

1

2

## The ecological drivers of variation in global language diversity

3

4

5 Xia Hua<sup>1,2</sup>, Simon J. Greenhill<sup>1,3</sup>, Marcel Cardillo<sup>2</sup>, Hilde Schneemann<sup>1,2,4</sup>, Lindell Bromham<sup>1,2</sup>

6

7

8 1. ARC Centre of Excellence for the Dynamics of Language, Australian National University,

9 Canberra ACT 0200, AUSTRALIA

10 2. Macroevolution and Macroecology Group, Division of Ecology & Evolution, Research School of

11 Biology, Australian National University, Canberra ACT 0200, AUSTRALIA

12 3. Max Planck Institute for the Science of Human History, Kahlaische Strasse 10, D-07743 Jena,

13 GERMANY

14 4. Meme Programme, University of Groningen, Nijenborgh 7, 9747 AG Groningen, THE

15 NETHERLANDS

1 Abstract

2

3 Language diversity is distributed unevenly over the globe. Why do some areas have so many  
4 different languages and other areas so few? Intriguingly, patterns of language diversity resemble  
5 biodiversity patterns, leading to suggestions that similar mechanisms may underlie both linguistic  
6 and biological diversification. Here we present the first global analysis of language diversity that  
7 identifies the relative importance of two key ecological mechanisms suggested to promote language  
8 diversification - isolation and ecological risk - after correcting for spatial autocorrelation and  
9 phylogenetic non-independence. We find significant effects of climate on language diversity  
10 consistent with the ecological risk hypothesis that areas of high year-round productivity lead to  
11 more languages by supporting human cultural groups with smaller distributions. Climate has a  
12 much stronger effect on language diversity than landscape features that might contribute to isolation  
13 of cultural groups, such as altitudinal variation, river density, or landscape roughness. The  
14 association between biodiversity and language diversity appears to be an incidental effect of their  
15 covariation with climate, rather than a causal link between the two. While climate and landscape  
16 provide strong explanatory signal for variation in language diversity, we identify a number of areas  
17 of high unexplained language diversity, with more languages than would be predicted from  
18 environmental features alone; notably New Guinea, the Himalayan foothills, West Africa, and  
19 Mesoamerica. Additional processes may be at play in generating higher than expected language  
20 diversity in these regions.

21

22 Keywords: language diversity, latitudinal gradient, phylogenetic non-independence, spatial  
23 autocorrelation, macroevolution, macroecology

## 1 Introduction

2

3 The geographic distribution of the world's >7500 languages is strikingly uneven<sup>1</sup> (Figure 1). For  
4 example, Papua New Guinea represents over 10% of the world's languages in less than 0.5% of the  
5 world's land area<sup>1</sup>. In contrast, the Russian federation covers 11% of the world land area but  
6 accounts for only 1.5% of the world's languages. The highly uneven distribution of languages  
7 remains a major unsolved problem in linguistics<sup>2,3</sup>. Yet, there are broad geographic patterns in the  
8 distribution of languages, suggesting a strong role for environmental determination of language  
9 diversity. The most notable of these patterns are latitudinal gradients: language diversity increases  
10 towards the equator<sup>4-11</sup>, and languages in the tropics tend to be restricted to smaller areas than  
11 languages at higher latitudes<sup>4,6,8</sup>.

12

13 Two broad kinds of ecological mechanism have been offered to explain geographic variation in  
14 language diversity: isolation and ecological risk. Isolation mechanisms are associated with  
15 landscape and geographic features that act as barriers to human movement. Such physical barriers  
16 may reduce interaction between groups, slowing the spread of linguistic variants among  
17 neighbouring populations, leading to the accumulation of language changes that distinguish each  
18 language from its neighbours<sup>12</sup>. In support of isolation mechanisms, studies have identified  
19 geographic correlates of language diversity such as river density<sup>13</sup>, landscape roughness<sup>13,14</sup>,  
20 elevation range<sup>15</sup>, and habitat diversity<sup>14</sup>.

21

22 The ecological risk mechanism for language diversity predicts an association between language  
23 diversity and climatic factors such as seasonal temperature variation and yearly rainfall<sup>16-19</sup>.

24 Smaller social groups are more likely to be stable and self-sufficient in areas with a more abundant  
25 and reliable year-round food supply. In contrast, areas of high seasonality or unpredictable rainfall  
26 may require communities to form social bonds across larger regions to obtain food and resources  
27 when they are scarce. In support of the ecological risk hypothesis, various studies have shown

1 correlations between language diversity and environmental productivity<sup>20</sup>, mean growing season<sup>16-</sup>  
2 <sup>18</sup>, rainfall<sup>4,20,21</sup>, and temperature<sup>4</sup>.

3

4 Areas of low ecological risk tend to occur in high-productivity environments that often also  
5 promote high plant and animal diversity<sup>22-24</sup>. Indeed, countries with high vertebrate and flowering  
6 plant diversity tend to have high language diversity<sup>25</sup>, and there are significant correlations between  
7 biological and language diversity at regional<sup>20,26</sup> and global scales<sup>27-29</sup>. It has also been suggested  
8 that high biodiversity could have a direct link with language diversity, if small social groups  
9 possess specialized knowledge of local biodiversity<sup>26,30</sup>, or if group boundaries are maintained to  
10 actively control local biodiversity resources<sup>31</sup>.

11

12 Although previous studies have reported correlations between language diversity and a number of  
13 environmental or landscape variables, there has yet to be a comprehensive global analysis. Many of  
14 these variables co-vary and tend to be clustered in space and more similar between related  
15 languages. Here we present the first global-scale analysis of language diversity that  
16 comprehensively deals with these statistical complexities. Our aims are to untangle and clarify the  
17 large-scale patterns of association between language diversity and key environmental variables, and  
18 thereby determine the relative strength of support for the alternative ecological mechanisms that  
19 may drive global patterns of language diversity.

20

21

## 22 Results and Discussion

23

24 To weigh the relative explanatory power of the isolation and ecological risk hypotheses, we tested  
25 the association between language diversity and six climatic variables, four geographic variables,  
26 two human population variables and four biodiversity measures. Our analyses were based on  
27 global-scale datasets of the geographic distribution of 6425 languages, high-resolution climatic and

1 geographic data layers, and global biodiversity datasets. We used the number of languages whose  
2 distribution overlaps each cell of a global equal-area grid as the measure of language diversity. A  
3 grid-based approach eliminates variation in language diversity and other variables due to  
4 differences in land area. It also allowed us to repeat analyses under three different spatial  
5 resolutions, as previous studies have shown that spatial resolution can influence tests for latitudinal  
6 gradients in language diversity<sup>7,26</sup>.

7

8 ***(i) Are patterns of language diversity influenced by phylogenetic and spatial non-independence?***

9

10 In order to untangle causal connections from incidental associations, we need to account for sources  
11 of covariation in our data. In particular, we need to address spatial autocorrelation and phylogenetic  
12 non-independence<sup>32,33</sup>. Grid cells that are located near each other are likely to have similar values of  
13 climatic and landscape variables, contain related human cultures and languages, and share much of  
14 their flora and fauna. If there are any particular features of cultures or languages that correlate with  
15 language diversity, then they will tend to covary with environmental features and biodiversity  
16 measures, whether or not there is any causal connection between them.

17

18 Our analyses show that regions in close geographic proximity, and with a high degree of language  
19 relatedness, tend to be more similar in their language diversity (Figures 1 and S1). Spatial  
20 autocorrelation and phylogenetic relatedness account for 15% of the variance in language diversity  
21 in analyses at low spatial resolution, 18% at medium resolution and 5% at high resolution, even  
22 after removing highly correlated grid cells (see Methods). This confirms the need to account for  
23 spatial proximity and phylogenetic relatedness between grid cells when testing for correlates of  
24 language diversity.

