinomial or ordered responses are an extension of logistic regression and as the 276 number of states increases the data requirements increase. Given this, the approach taken 277 here to treat the data as continuous is unlikely to seriously compromise the results (see 278 also supplementary results). Moreover existing methods for such responses do not 279 combine spatial and phylogenetic signals, and can be very difficult to implement and 280 tune. In the future, faster methods for fitting phylogenetic models are under development 281 and these should facilitate further methodological advances [50]. We have assumed that 282 the variance scales linearly with both phylogenetic and geographic distances. This is 283 supported by diagnostics (e.g. see [27] for a worked examples). The assumption of

2	
2	
3 4 5 6 7 8 9 10 11 23 14 15 16 17 18 0	
4	
5	
6	
7	
8	
9	
10	
11	
12	
12	
10	
14	
15	
16	
17	
18	
19	
20	
21	
20 21 22 23 24 25 26 27 28 29 30 31 32 33 4 35 36 37 38 39	
22	
20	
24	
25	
26	
27	
28	
29	
30	
31	
32	
33	
24	
34	
35	
36	
37	
38	
39	
40	
41	
42	
43	
44	
44 45	
46	
47	
48	
49	
50	
51	
52	
53	
54	
55	
55 56	
57	
58	
59	
60	

linearity is not terribly critical so long as variance increases with distance. In previous
work we suggested how the assumption could be varied (Table 1 in [27]). However it
should be noted that nonlinear transformations of variance matrices are potentially
difficult to work with. For example we have recently shown that a commonly used
transformation (the Ornstein Uhlenbeck) is severely biased under most circumstances for
even large datasets (Thomas et al. in review).

290 The models we developed are strongly dependent on range size as a predictor of IUCN

status, which reflects the importance of range size in the formal assessment process. It is

important to note that our predictive models are not aimed at *testing* the relevance of this

293 variable (which would require variable elimination to avoid circularity), but to use this

294 formally recognized association for prediction. In other words, we use A (assessed

295 species) modeled by B (novel framework and independent variables) to predict C (not yet

assessed species), not to make inference about A.

297

298 **Results**

299 Assessed species

For the 16 mammal orders analyzed the threat probabilities (whether a species is nonthreatened or threatened) predicted by the models successfully explain observed variation in assessed threat score (R^2 values were typically ~40% or greater; Table S1) and effectively predict species threat category (Fig. 1A). Range size and body mass were generally strong correlates of threat status, with smaller ranging and larger species typically being subject to greater threat (Table S1). Given the inherent role of range size

1	
2	
3	
1	
3 4 5 6 7 8	
5	
0	
1	
8	
9 10	
10	
11	
12	
13	
14	
15	
16	
17	
$\begin{array}{c} 11\\ 12\\ 13\\ 14\\ 15\\ 16\\ 17\\ 18\\ 92\\ 12\\ 23\\ 24\\ 25\\ 27\\ 28\\ 29\\ 31\\ 32\\ 33\\ 35\\ 37\\ 33\\ 35\\ 37\\ 38\\ 37\\ 38\\ 37\\ 38\\ 37\\ 38\\ 38\\ 38\\ 38\\ 38\\ 38\\ 38\\ 38\\ 38\\ 38$	
19	
20	
21	
22	
22 22	
23	
24	
25	
26	
27	
28	
29	
30	
31	
32	
33	
34	
35	
36	
37	
38	
39	
40	
40 41	
42	
42 43	
44	
45	
46	
47	
48	
49	
50	
51	
52	
53	
54	
55	
56	
57	
58	
59	
60	
$\mathbf{u}\mathbf{u}$	

