# nature portfolio

## Peer Review File

Inferring language dispersal patterns with velocity field estimation

Reviewers' Comments:

Reviewer #1:
Remarks to the Author:

The authors propose a vectorial framework to reconstruct the spatial dispersal of four language families around the world. The authors use a very wide range of methods that are borrowed from data science, physics and others from linguistics. I do not have the expertise to cover all of these methods, however the authors could help the reader understand if these methods are clustering algorithms, prediction methods, accuracy tests, etc. Some methods are called in the main text without further description, while some others are wrongly described, e.g. PCA is described in the main text as a similarity or clustering algorithm, actually PCA helps filtering out the least important features in order to describe a target variable in a space defined by superposition of few important features.
The methods section is a repetition of the vague description of the tools made in the main text and no further information is provided. The reader needs to get to Supplementary Information #3 to finally get a technical description of the methods that should actually appear in the Methods section. However, here the technical details are not clearly expressed and the physical meaning of the vector is unclear. Due to this, all the following results are unclear.
The text is hard to read, mostly due the presence of many typos and other grammar issues. Long sentences are used for speculative purposes, while key methodological descriptions are narrowed down to few vague sentences.
I realize that the authors did a very hard work and that the storytelling is not easy to unroll in a linear way. Still, I feel that the authors should make an effort to simplify, correct and make the text clearer in order to be readable by an interdisciplinary audience.

Here is a list of concerns:
- typo in the abstract, the sentence "And its effectiveness and robustness have been carefully verified by both simulated and empirical validations" starts with "and".
- line 87: again the sentence starts with "and"
- "And such relatedness could vary with time when languages continuously dispersal into new regions." sentence starts with "and" + dispersal is a noun, the verb is disperse. The same is repeated in many other sentences, please correct.
- line 106: "The Principal Component Analysis (PCA) is implemented to exhibit the linguistic relatedness of present languages." it is not clear on what variables the PCA is implemented. PCA identifies the most important variables to explain the variance of a target variable (in this case, I guess, the target variable is the languages relatedness?). Clustering classification is a forthcoming step.
- subsection "Simulated validations for language velocity field estimation". I really struggle here to understand what data did the authors use to validate their results. The dataset that is supposedly used as ground truth is also simulated by a phylogeographic algorithm. The authors claimed in the introduction that this method only captures vertical dependency of languages and not horizontal contacts and borrowings. I am confused about what is the contribution of this validation. Maybe the authors could add this discussion in the limitations of the study.
- I would avoid the usage of the word "true", unless there are striking evidences of the coordinates of the language dispersal origin.
- what is the delta score of tree-likeness?
- the authors do not describe the data they used accurately. For instance, what is a trait? What is a cognate? It is never stated.
- "Third, the changes in the state frequencies of linguistic traits are proportional to their sociolinguistic prestige in a certain area.". I don't get the logic of this sentence. What is the meaning of prestige here? The definition of prestige is expressed only in the next paragraph, it should be introduced before going into interpretations.
- "It is noted that the larger length of the velocity vector of a language denotes the more rapid change of this language during its evolution". The reader is provided with no tools to understand this

sentence. A schematic representation of a vector could really help. E.g. what are the elements of a vector?
- Does the PCA find only two components, or the authors found that more components did not lead to more variance explainability? Again, PCA here is presented as a tool to find similarities among datapoints, actually it is a rearrangements of the predictors of the model that tells what are the most important features in the model. The authors say nothing about all this. Projecting the points in to the PC space allows to visualize clusters, but actual clustering is performed by other tools, such as k-nearest-neighbors.
- it is not clear how the vectors are formed in the PC space. Up to my understanding the PCA describes the datapoint with two components, hence I expect to observe a single point with coordinates (PC1,PC2) in the PC space. By the way, we cannot build a vector with one point. I understand from SI-3 that the vectors are computed as the difference in the PC space of $X(0) - X(-m)$, where $t=0$ represents now and $t=-m$ represents a moment in the past. What is this moment in the past? Then I read "Therefore, Vl describes the change of the state frequencies of language l in a unit of time.". what is the unit of time? Years, centuries?
- what is the delta score and how is it computed? It is never stated in the text, nor in the SI
- Later on I read "In this study, we set $m = 1$.", but no reason is given, nor the unit of time is stated. One year? One century? Again, this is very opaque. I do not understand the physical meaning of this vectorial framework because no clear explanation is provided.
- the authors said that they study the spatial dispersal of languages along 10,000 years, to my understanding the vector field describes the change of the language between one exact moment of the past and $t=0$, which is supposed to be today.


Reviewer #2:
Remarks to the Author:
As I stated in my previous reviews of this paper, it is interesting, convincing, and historically significant in its conclusions. I am pleased to see that the authors have cut down the paper to deal with the four clearest examples, these being Indo-European, Sino-Tibetan, Bantu, and Arawak. The more troublesome Austroasiatic, Japonic and Oceanic examples have been removed, and I think this decision has added greatly to the clarity of the paper. It deserves to be published in Nature Communications.