25

26 After correcting for spatial autocorrelation and phylogenetic relatedness, there is a latitudinal  
27 gradient in language diversity, with more languages near the equator than at higher latitudes (Figure

1 1). The regression coefficient of absolute latitude against language diversity is significantly  
2 negative under all resolutions (low:  $t=-2.58$ ;  $p=0.011$ ; medium:  $t=-2.49$ ;  $p=0.014$ ; high:  $t=-5.22$ ;  
3  $p<0.001$ ).

4

5 ***(ii) Is language diversity influenced by climate?***

6

7 We tested six climatic variables for associations with language diversity: mean annual temperature,  
8 mean annual precipitation, temperature seasonality, precipitation seasonality, net primary  
9 productivity, and mean annual growing season (Figure S2). Among these climatic variables,  
10 precipitation seasonality has the strongest association with language diversity in low resolution  
11 analyses, and temperature seasonality has the strongest association with language diversity in  
12 medium and high resolution analyses, independently of their covariation with other climatic  
13 variables (Table 1). The six climatic variables also provide sufficient explanation for the latitudinal  
14 gradient in language diversity, because adding latitude as an explanatory variable in the model in  
15 addition to the six climatic variables does not significantly increase the model fit under any  
16 resolution (low:  $LR=0.94$ ,  $p=0.33$ ; medium:  $LR=1.15$ ,  $p=0.28$ ; high:  $LR=2.03$ ,  $p=0.15$ ).

17

18 The ecological risk hypothesis may provide an explanation for the association between language  
19 diversity and seasonality, which predicts higher language diversity in regions with longer periods of  
20 reliable food production, by allowing smaller cultural groups to be self-sufficient<sup>16-19</sup>. This  
21 hypothesis makes testable predictions about the associations between climate, population size and  
22 density, and language diversity. We followed previous studies<sup>16-19</sup> in using mean growing season  
23 (the number of days per year suitable for growing crops) as an indicator of ecological risk, although  
24 our results indicate that temperature seasonality may be a better predictor of the influence of  
25 environment on language diversity. The ecological risk hypothesis predicts that longer growing  
26 seasons will result in reduced area per language and smaller speaker population sizes. We find  
27 evidence to support both of these predictions. Longer growing seasons are associated with a greater

1 number of languages per grid cell at all three resolutions (Table 2), consistent with a reduction in  
2 range sizes of languages allowing tighter “packing” of languages. The increase in language  
3 diversity is not simply a result of areas with long growing seasons supporting a greater number of  
4 people, because mean growing season has a significant positive association with language diversity  
5 beyond its covariation with population density (Table 2). Mean growing season is negatively  
6 associated with minimum speaker population size (the population size of the smallest language in a  
7 grid cell) under medium and high resolutions (Table 2), consistent with the prediction that smaller  
8 cultural groups are more able to persist in areas of longer growing season. There is no association  
9 between mean growing season and the average speaker population size of all the languages in a grid  
10 cell, so increased packing is primarily a result of reduction in language range size in areas with  
11 longer growing seasons, rather than the formation of cultural groups with fewer members.

12

13 Our results are broadly consistent with the ecological risk hypothesis, because mean growing season  
14 is associated with both the minimum group size and the number of languages per grid cell.

15 However, we find that seasonality in temperature and precipitation have additional association with  
16 language diversity that is not due to mean growing season. This is consistent with a recent study  
17 that supports associations between language diversity and precipitation in the wettest quarter and  
18 temperature in the warmest quarter<sup>34</sup>. While growing season is defined by the number of days above  
19 a defined minimum temperature and moisture availability, seasonality will reflect both minimum  
20 and maximum of temperatures and moisture. Therefore, this result may suggest that climatic  
21 extremes across seasons shape language diversity, in addition to average length of growing season.

22

### 23 *(iii) Is language diversity influenced by landscape?*

24

25 To examine the effect of isolation on language diversity, we tested four landscape variables that  
26 may influence patterns of human movement and therefore contribute to the isolation of cultural  
27 groups: mean altitude, altitudinal range, landscape roughness, and river density (Figure S2). Higher

1 river density is associated with greater language diversity at low and medium resolutions, beyond  
2 its covariation with climatic variables and the other landscape variables (Table 3). This is consistent  
3 with previous proposals that rivers act to isolate populations into smaller language groups<sup>13</sup>.  
4 However, we find little additional support for this hypothesis. While river density is associated with  
5 smaller minimum speaker population size at medium resolution (Table 3), there is no association  
6 between river density and average speaker population size (controlling for the effects of population  
7 density). These observations suggest that the association between river density and language  
8 diversity is more akin to the ecological risk hypothesis than to the isolation hypothesis, because  
9 rivers seem to allow the persistence of smaller speaker populations, but not to divide human  
10 populations into smaller speaker populations. In this sense, rivers seem to act more as an ecological  
11 resource than a barrier to interaction.

12

13 Similarly, while altitudinal range is associated with language diversity at high resolution with  
14 marginal significance, there is no evidence that this is caused by isolation, as altitudinal range does  
15 not result in reduction in speaker population size, even when controlling for population density  
16 (Table 3). While landscape roughness is significantly associated with language diversity when  
17 altitudinal range is not included in the model ( $t=2.87$ ;  $p=0.004$ ), we find no significant association  
18 between landscape roughness and language diversity beyond its covariation with climatic variables  
19 and the other landscape variables under the three resolutions, and no statistically significant  
20 negative association between landscape roughness and speaker population size (Table 3).

21

22 In contrast to a previous study that described river density and landscape roughness as universal  
23 determinants of language diversity<sup>13</sup>, we find little evidence that landscape variables have a strong  
24 or consistent influence on language diversity. Although we use similar data to Axelson & Manrubia  
25 (2014), there are a number of differences in our analytical approach. To compare our results to  
26 theirs, we reanalyze our data using their method, fitting continent-specific parameter values and not  
27 including altitudinal range. Without correcting for spatial and phylogenetic non-independence



1 among grid cells, we get similar results to Axelson & Manrubia (2014), namely that river density  
2 and landscape roughness have significant associations with language diversity in most continents  
3 (Table S1). But when we correct the data for non-independence among grid cells, neither river  
4 density nor landscape roughness has a significant association with language diversity in any  
5 continent (Table S1). We therefore conclude that the previous result was driven primarily by spatial  
6 autocorrelation and phylogenetic non-independence, with the similarity in both landscape variables  
7 and language diversity between neighbouring grid cells generating spurious correlations.

8  
9 In conclusion, we find little consistent support for effect of isolation mechanisms on language  
10 diversity. While we find associations between language diversity and river density, altitudinal range  
11 and landscape roughness, these landscape factors have much less influence on language diversity  
12 than climatic factors, and there is little indication that this is caused by the division of human  
13 populations into smaller, isolated cultural groups. Instead, previous results suggesting river density  
14 and landscape roughness are universal determinants of language diversity<sup>13</sup> may have been driven  
15 by autocorrelation among grid cells.

16

17 ***(iv) Is language diversity significantly associated with biodiversity?***

18

19 We now ask if biodiversity provides any additional explanation for language variation beyond  
20 covariation with climate and landscape factors. Adding mammal or bird diversity as additional  
21 predictors to the climatic and landscape variables significantly improves model fit, but adding  
22 vascular plant and amphibian diversity do not provide additional explanatory power (Table 4).  
23 Adding biome to the analysis increases model fit above climate variables at low resolution,  
24 suggesting that ecosystem structures may influence language diversity, however it does not provide  
25 significant explanatory power above the effect of climate at medium and high resolutions (low:  
26  $LR=27.01$ ,  $p=0.02$ ; medium:  $LR=14.91$ ,  $p=0.38$ ; high:  $LR=11.83$ ,  $p=0.62$ ).