306	in the IUCN assessment process [11] these associations are not altogether unexpected and
307	confirm previous findings [19, 20, 51, 52]. Less consistently than recently observed in all
308	terrestrial birds [20], land-cover encroachment and human influence measures are
309	strongly positively correlated with IUCN threat category in several orders. This
310	contributes to the overall predictive ability of the models and confirms the relevance of
311	such variables for threat predictions (Table S1).
312	
313	In addition to the strong phylogenetic dependence of body mass (Table S2), 9 of the
314	orders showed phylogenetic or spatial dependence in the residuals of the models for
315	IUCN threat. The degree of net phylogenetic signal in the residuals of the final models is
316	generally low, with the phylogenetic effect estimated as zero for seven and very low (0.1
317	or less) for five orders. Notably higher estimates are obtained for primates (0.66). Six
318	orders showed strong spatial signals, with estimates of the spatial coefficient, ϕ , as high
319	as 0.6-1.0 (Table S1). The threat probabilities resulting from our models are the
320	probabilities that each species is in one of the threatened states rather than not threatened,
321	given the mean and variance predicted by the model (see methods). The Receiver
322	Operating Characteristic plot (Fig 1B) indicates a very strong discrimination of
323	threatened from non-threatened species with an AUC of 0.90 for the whole dataset and a
324	median of 0.91 for all orders. These were associated with high degrees of sensitivity and
325	specificity (typically ca. 0.8 - 0.9, Table S1). Predicted threat probabilities are remarkably
326	effective in delimiting threat status, as especially illustrated by the most and least
327	threatened IUCN classes: only 4% of species assessed to be of 'least concern' were
220	

328 predicted to have a threat probability of 0.5 or greater (Fig. 1C; see Fig. S1 for order-

level plots) and only 11% of species assessed to be 'critically endangered' were assigned
a threat probability lower than 0.5 (Fig. 1C). Across all threat categories, 61% of species
assessed as being under some degree of threat had estimated threat probabilities greater
than 0.6 and (with nearly 31% greater than 0.8) (Fig. 1C). Overall, our predicted threat
probabilities are a strong discriminator of threat status with particularly high values
(>0.8) extremely unlikely for species that are not actually threatened.

The relative richness of species assessed as being threatened is geographically very uneven (Fig. 3*A*). Applied to assessed species, our model predicts this observed pattern very well (Fig. 3B). Overall, however, there is a strong association between the predicted average probability or predicted proportion of species threatened and the observed proportion of species assessed threatened (r = 0.74 and r = 0.68, respectively; Fig. 4A, B) as well as, expectedly, between predicted and observed threatened richness (r = 0.82, Fig. 4C). This suggests that our models successfully capture the biogeography of assessed threatened species.

Data-deficient species

Data-deficient species are predicted to be substantially more likely to be threatened than assessed species (Fig. 1D), with an average predicted threat probability of 0.40 compared to 0.21 in assessed species. For data-deficient species 28% of threat probability estimates were greater than 0.6, whereas for assessed species it was only 11%. Overall threat probabilities were higher for data-deficient species in 10 orders, and statistically

significantly so for 7 orders (Fig. 2). Classifying data deficient species into binary threat categories using a standardized order-specific threshold (the value at which sensitivity equals specificity) results in a total of 331 of 483 species predicted threatened, i.e. 69% of species threatened compared to 29% among assessed species. This difference is repeated among almost all orders, with a total of 298 potentially threatened data-deficient species identified among the Chiroptera, Rodentia and Eulipotyphla alone (for species-level results see Table S3).

Geographically, data-deficient species are predicted to exhibit substantially higher average probabilities and proportions of species threatened (grid cell assemblage values of 0.12 and 0.17, respectively) than assessed species (0.06 and 0.06, respectively). At the grid cell level, the predicted average threat probabilities or proportion of data-deficient species shows barely any relationship with the proportion of species assessed to be threatened (Spearman rank correlations: r = 0.30 and r = 0.29, respectively; Figs. 4D,E). Equally, the richness of data-deficient species predicted to be threatened is only weakly correlated with that of species assessed to be threatened (r = 0.30, Fig. 4F). The discordance in geographic 'hotspots' of predicted assessed and data-deficient threat is apparent when comparing the maps of their predicted threat probabilities and species richness in Figure 3. Threat levels predicted for data-deficient species substantially exceed those of assessed species in many locations (note different color scales). Data-deficient species hold much higher predicted threat levels than assessed species in Colombia and Central America, Southern South America, and parts of Southeast Asia. In terms of species richness (Fig 3F, 4F), data-deficient species are predicted to strongly

increase the number of at-risk species in Sumatra, New Guinea, Colombia, and especially
Sulawesi, where in grid cell 10 likely threatened mammal species add to the known 16.
This suggests that these regions are even more important for conservation than previous
global conservation prioritization analyses may have suggested [53, 54]. In contrast, datadeficient species predicted *not* to be threatened occur both outside (e.g. Southern South
America) and inside (e.g. Borneo, Central and West African forests) some main areas of
known (assessed) high prevalence of threatened species (Fig 3*D*, *E*).