27

1 Why are bird and mammal diversity associated with language diversity? There is no evidence that  
2 this is due to a direct causal relationship between biodiversity and language diversity, because there  
3 is no consistent relationship between these biodiversity measures and residual variation in language  
4 diversity, above and beyond that explained by climate and landscape (Table S2). Instead, the  
5 increase in model fit when bird and mammal diversity are added to the model of language diversity,  
6 climate and landscape, seems to be driven primarily by regions that have both low language  
7 diversity and low species diversity, particularly the Sahara, the Arabian Peninsula, and the Tibetan  
8 Plateau (Figure 2 and S3), which present harsh environmental conditions for birds and mammals  
9 (including humans). These are not the only regions of low diversity but they seem to have a  
10 disproportionate influence on the relationship between mammal and bird diversity and language  
11 diversity (Figure S4). Running the high resolution analysis without these low diversity areas, we  
12 find that adding mammal or bird diversity as additional predictors to the climatic and landscape  
13 variables no longer increases model fit ( $n = 334$ , mammal:  $LR=1.92$ ,  $p=0.17$ ; bird:  $LR=3.67$ ,  
14  $p=0.07$ ), but results for the climatic and landscape effects are similar to the complete dataset.  
15 Temperature seasonality is still the strongest predictor for language diversity in the climatic  
16 variables ( $t=-2.34$ ,  $p=0.02$ ) and so is altitudinal range in the landscape variables ( $t=2.27$ ,  $p=0.02$ ).  
17 These results suggest that the low diversity areas have a significant effect on the association  
18 between biodiversity and language diversity, but they are not responsible for the broader association  
19 between language diversity and climatic and landscape effects.

20

21 In conclusion, we find that the association between language diversity and biodiversity appears to  
22 be largely a result of their covariation with common climatic and landscape factors, and any  
23 additional increase in model fit between language diversity and mammal and bird diversity is likely  
24 due to the disproportionate effect of a few regions of harsh environment that reduce both  
25 biodiversity and language diversity.

26

27 ***(v) What explains the residual variation in language diversity?***

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

The six climatic variables and the four landscape variables together explain 45% of the variance in language diversity under low resolution, 31% under medium resolution, and 27% under high resolution (after correction for phylogenetic and spatial non-independence). About 80% of the explanatory power is contributed by the six climatic variables under the three resolutions. Measures of biodiversity do not appear to add additional explanatory power beyond their covariation with climatic factors, above and beyond the influence of several key areas of low diversity.

What accounts for the remaining variation in language diversity? Figure 3 shows the distribution of the residuals in language diversity after removing the climatic and landscape effects on language diversity. We can identify areas of high unexplained language diversity as the red grid cells with residuals  $\geq 1.96$  standard deviations higher than predicted by the climatic and landscape variables alone. These grid cells are concentrated in four regions – New Guinea, the Himalaya foothills of Assam and Bangladesh, West Africa, and Mesoamerica. Language diversity in grid cells with residuals  $\leq -1.96$  (blue) is lower than we would predict based on the climatic and landscape variables, most notably in the lower Amazon Basin of South America.

There are two possible explanations for these areas of relative excess or paucity of languages, beyond that predicted by climate and landscape. One is that they reflect relative completeness of language documentation. For example, Amazonia is considered an area of high language diversity<sup>27</sup>, but incomplete documentation in the central areas of this region have led to it being described as the “least known and least understood linguistic region”<sup>35</sup>. Therefore, the true number of languages may be higher than the documented number of languages. However, it seems unlikely that the opposite effect (over-reporting of language diversity) would explain the areas of high unexplained language diversity.

1 Alternatively, it may be that other factors contribute significantly to shaping language diversity that  
2 are not captured by climate variables (representing the ecological risk hypothesis) nor by landscape  
3 variables (representing isolation mechanisms). For example, regions of higher than expected  
4 language diversity may have had a longer period of *in situ* language diversification, or have  
5 undergone a higher rate of diversification, leading to a greater accumulation of languages in these  
6 regions than in other regions of similar climate. One way to investigate the influence of time or  
7 diversification rate on diversity is to use a phylogeny that contains information on the relative  
8 timing of diversification events in order to compare the timescale and rate of diversification in  
9 different regions<sup>24,36</sup>. While phylogenies are available for the languages within some language  
10 families<sup>37-42</sup>, there is currently no global dated phylogeny of languages, nor is there general  
11 agreement on the relationships or age of language families. Therefore we lack the means to make a  
12 quantitative comparison of duration or rates of diversification between the majority of grid cells  
13 (those that contain languages from different families or languages not contained in comprehensive  
14 phylogenies).

15  
16 Nevertheless, we can make a qualitative comparison of the relative depth of divergence represented  
17 in each grid cell if we make the simple assumption that languages from the same language family  
18 diverged more recently than languages from different families. Number of language families per  
19 grid cell is a significant predictor of residuals in language diversity under the three resolutions (low:  
20  $t=4.65, p=<0.001$ ; medium:  $t=6.27, p=<0.001$ ; high:  $t=8.83, p=<0.001$ ; Figure S5). However, we are  
21 hesitant to draw strong conclusions from this pattern. For example, while New Guinea has more  
22 language families per grid cell than most other regions, the other areas of high unexplained  
23 language diversity do not have unusually high language family diversity, and some areas with many  
24 language families do not have high language richness (Figure 4; Figure S5). Clearly, this is not an  
25 ideal analysis of variation in time for diversification, as we cannot standardize time or rate of  
26 language evolution across families without a global dated phylogeny. But it suggests that time to

1 diversification may be a profitable area of enquiry once complete language phylogenies become  
2 available.

3

#### 4 **Conclusion**

5 The overall picture supported by our analyses is that environmentally-driven ecological processes  
6 are a major determinant of global variation in the diversity of human languages, as they are for  
7 global variation in biodiversity. Associations between global patterns of language diversity and  
8 climate are consistent with the ecological risk hypothesis, that stable productive climates allow  
9 human cultures to persist in smaller, more localized groups. Our results offer less support for  
10 isolation mechanisms as drivers of language diversity. While there are significant associations  
11 between language diversity and river density, altitudinal extent and landscape roughness, landscape  
12 factors have less explanatory power than climate. The association between biodiversity and  
13 language diversity is likely due to an incidental association between language and species richness  
14 driven by shared causal factors such as climate and landscape. The importance of influences such as  
15 time to accumulate diversity or the rate of language diversification are yet to be explored in detail.

16

17

#### 18 **Methods**

19 All analyses are based on grid cells for a global equal-area projection at three different resolutions:  
20 low resolution with grid cell size 1000x1000km, medium resolution with grid cell size 500x500km,  
21 and high resolution with grid cell size 200x200km. We excluded grid cells where no language  
22 distributions overlap that cell.

23

#### 24 **Data collection.**

25 *Languages.* Language distributions were compiled from World Language Mapping System v17<sup>1</sup>  
26 that includes information on the geographic distribution of 6425 languages. Language diversity was  
27 calculated by overlaying the language range polygons with the global grid using the R packages

1 ‘sp’ and ‘raster’<sup>43-45</sup> and counting the number of languages whose distribution overlaps all or part of  
2 a grid cell. We included the percentage of a grid cell covered by land as an independent variable in  
3 each regression model. For islands that cover less than 1% of the area of the grid cell, land coverage  
4 was set to 0.01 unless the exact number could be derived from the language data.

5  
6 *Population.* Both the ecological risk and isolation hypotheses make predictions about the  
7 relationship between speaker population size and climate or landscape factors. We included both the  
8 minimum and the average size of speaker populations of all the languages present in each grid cell,  
9 based on estimates of the number of native (“L1”) speakers of a language as recorded in World  
10 Language Mapping System v17<sup>1</sup>. To control for regional variations in number of people per grid  
11 cell, we used total human population density from the Gridded Population of the World database<sup>46</sup>.

12  
13 *Climate.* We included four climatic variables in our analysis: annual temperature, annual  
14 precipitation, temperature seasonality, and precipitation seasonality, averaged over each grid cell.  
15 We also included two variables derived from eco-climatic factors: net primary productivity and  
16 mean growing season of a grid cell. Net primary productivity data were derived from the  
17 Socioeconomic Data and Applications Center<sup>47</sup>. Data on growing season were obtained from the  
18 Global Agro-ecological Zones Data Portal version v3.0<sup>48</sup>, which is calculated as the number of days  
19 per year suitable for growing crops based on precipitation, evapotranspiration and soil moisture  
20 holding capacity. The other climatic variables were obtained from the Worldclim global climate  
21 data set v1.4<sup>49</sup>.

22  
23 *Landscape.* We included four variables describing landscape factors that potentially influence  
24 population movement and range expansion: average altitude, altitudinal range, landscape roughness,  
25 and river density in each grid cell. Altitude data were obtained from Worldclim<sup>49</sup>. Landscape  
26 roughness data were calculated as the autocorrelation in altitude<sup>13</sup>, derived from the SRTM30

1 elevation dataset<sup>50</sup>. River density was calculated as the number of river branches within each grid  
2 cell<sup>13</sup>, derived from the Global Self-consistent, Hierarchical, High-resolution Shoreline database<sup>51</sup>.  
3  
4 *Biodiversity*. Vascular plant richness data was from Kreft & Jetz (2006)<sup>52</sup>. Species richness of all  
5 the amphibian, mammal, and bird taxa was from BiodiversityMapping<sup>53,54</sup> that produces maps of  
6 species richness from species distribution data obtained from IUCN, BirdLife International and  
7 NatureServe databases. These maps were resampled to the grid resolutions we used in our analyses.  
8 To capture broad scale variation in ecosystem structure and composition, we also compiled data on  
9 the world's biomes from WWF<sup>55</sup>. Biomes are discrete regions with a distinct ecological character  
10 that is determined by a combination of climate, geomorphology and vegetation types<sup>55</sup>.

11

## 12 **Statistical analysis.**

13 *Model structure*. We applied generalized least squares (GLS) analysis, implemented in the R  
14 package *nlme*<sup>56</sup>, to fit regression models to log-transformed language diversity. Log-transforming  
15 predictor variables did not change the significant results, so we report results using untransformed  
16 predictors. We accounted for spatial autocorrelation and phylogenetic relatedness by constraining  
17 the residual correlation in language diversity between each pair of grid cells to be a linear function  
18 of the spatial proximity and phylogenetic similarity between the two cells. The correlation matrix  
19 has the form:  $(1 - \alpha)I + \alpha[\beta P + (1 - \beta)D]$ , where  $I$  is an identity matrix,  $P$  is the phylogenetic  
20 similarity matrix, and  $D$  is the spatial proximity matrix,  $\alpha$  represents the relative contribution of  
21 spatial and phylogenetic versus other residual effects,  $\beta$  represents the relative contribution of  
22 spatial versus phylogenetic effects<sup>57</sup>. Because our analysis controls for non-independence of grid  
23 cells, we can be more confident that the results are not driven by pseudoreplication. For example,  
24 without such correction, grid cells in the Arctic that repeatedly sample the same widely distributed  
25 languages (e.g. Russian and Yakut) may have a disproportionate influence on global language  
26 diversity correlations (Figure S1).

27

1 *Phylogenetic relatedness.* In order to correct for non-independence due to descent, we need a matrix  
2 of covariation representing expected patterns of similarity. There is no accepted universal  
3 phylogeny for the world's languages, so we constructed a global hierarchy of language relationships  
4 from the World Language Mapping System<sup>1</sup> taxonomy using the python library Treemaker<sup>58</sup>. This  
5 hierarchy is a proxy for the expected patterns of similarity due to relatedness and does not represent  
6 a phylogenetic history of descent. It represents the best available estimate of the relationships within  
7 language families and therefore provides a way to generate a matrix of expected similarity due to  
8 descent<sup>33</sup>. The global language taxonomy is only resolved to the language family level, so we  
9 assume that any pair of languages from different families represent the maximum distance from  
10 each other. This hierarchy is therefore completely unresolved at the base. The expected similarity  
11 due to relatedness of languages was calculated for each pair of grid cells using the *PhyloSor*  
12 metric<sup>59</sup>. This measure compares the sum of distances on the language hierarchy that connect all the  
13 languages that occur in a pair of grid cells to the sum of distances that connect all the languages  
14 occurring in each grid cell. This measure ranges from 0 to 1, with 0 for two grid cells that do not  
15 share any language families in common and 1 for two grid cells that have an identical set of  
16 languages.

17  
18 *Spatial proximity.* The spatial proximity matrix was derived from the great-circle distances between  
19 the centroids of each pair of grid cells. We modeled the decay in similarity of language diversity  
20 with distance as the Gaussian function  $e^{-(d/\gamma)^2}$ , where  $d$  is the great-circle distance between the two  
21 grid cells and  $\gamma$  is the coefficient describing how fast similarity decays over the distance between  
22 grid cells.

23  
24 *Subsampling.* Adjacent grid cells can share similar or identical values for environmental variables,  
25 as well as sharing many of the same species and languages, making their correlation coefficient at  
26 or close to 1. A large number of self-similar values lead to degeneracy of the matrix (with much less  
27 information than the number of entries in the matrix). Under medium and high resolutions,



1 correlation between adjacent grid cells is so high that the correlation matrix is nearly singular,  
2 leading to a high level of error when taking the inverse of a large matrix. We limit self-similarity  
3 across the correlation matrix by subsampling grid cells to avoid adjacent cells with highly similar  
4 values. For medium resolution, we avoided sampling adjacent cells by first removing the nine  
5 surrounding grid cells, i.e., sampling a grid cell every two rows and columns. This was insufficient  
6 to allow convergence of likelihood estimation for the high resolution grids, so we then removed the  
7 24 surrounding grid cells, i.e., sampling a grid cell every three rows and columns (Figure S6). This  
8 resulted in 216 grid cells under low resolution, 192 grid cells under medium resolution, and 366  
9 grid cells under high resolution. This subsampling procedure also has the effect of reducing the  
10 disparity in number of datapoints at different resolutions.

11

12 *Implementation.* We used the *subplex* method in the R package *nloptr*<sup>60</sup> to find the maximum-  
13 likelihood estimates for the coefficients in our regression models. To test if a variable is associated  
14 with language diversity above its covariation with the other variables, the variable was dropped  
15 from the full model that included all the variables, then a likelihood ratio test was used to test if  
16 dropping the variable significantly decreased model fit. To assess how much variance in language  
17 can be explained by the climatic and landscape variables, we calculated the predicted  $R^2$  of the  
18 regression model that included all the climatic and landscape variables as predictors. To evaluate  
19 the contribution of phylogenetic non-independence and spatial autocorrelation, we refitted the  
20 regression model using the method of ordinary least squares (OLS), which does not account for  
21 correlation structure in language diversity among grid cells. Difference in the predicted  $R^2$  between  
22 the GLS method and the OLS method quantifies the impact of spatial autocorrelation and  
23 phylogenetic non-independence to the results.