383 Discussion

In this study we have shown that data-deficient species are much more likely to be under threat than those that have already been assessed and that the geographic distribution of data-deficient species that are likely threatened is different to that of assessed threatened species. This may have important implications for global mammal conservation strategies [55]. According to our analysis it is extremely likely that well over three hundred additional mammal species (69% of those data-deficient) are threatened, many of them potentially severely so. This is over an order of magnitude more than suggested by Davidson et al [19] which identified 28 data-deficient mammal species as potentially threatened, but did not use environmental, spatial or phylogenetic information. Using eigenvectors, no encroachment data and model validation with only bats, Jones and Safi [18] estimated 35% of 481 data-deficient mammal species to be potentially threatened. Our statistically more robust approach [28, 29] that additionally uses remotely sensed encroachment information thus suggests much greater levels of threat in data-deficient

species than previously thought. The relatively low degree of phylogenetic signal of
IUCN status we found here contrasts with previous related results in carnivores [17]. This
difference has two sources: firstly from the inclusion of species' body masses in our
analyses, and secondly from the inclusion of spatial effects, which also has phylogenetic
signal. In particular mean mass is both strongly phylogenetically determined in all orders
and strongly related to IUCN status (Table S1). Accounting for body mass thus decreases
the detectable phylogenetic signal.

Our findings suggest that data-deficient species cannot be ignored in conservation threat assessments and in interpretation of threat status for policy setting. In mammals data-deficient species are clearly more likely than non data-deficient species to be under significant threat. The association between threat status and data deficiency arises, because narrow-ranged (and thus often scarce), large-bodied (that thus often low-density, long generation time) species, are also very likely to be those for which little data exist (Table S1, [see also 45]). There are notable exceptions: for instance, the threat probability of data-deficient primates is no higher than that of those that have been assessed, likely reflecting the relatively higher research effort directed at primates in the past. In contrast, rodents are much more difficult to study (they are small, live in inaccessible habitat and are frequently nocturnal) and for them over half of data-deficient species are predicted to be threatened whereas only 16% have been assessed as threatened (Fig 2). Our findings contrast with recent results for birds, where just 0.6% of species are data-deficient and where species that were recently moved from this category were found to be less threatened than non-data-deficient species [13]. However, these only recently assessed

bird species are likely not a representative sample of data-deficient species as whole, as data-deficient species assessed first will likely be ones that are more easily studied (and thus face different, potentially lesser threats) than those assessed last. The statistical results gathered from all species may offer more reliable guidance.

Our general aim was to demonstrate how readily-available information can be used to make initial predictions about the likely conservation status of species for which a formal assessment has not yet been possible. If a similar proportion of data-deficient yet threatened mammal species (69%) was to be found amongst data-deficient amphibians (1,600 out of 6,312 species are data-deficient, [3]), it would represent a very large increase of amphibian species at-risk. Such a scenario would add many new species to the threated categories in the Red List with strong potential consequences for geographic conservation prioritization. The transferability of mammal-based estimates to other taxa is of course unclear, but this realization highlights the importance of expanding assessment work and seizing the increasing opportunities for rigorous statistical inference of threat status.

The strong importance of select life history traits and range size for predicting threat
status has previously been illustrated [16]. Recently the complementary power of
remotely-sensed measures of human land encroachment to predict threat status has also
been demonstrated for birds [20]. Combined with an increasingly thorough understanding
of the spatial context of species [15] and ever- improving data on the phylogenetic signal,
a general predictive framework is emerging that may be instrumental for statistically

1
2
3
4
т Б
5
6
7
8
9
10
11
10
10
13
14
15
16
17
18
19
3 4 5 6 7 8 9 10 11 2 3 14 5 16 7 8 9 10 11 2 3 14 5 16 7 8 9 10 11 2 3 14 5 16 7 8 9 10 11 2 3 14 5 16 7 8 9 2 2 2 3 2 4 2 5 2 7 8 9 3 0 1 3 2 3 3 4 5 5 6 7 8 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
20
21
22
23
24
25
26
27
28
20
29
30
31
32
33
34
35
36
07
3/
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
51 52 53
53 54
04 55
55
56
57
58
59

59 60 443 assessing the thousands of species for which an individual evaluation is time- or cost-444 prohibitive. By identifying already assessed species with highly over- or under-predicted 445 threat status for further scrutiny, it may also someday help improve the Red Listing 446 process which is not without human error. Clearly, the presented framework is no silver 447 bullet to replace the need for expert assessment based on field ecological data. We expect 448 that assessment data for at least 50% of species, depending on representativeness, is 449 needed to provide reasonably reliable threat predictions, but this will vary by group and 450 likely often be higher. But this does potentially free up resources and lower completion 451 thresholds [56, 57] that would benefit the assessment of neglected taxa such as 452 invertebrates and plants. More generally, a complementary approach to traditional expert-453 based assessment may emerge that combines available phylogenetic/biological data with 454 improved species distribution knowledge linked to a remotely-sensed monitoring of land 455 cover [15] – all facilitating a dynamic and continuous baseline assessment of the state of 456 species. 457 458

459 Acknowledgements

We thank Arne Mooers, Tien Ming Lee, and members of the Jetz Lab for feedback on the
manuscript. RPF was funded by a Royal Society University Research Fellowship for this
project. WJ acknowledges support from NSF grants DBI 0960550 and DEB 1026764 and
NASA Biodiversity Grant NNX11AP72G. We are grateful to Felisa Smith and the
NESCent body size group for sharing mammal body mass data.

465	
466	
F	References
1	 Pereira H.M., Leadley P.W., Proença V., Alkemade R., Scharlemann J.P.W., Fernandez-Manjarrés J.F., Araújo M.B., Balvanera P., Biggs R., Cheung W.W.L., et al. 2010 Scenarios for Global Biodiversity in the 21st Century. <i>Science</i> 330(6010), 1496-1501. (doi:10.1126/science.1196624).
2	. Wilson K.A., McBride M.F., Bode M., Possingham H.P. 2006 Prioritizing global conservation efforts. <i>Nature</i> 440 , 337-340.
3	 IUCN. 2011 IUCN Red List of Threatened Species 2011.1 - <u>http://www.iucnredlist.org</u>. Downloaded on 29 Oct 2011. (Gland, Switzerland, IUCN - <u>http://www.iucnredlist.org</u>.
4	 Schipper J., Chanson J.S., Chiozza F., Cox N.A., Hoffmann M., Katariya V., Lamoreux J., Rodrigues A.S.L., Stuart S.N., Temple H.J., et al. 2008 The Status of the World's Land and Marine Mammals: Diversity, Threat, and Knowledge. <i>Science</i> 322(5899), 225-230. (doi:10.1126/science.1165115).
5	. Stattersfield A.J., Capper D.R., Dutson G.C.L., BirdLife International, IUCN. 2000 <i>Threatened birds of the world : the official source for birds on the IUCN Red List.</i> Cambridge, Barcelona, BirdLife International; Lynx Edicions; xii, 852 p.
6	 Stuart S.N., Chanson J.S., Cox N.A., Young B.E., Rodrigues A.S.L., Fischman D.L., Waller R.W. 2004 Status and Trends of Amphibian Declines and Extinctions Worldwide. <i>Science</i> 306(5702), 1783-1786. (doi:10.1126/science.1103538).
7	 Hoffmann M., Hilton-Taylor C., Angulo A., Böhm M., Brooks T.M., Butchart S.H.M., Carpenter K.E., Chanson J., Collen B., Cox N.A., et al. 2010 The Impact of Conservation on the Status of the World's Vertebrates. <i>Science</i> 330(6010), 1503-1509. (doi:10.1126/science.1194442).
8	. Ceballos G., Ehrlich P.R. 2009 Discoveries of new mammal species and their
	22

implications for conservation and ecosystem services. *Proceedings of the National Academy of Sciences* **106**(10), 3841-3846. (doi:10.1073/pnas.0812419106).