24

25 Acknowledgements:

26 We thank Holger Kreft & Walter Jetz for supplying species richness data for vascular plants.

27

1 Data availability statement:

2 All the data are from published and publicly available databases. Publications and web links for  
3 these datasets are reported in the references. Figure 1, 2, 4, S2-4 have associated raw data. No  
4 restriction on data availability.

5

6 References:

- 7 1. Lewis, P. M. et al. *Ethnologue: Languages of the World*. (SIL International, Dallas, 2014).
- 8 2. Gavin, M. C. et al. Toward a Mechanistic Understanding of Linguistic Diversity. *BioScience* **63**,  
9 524–535 (2013).
- 10 3. Greenhill, S. J. Demographic correlates of language diversity. In *The Routledge Handbook of*  
11 *Historical Linguistics* (eds Bower, C. & Evans, B.) 555–578 (Routledge, London, 2014).
- 12 4. Collard, I. F. & Foley, R. A. Latitudinal patterns and environmental determinants of recent  
13 human cultural diversity: do humans follow biogeographical rules? *Evol. Ecol. Res.* **4**, 371–383  
14 (2002).
- 15 5. Currie, T. E. & Mace, R. The Evolution of Ethnolinguistic Diversity. *Adv. Complex Syst.* **15**,  
16 1150006 (2012).
- 17 6. Gavin, M. C. & Stepp, J. R. Rapoport's Rule Revisited: Geographical Distributions of Human  
18 Languages. *PloS One* **9**, e107623 (2014)
- 19 7. Hillebrand, H. On the generality of the latitudinal diversity gradient. *Am. Nat.* **163**, 192–211  
20 (2004).
- 21 8. Mace, R. & Pagel, M. A latitudinal gradient in the density of human languages in North  
22 America. *Proc. Royal. Soc. B.* **261**, 117–121 (1995).
- 23 9. Nichols, J. *Linguistic diversity in space and time*. University of Chicago Press (1992)
- 24 10. Rohde, K. Latitudinal gradients in species diversity and Rapoport's rule revisited: a review of  
25 recent work and what can parasites teach us about the causes of the gradients? *Ecography* **22**,  
26 593–613 (1999).

- 1 11. Sax, D. F. Latitudinal gradients and geographic ranges of exotic species: implications for  
2 biogeography. *J Biogeogr.* **28**, 139–150 (2001).
- 3 12. Labov, W. Transmission and Diffusion. *Language* **83**, 344–387 (2007).
- 4 13. Axelsen, J. & Manrubia, S. River density and landscape roughness are universal determinants of  
5 linguistic diversity. *Proc. Royal. Soc. B.* **281**, 20133029 (2014).
- 6 14. Cashdan, E. Ethnic diversity and its environmental determinants: Effects of climate, pathogens,  
7 and habitat diversity. *Am. Anthropol.* **103**, 968–991 (2001).
- 8 15. Amano, T. et al. Global distribution and drivers of language extinction risk. *Proc. Royal. Soc. B.*  
9 **281**, 20141574 (2014).
- 10 16. Coupé, C. et al. Investigations into determinants of the diversity of the world's languages. In  
11 *Eastward flows the great river: Festschrift in honor of Professor William S-Y. Wang on his 80th*  
12 *birthday.* (eds Peng, G. & Shi, F.) 75–108 (City University of HK Press, Hong Kong, 2013).
- 13 17. Nettle, D. Explaining global patterns of language diversity. *J. Anthropol. Archaeol.* **17**, 354–374  
14 (1998).
- 15 18. Nettle, D. *Linguistic Diversity.* Oxford University Press, Oxford (1999).
- 16 19. Nettle, D. Language Diversity in West Africa: An ecological approach. *J. Anthropol. Archaeol.*  
17 **15**, 403–438 (1996).
- 18 20. Moore, J. L. et al. The distribution of cultural and biological diversity in Africa. *Proc. Royal.*  
19 *Soc. B.* **269**, 1645–1653 (2002).
- 20 21. Nichols, J. Modeling ancient population structures and movement in linguistics. *Annu. Rev.*  
21 *Anthropol.* **26**, 359–384 (1997).
- 22 22. Quintero, I. & Jetz, W. Global elevational diversity and diversification of birds. *Nature* **555**,  
23 246–250 (2018).
- 24 23. Jetz, W. & Fine, P. V. A. Global gradients in vertebrate diversity predicted by historical area-  
25 productivity dynamics and contemporary environment. *PLoS Biol.* **10**, e1001292 (2012).
- 26 24. Mittelbach G. et al. Evolution and the latitudinal diversity gradient: speciation, extinction and  
27 biogeography. *Ecol. Lett.* **10**, 315–331 (2007).

- 1 25. Harmon, D. Losing species, losing languages: connections between biological and linguistic  
2 diversity. *Southwest J. Linguist.* **15**, 89–108 (1996).
- 3 26. Manne, L. L. Nothing has yet lasted forever: current and threatened levels of biological and  
4 cultural diversity. *Evol. Ecol. Res.* **5**, 517–527 (2003).
- 5 27. Gorenflo, L. J. et al. Co-occurrence of linguistic and biological diversity in biodiversity hotspots  
6 and high biodiversity wilderness areas. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 8032–8037 (2012)
- 7 28. Stepp, J. R. et al. Development of a GIS for global biocultural diversity. *Policy Matters* **13**,  
8 267–270 (2004).
- 9 29. Sutherland, W. J. Parallel extinction risk and global distribution of languages and species.  
10 *Nature* **423**, 276–279 (2003).
- 11 30. Maffi, L. Linguistic, cultural, and biological diversity. *Annu. Rev. Anthropol.* **34**, 599–617  
12 (2005).
- 13 31. Pagel, M. & Mace, R. The cultural wealth of nations. *Nature* **428**, 275–278 (2004).
- 14 32. Cardillo, M. et al. Links between language diversity and species richness can be confounded by  
15 spatial autocorrelation. *Proc. Royal Soc. B.* **282**, 20142986 (2015).
- 16 33. Bromham, L. et al. Parasites and politics: why cross-cultural studies must control for  
17 relatedness, proximity and covariation. *Royal Soc. Open Sci.* **5**, 181100 (2018).
- 18 34. Derungs, C. et al. Environmental factors drive language density more in food-producing than in  
19 hunter–gatherer populations. *Proc. Royal. Soc. B.* **285**, 20172851 (2018).
- 20 35. Dixon, R. M. W. & Aikhenvald, A. Y. The Amazonian languages. Cambridge University Press,  
21 Cambridge (1998).
- 22 36. Wiens, J. J. & Donoghue, M. J. Historical biogeography, ecology and species richness. *Trends*  
23 *Ecol. Evol.* **19**, 639–644 (2004).
- 24 37. Bouckaert, R. R. et al. The origin and expansion of Pama–Nyungan languages across Australia.  
25 *Nature Ecol. Evol.* **2**, 741–749 (2018).
- 26 38. Bouckaert, R. R. et al. Mapping the Origins and Expansion of the Indo-European Language  
27 Family. *Science* **337**, 957–960 (2012).