- IUCN. 2001 IUCN Red List Categories & Criteria (version 3.1). (p. 30. Gland, Switzerland, IUCN.
- IUCN. 2006 IUCN Standards and Petitions Working Group: Guidelines for Using the IUCN Red List Categories and Criteria. Version 6.2. Prepared by the Standards and Petitions Working Group of the IUCN SSC Biodiversity Assessments Sub-Committee in December 2006. Downloadable from http://app.iucn.org/webfiles/doc/SSC/RedList/RedListGuidelines.pdf. (
- Mace G., Collar N., Gaston K., Hilton-Taylor C., Akcakaya H., Leader-Williams N., Milner-Gulland E., Stuart S. 2008 Quantification of extinction risk: IUCN's system for classifying threatened species. *Conserv Biol* 22(6), 1424-1442.
- 12. Mace G.M., Lande R. 1991 Assessing Extinction Threats: Toward a Reevaluation of IUCN Threatened Species Categories. *Conserv Biol* **5**(2), 148-157.
- Butchart S.H.M., Bird J.P. 2010 Data Deficient birds on the IUCN Red List: What don't we know and why does it matter? *Biol Conserv* 143(1), 239-247. (doi:10.1016/j.biocon.2009.10.008).
- Hurlbert A.H., Jetz W. 2007 Species richness, hotspots, and the scale dependence of range maps in ecology and conservation. *PNAS* 104, 13384-13389. (doi:10.1073/pnas.0704469104).
- Jetz W., McPherson J.M., Guralnick R.P. 2012 Integrating biodiversity distribution knowledge: toward a global map of life. *Trends in Ecology and Evolution* 27(3), 151-159. (doi:10.1016/j.tree.2011.09.007).
- Cardillo M., Mace G.M., Jones K.E., Bielby J., Bininda-Emonds O.R.P., Sechrest W., Orme C.D.L., Purvis A. 2005 Multiple Causes of High Extinction Risk in Large Mammal Species. *Science* 309(5738), 1239-1241.
- 17. Safi K., Pettorelli N. 2010 Phylogenetic, spatial and environmental components of extinction risk in carnivores. *Glob Ecol Biogeogr* **19**(3), 352-362.

 Jones K.E., Safi K. 2011 Ecology and evolution of mammalian biodiversity. *Philosophical Transactions of the Royal Society B: Biological Sciences* 366(1577), 2451-2461. (doi:10.1098/rstb.2011.0090).

- Davidson A.D., Hamilton M.J., Boyer A.G., Brown J.H., Ceballos G. 2009
 Multiple ecological pathways to extinction in mammals. *Proceedings of the National Academy of Sciences* 106(26), 10702-10705. (doi:10.1073/pnas.0901956106).
- Lee T.M., Jetz W. 2011 Unravelling the structure of species extinction risk for predictive conservation science. *Proceedings of the Royal Society B: Biological Sciences* 278(1710), 1329-1338. (doi:10.1098/rspb.2010.1877).
- Cardillo M., Mace G.M., Gittleman J.L., Jones K.E., Bielby J., Purvis A. 2008 The predictability of extinction: biological and external correlates of decline in mammals. *Proceedings of the Royal Society B: Biological Sciences*.
- 22. Pagel M. 1997 Inferring evolutionary processes from phylogenies. *Zoologica Scripta* **26**, 331-348.
- 23. Garland T., Midford P.E., Ives A.R. 1999 An introduction to phylogeneticallybased statistical methods with a new method for confidence intervals on ancestral values. *American Zoologist* **39**, 374-388.
- Garland T.J., Ives A.R. 2000 Using the past to predict the present: confidence intervals for regression equations in phylogenetic comparative methods. *American Naturalist* 155, 346-364.
- Peres-Neto P.R. 2006 A unified strategy for estimating and controlling spatial, temporal and phylogenetic autocorrelation in ecological models. *Oecologia Brasiliensis* 10, 105-119.
- 26. Martins E.P., Hansen T.F. 1997 Phylogenies and the comparative method: A general approach to incorporating phylogenetic information into the analysis of interspecific data. *American Naturalist* **149**(4), 646-667.
- 27. Freckleton R.P., Jetz W. 2009 Space versus phylogeny: disentangling phylogenetic and spatial signals in comparative data. *Proceedings of the Royal Society B: Biological Sciences* 276(1654), 21-30.