- 1 39. Gray, R. D. et al. Language phylogenies reveal expansion pulses and pauses in Pacific  
2 settlement. *Science* **323**, 479–483 (2009).
- 3 40. Grollemund, R. et al. Bantu expansion shows habitat alters the route and pace of human  
4 dispersals. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 13296–13301 (2015).
- 5 41. Kolipakam, V. et al. A Bayesian phylogenetic study of the Dravidian language family. *Royal*  
6 *Soc. Open Sci.* **5**, 171504 (2018).
- 7 42. Chang, W. et al. Ancestry-constrained phylogenetic analysis supports the Indo-European steppe  
8 hypothesis. *Language* **91**, 194–244 (2015).
- 9 43. Hijmans, R. raster: Geographic data analysis and modeling. R package version 2.3-40. (2015).
- 10 44. Pebesma, E. J. & Bivand, R. S. Classes and methods for spatial data in R. *R news* **5**, 9–13  
11 (2005).
- 12 45. R Core Team. R: A language and environment for statistical computing. In: Team RC (ed) R  
13 Foundation for Statistical Computing, Vienna (2014).
- 14 46. Center for International Earth Science Information Network, Columbia University. Gridded  
15 Population of the World, Version 4 (GPWv4): Population Density Adjusted to Match 2015  
16 Revision of UN WPP Country Totals. (Palisades, NY: NASA Socioeconomic Data and  
17 Applications Center, 2016).
- 18 47. Imhoff, M. L. et al. HANPP Collection: Global Patterns in Net Primary Productivity (NPP).  
19 (Palisades, NY: NASA Socioeconomic Data and Applications Center, 2004).
- 20 48. FAO/IIASA. Global Agro-ecological Zones (GAEZ v3.0). <http://www.fao.org/nr/gaez>. (2010).
- 21 49. Hijmans, R. J. et al. Very high resolution interpolated climate surfaces for global land areas. *Int.*  
22 *J. Climatol.* **25**, 1965–1978 (2005).
- 23 50. Becker, J. J. et al. Global bathymetry and elevation data at 30 arc seconds resolution:  
24 SRTM30\_PLUS. *Mar. Geodesy* **320**, 355–371 (2009).
- 25 51. Wessel, P. & Smith, W. H. F. A global, self-consistent, hierarchical, high-resolution shoreline  
26 database. *J. Geophys. Res.* **101**, 8741–8743 (1996).

- 1 52. Kreft, H. & Jetz, W. Global patterns and determinants of vascular plant diversity. *Proc. Natl.*  
2 *Acad. Sci. U. S. A.* **104**, 5925–30. (2007)
- 3 53. Pimm, S. L. et al. The biodiversity of species and their rates of extinction, distribution, and  
4 protection. *Science* **344**, 987–998 (2014).
- 5 54. Jenkins, C. N. et al. Global patterns of terrestrial vertebrate diversity and conservation. *Proc.*  
6 *Natl. Acad. Sci. U. S. A.* **110**, E2602–E2610 (2013).
- 7 55. Olson, D. M. et al. Terrestrial ecoregions of the world: a new map of life on earth: a new global  
8 map of terrestrial ecoregions provides an innovative tool for conserving biodiversity. *BioScience*  
9 **51**, 933–938 (2001).
- 10 56. Pinheiro, J. et al. nlme: Linear and Nonlinear Mixed Effects Models. R package version 3.1-  
11 131.1 (2018).
- 12 57. Freckleton, R. P. & Jetz, W. Space versus phylogeny: disentangling phylogenetic and spatial  
13 signals in comparative data. *Proc. Royal. Soc. B.* **276**, 21– 30 (2009).
- 14 58. Greenhill, S. J. Treemaker v1.0.A Python tool for constructing a newick formatted tree from a  
15 set of classifications. <https://github.com/SimonGreenhill/Treemaker>. (2016).
- 16 59. Bryant, J. A. et al. Microbes on mountainsides: contrasting elevational patterns of bacterial and  
17 plant diversity. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 11505–11511 (2008).
- 18 60. Johnson, S. G. The NLOpt nonlinear-optimization package. <http://ab-initio.mit.edu/nlopt>.  
19 (2014).
- 20

1 **Table 1.** Climatic effects on language diversity, at high, medium, and low resolution ( $n$  is the  
2 number of grid cells used in the analysis at each resolution). We list the  $t$  value and the  $p$  value of  
3 each predictor in a generalized least squares regression that includes all the six eco-climatic  
4 predictors. Two additional parameters are the intercept and the coefficient for land coverage.  
5 Because collinearity can inflate the standard error of regression coefficient, we also conduct  
6 likelihood ratio ( $LR$ ) tests to assess if adding a predictor significantly increases model fit. If so, the  
7 predictor has a significant effect on language diversity beyond its covariation with other predictors.  
8 Significant results are in bold.

9

Predictor	Low ( $n=216$ )			Medium ( $n=192$ )			High ( $n=366$ )		
	$t$	$p$	$LR$	$t$	$p$	$LR$	$t$	$p$	$LR$
Annual mean precipitation	0.56	0.577	0.32	1.32	0.190	1.67	1.74	0.083	1.27
Annual mean temperature	0.04	0.968	0.00	-0.90	0.369	0.69	1.15	0.251	0.80
Precipitation seasonality	<b>2.15</b>	<b>0.032</b>	<b>4.54</b>	1.54	0.126	2.01	0.10	0.920	0.55
Temperature seasonality	-1.14	0.257	1.30	<b>-2.43</b>	<b>0.016</b>	<b>4.76</b>	<b>-2.24</b>	<b>0.026</b>	<b>4.12</b>
Net primary productivity	1.86	0.064	3.54	1.65	0.101	2.69	1.58	0.114	3.38
Mean annual growing season	1.29	0.199	1.72	1.06	0.290	1.16	0.85	0.395	0.20

10

11

12

1 **Table 2.** Predictions of the ecological risk hypothesis ( $n$  is the number of grid cells used in the  
 2 analysis at each resolution). We list the  $t$  value, its standard error and the  $p$  value of the predictor in  
 3 each generalized least squares regression for the response variable. Each model includes three  
 4 parameters: intercept, coefficient of land coverage, and coefficient of the predictor. Significant  
 5 results are in bold. To test if mean growing season shows significant association with the response  
 6 variable beyond its covariation with the predictor, we conduct a likelihood ratio test on whether  
 7 adding mean growing season as an additional predictor significantly increase the model fit.

8  
 9

	Response	Predictor	$t$	$p$	Does mean growing season show significant association with response beyond its covariation with predictor?
Low resolution ( $n=216$ )	Language diversity	Mean growing season	<b>3.96</b>	<b>&lt;0.001</b>	NA
	Average population size	Mean growing season	-1.02	0.309	NA
	Min. population size	Mean growing season	-1.72	0.087	NA
	Population density	Mean growing season	<b>3.22</b>	<b>0.002</b>	NA
	Language diversity	Latitude	<b>-2.58</b>	<b>0.011</b>	Yes: $LR = 14.00, p < 0.001$
	Language diversity	Population density	-0.59	0.556	Yes: $LR = 15.85, p < 0.001$
Medium resolution ( $n=192$ )	Language diversity	Mean growing season	<b>4.70</b>	<b>&lt;0.001</b>	NA
	Average population size	Mean growing season	-1.92	0.056	NA
	Min. population size	Mean growing season	<b>-3.31</b>	<b>0.001</b>	NA
	Population density	Mean growing season	<b>3.97</b>	<b>&lt;0.001</b>	NA
	Language diversity	Latitude	<b>-2.49</b>	<b>0.014</b>	Yes: $LR = 15.38, p < 0.001$
	Language diversity	Population density	-0.30	0.763	Yes: $LR = 18.53, p < 0.001$



---

	Language diversity	Mean growing season	<b>5.57</b>	<b>&lt;0.001</b>	NA
High resolution ( $n=366$ )	Average population size	Mean growing season	-1.83	0.068	NA
	Min. population size	Mean growing season	<b>-3.13</b>	<b>0.002</b>	NA
	Population density	Mean growing season	<b>4.83</b>	<b>&lt;0.001</b>	NA
	Language diversity	Latitude	<b>-5.22</b>	<b>&lt;0.001</b>	Yes: $LR = 29.68, p < 0.001$
	Language diversity	Population density	0.66	0.512	Yes: $LR = 29.07, p < 0.001$