	Freckleton R.P., Cooper N., Jetz W. 2011 Comparative Methods as a Statistical :: The Dangers of Ignoring an Evolutionary Model. <i>American Naturalist</i> 178 (1), 0-E17. (doi:10.1086/660272).
29. Reg	Beale C.M., Lennon J.J., Yearsley J.M., Brewer M.J., Elston D.A. 2010 gression analysis of spatial data. <i>Ecol Lett</i> 13 (2), 246-264.
30. 202	IUCN. 2009 IUCN Red List of Threatened Species. Version 2009, Version 10.4. <u>http://www.iucnredlist.org</u> . Downloaded on 27 October 2010. (
	Bininda-Emonds O.R.P., Cardillo M., Jones K.E., MacPhee R.D.E., Beck R.M.D., enyer R., Price S.A., Vos R.A., Gittleman J.L., Purvis A. 2007 The delayed rise of sent-day mammals. <i>Nature</i> 446 , 507-512.
life	Jones K.E., Bielby J., Cardillo M., Fritz S.A., O'Dell J., Orme C.D.L., Safi K., chrest W., Boakes E.H., Carbone C. 2009 PanTHERIA: a species-level database of history, ecology, and geography of extant and recently extinct mammals. <i>Ecology</i> (9), 2648-2648.
	Bartholomé E., Belward A. 2005 GLC2000: a new approach to global land cover pping from Earth observation data. <i>International Journal of Remote Sensing</i> 26 (9), 59-1977.
exi	Herold M., Mayaux P., Woodcock C.E., Baccini A., Schmullius C. 2008 Some allenges in global land cover mapping: An assessment of agreement and accuracy in sting 1 km datasets. <i>Remote Sensing of Environment</i> 112 (5), 2538-2556. (doi:DOI: 1016/j.rse.2007.11.013).
35. The	Sanderson E., Jaiteh M., Levy M., Redford K., Wannebo A., Woolmer G. 2002 e human footprint and the last of the wild. <i>Bioscience</i> 52 (10), 891-904.
Da Cei	WCS. 2005 Last of the Wild Project, Version 2, 2005 (LWP-2): Global Human luence Index (HII) Dataset (Geographic). Palisades, NY: NASA Socioeconomic ta and Applications Center (SEDAC) (Wildlife Conservation Society (WCS); nter for International Earth Science Information Network (CIESIN), Columbia iversity. 2005
37.	Pagel M. 1999 Inferring the historical patterns of biological evolution. <i>Nature</i> 25
	http://mc.manuscriptcentral.com/issue-ptrsb

, 877-884.

- Haining R. 1990 Spatial data analysis in the social and environmental sciences§.
 Cambridge, Cambridge University Press.
- 39. Freckleton R.P., Harvey P.H., Pagel M. 2002 Phylogenetic dependence and ecological data: a test and review of evidence. *Amercian Naturalist* **160**, 716-726.
- 40. Cooper N., Purvis A. 2010 Body Size Evolution in Mammals: Complexity in Tempo and Mode. *The American Naturalist* 175(6), 727-738. (doi:doi:10.1086/652466).
- 41. Nakagawa S., Freckleton R.P. 2008 Missing inaction: the dangers of ignoring missing data. *Trends in Ecology and Evolution* **23**, 592-596.
- 42. Rubin D.B. 1987 *Multiple imputation for nonresponse in surveys*. New York, John Wiley & Sons.
- 43. Schafer J.L. 1997 Analysis of incomplete multivariate data. Chapman & Hall.
- 44. McKnight P.E., Mcknight K.M., Sidani S., Figueredo A.J. 2007 *Missing data: a gentle introduction*. New York, Guilford Press.
- González-Suárez M., Lucas P.M., Revilla E. 2012 Biases in comparative analyses of extinction risk: mind the gap. *J Anim Ecol* 81(6), 1211–1222. (doi:10.1111/j.1365-2656.2012.01999.x).
- Whittingham M.J., Stephens P.A., Bradbury R.B., Freckleton R.P. 2006 Why do we still use stepwise modelling in ecology and behaviour? *J Anim Ecol* 75(5), 1182-1189.
- 47. Hanley J.A., McNeil B.J. 1982 The Meaning and Use of the Area under a Receiver Operating Characteristic (ROC) Curve. *Radiology* **143**(1), 29-36.
- Hadfield J.D., Nakagawa S. 2010 General quantitative genetic methods for comparative biology: phylogenies, taxonomies and multi-trait models for continuous and categorical characters. *J Evolution Biol* 23(3), 494-508. (doi:10.1111/j.1420-9101.2009.01915.x).
- 49. Ives A.R., Helmus M.R. 2011 Generalized linear mixed models for phylogenetic

analyses of community structure. *Ecological Monographs* 81(3), 511-525.
(doi:10.1890/10-1264.1).
50. Freckleton R.P. 2012 Fast likelihood calculations for comparative analyses.