---

1

2

1 **Table 3.** Landscape effects on language diversity and speaker population size after accounting for  
 2 their covariation with climatic variables, at high, medium, and low resolution ( $n$  is the number of  
 3 grid cells used in the analysis at each resolution). Models with population size also control for  
 4 population density in addition to climatic variables. We list the  $t$  value and the  $p$  value of each  
 5 landscape variable in a generalized least squares regression that includes all the six climatic and  
 6 four landscape variables. Two additional parameters are the intercept and the coefficient for land  
 7 coverage. We also conduct likelihood ratio test to test if adding a landscape variable significantly  
 8 increases model fit. If so, the variable has a significant effect on language diversity beyond its  
 9 covariation with climatic variables and the other landscape variables. Significant results are in bold.  
 10

Response	Predictor	Low ( $n=216$ )			Medium ( $n=192$ )			High ( $n=366$ )		
		$t$	$p$	$LR$	$t$	$p$	$LR$	$t$	$p$	$LR$
	Average altitude	-1.21	0.228	1.16	0.70	0.483	0.52	0.21	0.830	0.92
Language diversity	Altitudinal range	1.67	0.096	2.93	1.01	0.312	1.09	<b>1.98</b>	<b>0.049</b>	<b>4.03</b>
	Landscape roughness	0.13	0.900	0.02	-0.00	0.999	0.00	0.98	0.328	0.99
	River density	<b>3.02</b>	<b>0.003</b>	<b>9.22</b>	<b>2.44</b>	<b>0.016</b>	<b>6.21</b>	1.42	0.157	2.07
Average speaker population size	Average altitude	0.44	0.662	0.20	0.74	0.463	0.56	0.73	0.463	0.55
	Altitudinal range	1.21	0.229	1.44	-0.65	0.515	0.44	-1.20	0.232	1.45
	Landscape roughness	-0.82	0.412	0.70	1.24	0.215	1.59	1.94	0.053	3.76
	River density	-1.02	0.310	1.06	-1.20	0.234	1.44	-0.47	0.641	0.21
Minimum speaker population size	Average altitude	0.50	0.616	0.27	0.81	0.420	0.62	1.20	0.232	1.46
	Altitudinal range	0.22	0.824	0.05	-1.19	0.238	1.22	-1.74	0.083	3.06
	Landscape roughness	-1.25	0.214	1.64	1.57	0.118	2.35	1.04	0.301	1.08
	River density	-0.43	0.668	0.20	<b>-2.24</b>	<b>0.026</b>	<b>5.04</b>	-1.32	0.187	1.68

11

12

1

2 **Table 4.** Association between biodiversity and language diversity after accounting for their  
3 covariation with all the climatic and landscape variables, at low, medium and high resolution ( $n$  is  
4 the number of grid cells used in the analysis at each resolution). We list the  $t$  value and the  $p$  value  
5 of a biodiversity variable in a generalized least squares regression that includes the biodiversity  
6 variable and all the six climatic and four landscape variables. Two additional parameters are the  
7 intercept and the coefficient for land coverage. We also conduct likelihood ratio (LR) test to test if  
8 adding the biodiversity variable significantly increases model fit. Significant results are in bold.

9

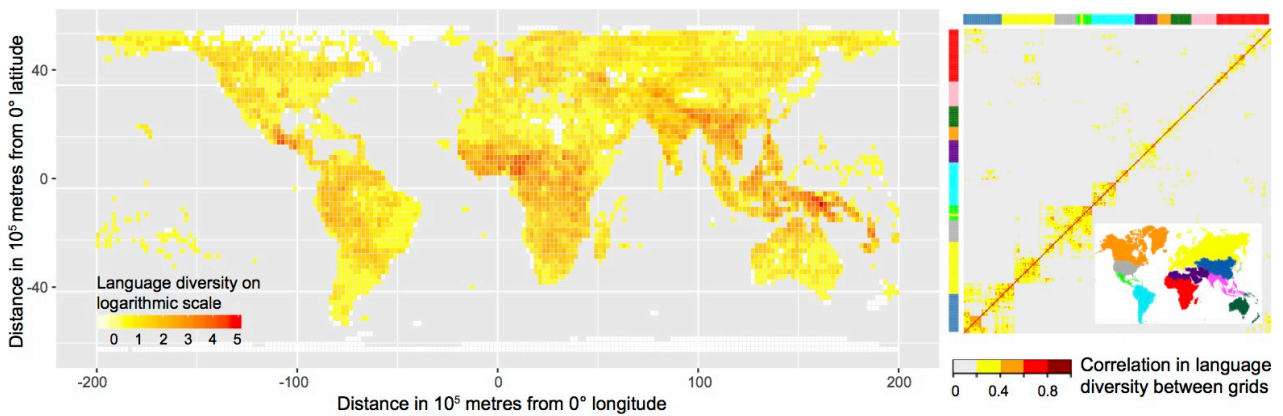
Biodiversity	Low ( $n=216$ )			Medium ( $n=192$ )			High ( $n=366$ )		
	$t$	$p$	$LR$	$t$	$p$	$LR$	$t$	$p$	$LR$
Plant diversity	1.32	0.188	1.78	0.76	0.449	0.60	0.28	0.783	0.10
Amphibian diversity	0.38	0.705	0.14	-0.67	0.507	0.41	-0.30	0.761	0.10
Mammal diversity	<b>3.12</b>	<b>0.002</b>	<b>9.09</b>	1.79	0.075	3.00	<b>2.23</b>	<b>0.027</b>	<b>11.32</b>
Bird diversity	<b>3.74</b>	<b>&lt;0.001</b>	<b>12.23</b>	<b>2.66</b>	<b>0.009</b>	<b>6.65</b>	<b>2.99</b>	<b>0.003</b>	<b>15.37</b>

10

11

1 **Figure 1.** Left panel: Global distribution of language diversity. Values on a logarithmic scale of  
2 number of languages are shown for 200x200km cells of an equal-area grid. Right panel: Correlation  
3 in language diversity between each pair of grid cells due to spatial autocorrelation and phylogenetic  
4 relatedness. Correlation coefficient is estimated from our generalized least squares model that  
5 includes all the climatic and landscape variables as predictors (see Methods). Correlated grid cells  
6 are roughly clustered into 9 geographic regions, so we colour code the rows and columns by these  
7 regions. Grid cells within East Asia, Europe, and the Americas are more autocorrelated than grid  
8 cells within the other regions.