- Methods in Ecology and Evolution **3**(5), 940-947. (doi:10.1111/j.2041-210X.2012.00220.x).
- 51. Purvis A., Gittleman J.L., Cowlishaw G., Mace G.M. 2000 Predicting extinction risk in declining species. *Proc R Soc Lond Ser B-Biol Sci* **267**(1456), 1947-1952.
- 52. Cooper N., Bielby J., Thomas G.H., Purvis A. 2008 Macroecology and extinction risk correlates of frogs. *Glob Ecol Biogeogr* **17**(2), 211-221.

Ceballos G., Ehrlich P.R. 2006 Global mammal distributions, biodiversity hotspots, and conservation. *PNAS* 103(51), 19374-19379. (doi:10.1073/pnas.0609334103).

- Grenyer R., Orme C.D.L., Jackson S.F., Thomas G.H., Davies R.G., Davies T.J., Jones K.E., Olson V.A., Ridgely R.S., Rasmussen P.C., et al. 2006 Global distribution and conservation of rare and threatened vertebrates. *Nature* 444(7115), 93-96.
- Rondinini C., Rodrigues A.S.L., Boitani L. 2011 The key elements of a comprehensive global mammal conservation strategy. *Philosophical Transactions of the Royal Society B: Biological Sciences* 366(1578), 2591-2597. (doi:10.1098/rstb.2011.0111).
- Stuart S.N., Wilson E.O., McNeely J.A., Mittermeier R.A., Rodriguez J.P. 2010 The Barometer of Life. *Science* 328(5975), 177-. (doi:10.1126/science.1188606).
- Baillie J.E.M., Collen B., Amin R., Akcakaya H.R., Butchart S.H.M., Brummitt N., Meagher T.R., Ram M., Hilton-Taylor C., Mace G.M. 2008 Toward monitoring global biodiversity. *Conservation Letters* 1(1), 18-26.
- Fielding A.H., Bell J.F. 1997 A review of methods for the assessment of prediction errors in conservation presence/absence models. *Environ Conserv* 24(1), 38-49.

Figure Legends

Fig. 1. Explanatory and discriminative power of the fitted models of threat status for assessed species. (A), the relationship between fitted threat probability and IUCN status for assessed species. Threat probability is the probability a species is in one of the 'threatened' categories according to our spatial-phylogenetic multi-predictor model. (B), Receiver Operator Characteristic (ROC) curve, showing the relationship between true positive (sensitivity) and false positive (1 minus specificity) rate. The dashed line is the expected pattern if the threat probabilities were no better than random at discriminating threatened species. The AUC, which varies between 0.5 and 1, is the area highlighted in grey and is a measure of explanatory power. (C), the frequency distribution of fitted probabilities for species of contrasting conservation status. The green bars refer to species of 'least concern' (IUCN status 1 in *A*), whilst the red bars refer to species which are 'critically endangered' (IUCN status 5 in *A*). (D) Fitted / predicted threat probabilities shown separately for assessed species (grey) and data-deficient species (red).

Fig. 2. Prediction of threats for individual mammal orders. For each order the average fitted threat probabilities for assessed species (black points) and predicted threat probabilities for data-deficient species (red points) are shown (± standard errors). F-ratios and p-values refer to tests of differences between the mean fitted threat probabilities of assessed and data-deficient species. The numbers of assessed species is given for each order, together with the number of data-deficient and threatened (i.e. not 'least concern') species. The grey vertical bars show the threshold threat probability for each order (see Table S1), which is used to denote which data-deficient species are predicted to be threatened. The threshold is the point at which sensitivity = specificity (where