9

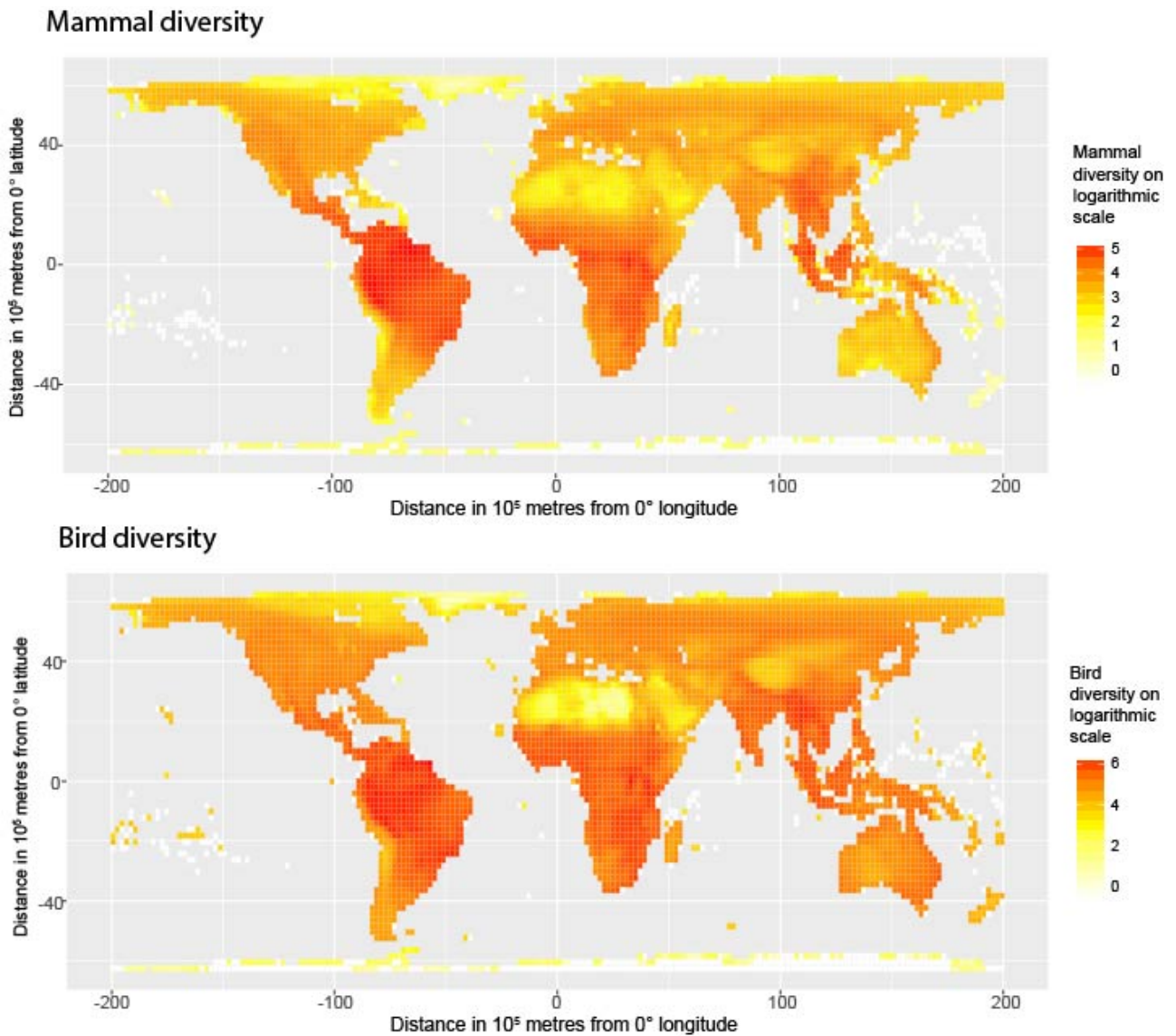


10

11

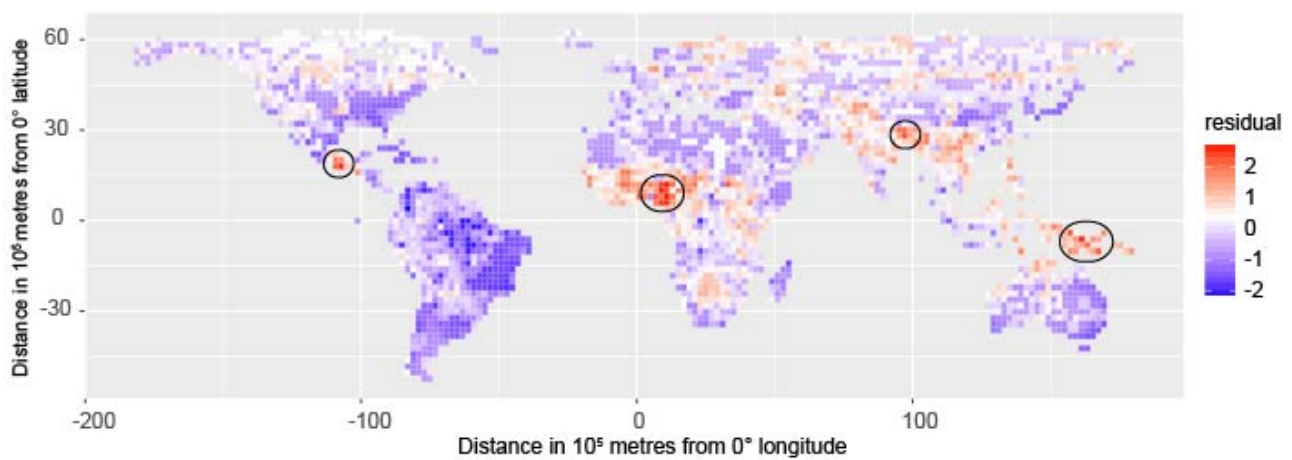
1 **Figure 2.** Global distribution of mammal diversity and bird diversity. Values on logarithm scale of  
2 number of species are shown for 200x200km cells of an equal-area grid. For amphibian and plant  
3 diversity see figure S3.

4



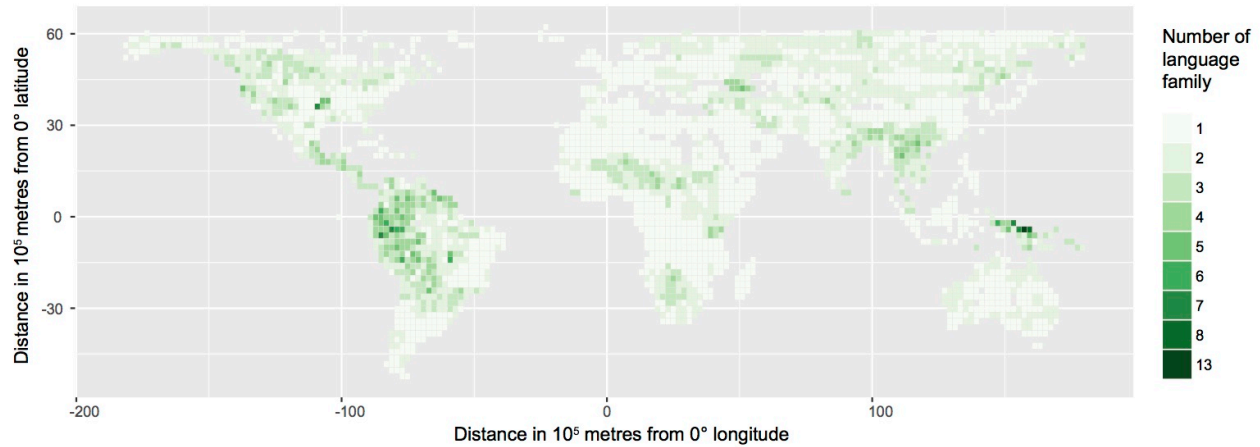
1 **Figure 3.** Global distribution of residuals in language diversity within 200x200km grid cells, after  
2 accounting for the climatic and landscape effects on language diversity across all grid cells.  
3 Aggregations of grid cells with residuals  $\geq 1.96$  (red) are circled. These indicate four regions of  
4 higher than expected language diversity, compared to regions of similar climate and landscape  
5 (New Guinea, the Himalaya foothills of Assam and Bangladesh, West Africa, and Mesoamerica).  
6 Areas of lower than expected language diversity with residuals  $\leq -1.96$  (blue) are distributed in  
7 South America, mostly in the Amazon basin.

8



1 **Figure 4.** Global distribution of the number of language families within 200x200km grid cells.  
2 Language family is defined by the World Language Mapping System taxonomy<sup>1</sup>. Language isolates  
3 are treated as distinct families. Number of language families within a grid cell is calculated as the  
4 number of language families that include at least one language distributed in the grid cell.

5



6

7