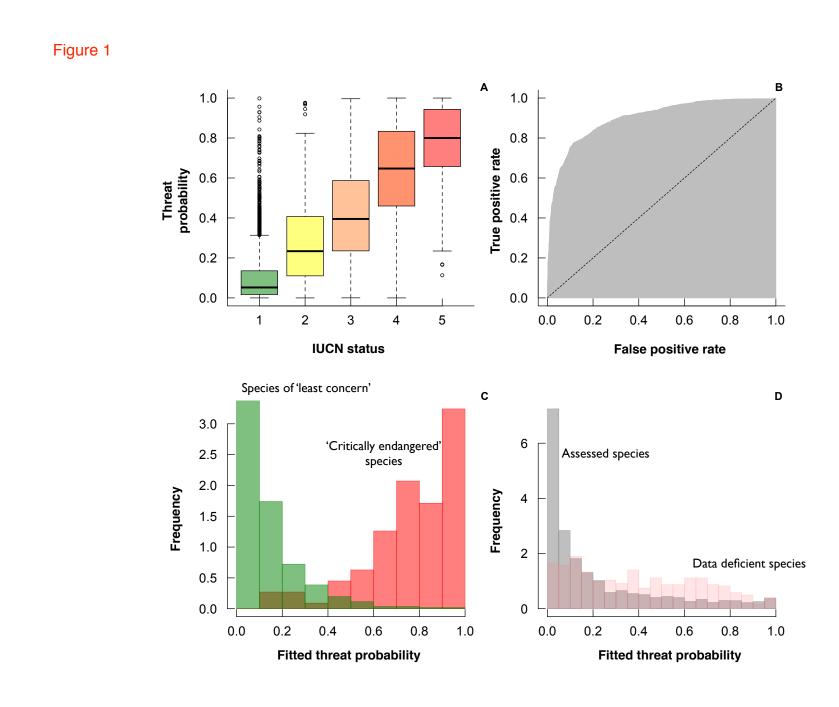
threatened and non-threatened status have equal chance of being correctly predicted [58]). Based on this probability, the final column gives the number of data-deficient species which are predicted to be threatened. Note that there are no data-deficient species in Perrissodactyla and the orders is thus not included here.

Fig. 3. The geography of observed and predicted mammal threat levels and richness. Panels illustrate the observed and predicted grid-cell proportions of all species assessed by IUCN to be threatened and analyzed here (A,B, 3,703 species, for model details see Table S1), and the predicted proportion of data-deficient species threatened (C, 483 species, for details see Fig 2). (D-F) show observed and predicted richness patterns: the richness of observed (assessed) threatened species (D, 812 of 3,703 species in analysis, i.e. all those assessed "vulnerable", "endangered" or "critically endangered"), those data-deficient species predicted by the combined spatial, phylogenetic and environmental model to be non-threatened (E, N = 152 of 483 species), and those predicted to be threatened (F, 331 of 483 species). Quantile classification of values across 110km equal area grid cells in Behrman projection. Note that color scales vary to emphasize geographic differences.

Fig. 4. Relationships between observed and predicted threat levels of grid cell assemblages.

The model-based predictions of average probability (A) and total proportion (B) of species threatened successfully captures the observed variation in proportion species threatened (A, Spearman r = 0.72; B, r = 0.66; 3,703 species; cf. Fig 3A, B). Observed and predicted richness of threatened assessed species is tightly associated (C, r = 0.77, cf. Fig 3D). In contrast, the predicted average threat levels and proportions of data-deficient species (cf. Fig 3C) show only very weak

association with the proportional threat patterns of assessed species (D, r = 0.33, E, r = 0.23; 843 data-deficient species). Equally, the areas with high richness of data-deficient species predicted to be threatened shows little covariance with those of high assessed threatened richness (F, r = 0.31, cf. Fig 3F). Darker gray represent higher density of points, line indicates 1:1 relationship. A total of 11,331 110km equal area grid cells that had \geq 50% dry land or were oceanic islands and had \geq 2 assessed species were analyzed.



	# of spe		Threate	
DIDELPHIMORPHIA F = 15.58** Opossums	Assessed 50	DD 11	Assessed 2	3
PERAMELEMORPHIA F = 25.92** Bandicoots & bilbies	16	2	6	2
DASYUROMORPHIA Carnivorous marsupials	59	1	10	1
DIPROTODONTIA Marsupial mammals	106	2	27	0
CINGULATA Armadillos	18	2	4	2
AFROSORICIDA Golden Moles / Tenrec / Otter Shrews	32	4	9	3
MACROSCELIDEA _F = 3.15(ns)	12	3	2	0
EULIPOTYPHLA	300	45	63	35
CHIROPTERA Bats F = 119.29***	759	126	126	89
CARNIVORA	212	14	55	7
CETARTIODACTYLA	193	7	82	6
PRIMATES F = 0.15(ns) Primates	220	8	116	5
SCANDENTIA F = 0.17(ns)	17	2	2	0
RODENTIA F = 273.34***	1618	252	281	174
LAGOMORPHA Hares, rabbits & pikas	75	4	14	4
0.0 0.2 0.4 0.6 0.8 1. Mean threat probabilty	0			