My first comment is that the paper still needs a light level of English editing. I do not have time to do this on behalf of the authors, but perhaps I can use the abstract as an example of how some light editing might increase its clarity:

Here is the original abstract:

Reconstructing the spatial evolution of worldwide languages could shed light on understanding the global demic diffusions and cultural spreads. The phylogeographic approaches have been frequently used to infer the dispersal patterns of languages. However, they have shown some limitations primarily because the phylogenetic tree cannot properly capture the complex socio-cultural scenarios like contact-induced borrowings and areal diffusions of languages. Here, we introduced the language velocity field, which could be estimated directly from linguistic data without phylogenetic reconstruction, to enable the inference of the dispersal routes and centers of language families and groups in the geographic space. And its effectiveness and robustness have been carefully verified by both simulated and empirical validations. With the language
38 velocity field estimation, we made inferences on the dispersal patterns of four language families and groups worldwide including around 700 languages. Our results showed that the dispersal routes of these languages were primarily compatible with the population activities inferred from ancient DNA and archaeological materials, and their dispersal centers were geographically proximate to the ancient homelands of agricultural or Neolithic cultures. Our findings highlight that the agricultural languages

dispersed along with demic diffusions and cultural spreads globally in the past 10,000 years. We expect that language velocity field estimation could greatly aid the spatial analysis of language evolution, and many more studies of demographic and cultural dynamics.

And here is how I would edit it:

Reconstructing the spatial evolution of languages worldwide can shed light on understanding global demic diffusions and cultural spreads. The phylogeographic approaches that have been frequently used to infer the dispersal patterns of languages show limitations, primarily because a phylogenetic tree cannot properly capture complex socio-cultural scenarios that involved contact-induced borrowing and areal diffusion of languages. Here, we introduce the language velocity field, which can be estimated directly from linguistic data without phylogenetic reconstruction, as a resource that can enable the inference of the dispersal routes and centers of language families and groups in geographic space. Its effectiveness and robustness have been carefully verified by both simulated and empirical validations. Using language velocity field estimations, we infer the dispersal patterns of four language families and groups worldwide, covering around 700 languages. Our results show that the dispersal routes of these languages were primarily compatible with human population spreads inferred from ancient DNA and archaeological materials, and their dispersal centers were geographically proximate to ancient homelands of agricultural (or Neolithic) cultures. Our findings highlight that agricultural languages dispersed with demic diffusions and cultural spreads on a global scale during the past 10,000 years. We expect that language velocity field estimation will aid greatly the spatial analysis of language evolution, with implications for studies of demographic and cultural dynamics.

Back to my commentary:

Figure 2 shows the proposed agricultural homeland in northern Amazonia for Arawak. This conflicts with text lines 184-186, where it is stated that " In addition, the language velocity field posited the dispersal of Arawak languages originated from the border of Peru, Brazil, and Bolivia in Western Amazonia, which was geographically close to the known ancient agricultural homeland of South America in the Andes". This statement implies a homeland much further to the south than shown on the map, which is what the archaeology would suggest. The map shows an area too far north. I note in Supplementary Notes 1 Table S2 that the Arawak homeland is put in the northern lowlands of Bolivia (upper Madeira River), which is precisely where I would expect it to be!

Likewise, lines 187-189 state " Moreover, in the case of Sino-Tibetan languages, their dispersal center was inferred in the Gansu province of China (Figure 2b). It was approximate to the geographic ranges of the Yangshao (7,000-5,000 years BP) and/or Majiayao (5,500-4,000 years BP) Neolithic cultures, although it was far from the ancient agricultural homelands known in the Yangzi and Yellow River Basins of China." Surely, Yangshao and Majiayao were centrally located in the Yellow River homeland of millet and pig agriculture? I cannot understand what is meant here, although, of course, the Yangzi is a different matter.

The discussion from lines 197 to 298 is highly technical, and I have no observations on it. Much the same applies to the materials and methods section. I can understand from lines 301-9 that the basic data come from a geographical plotting of cognate presences and absences, but I was puzzled by the statement (lines 304-6) "Lexical cognates of these language samples in each language family or group were binary-coded traits..." This sentence seems to confuse the concepts of cognate and language. How many cognate terms were used in the analysis, and from which proto-language levels were these cognates derived? In other words, how was a cognate defined? This might be explained in the supplementary data, but I think it should be clearer here in the main text.

Lines 449-40 state: "The diversity approach is an alternative phylogenetic tree-free approach and simply infers the location of the language homeland to the areas with the highest linguistic diversity." What is meant here by linguistic diversity? Does it relate to relative times of splitting from an inferred

phylogenetic family tree? (i.e., deeper-splitting subgroups are older)? I presume it is not simply related to number of languages.

I noticed in Supplementary Note 1 that phylogenetic discussions of Austroasiatic, Japonic and Oceanic are still mentioned, even through these groupings are no longer discussed in the main text.

Supplementary Notes 2: it is not clear to me that Supplementary sections 2 and 3 are really necessary (The interdisciplinary alignment of Genetics, Archaeology, and Linguistics; The Age-Area Hypothesis for inferring the language homeland). I think the observations made in this paper can stand quite well without them.

Peter Bellwood


Reviewer #3:
Remarks to the Author:
I find this study generally quite interesting, since the authors claim that they have developed a new method that allows to represent historical dynamics of individual languages in comparison with neighboring languages by multidimensional vectors, which can then be projected in lower-dimensional space in order to even infer the original locations from which the language family as a whole dispersed.

While interesting, I see some general problems with the study, mainly its fit with the journal where it was submitted to, and as a result, I recommend it to be rejected -- not because it is too low in quality, but rather because it is not a good fit with the journal, as I'll explain below.

Apart from this, I see some major and minor flaws, which I'll discuss below.

First, regarding the fit of the approach: What the authors propose is a methodological study, a new methodology of which they claim it outperforms established -- albeit controversial -- methods. In such a case, the journal where they submitted their study to, does not really qualify as a good fit, since we do not deal with new findings (they cannot be made until the method has been thoroughly evaluated) but rather with a new method that needs to be shown to work. For this reason, I think some journal like "Nature Methods" would be a much better fit here.

Second, if the authors accept that they need to convince us first that their method is useful and will enlarge our future knowledge about the spread of language families over time, they should please provide their method in a way that it can be replicated. As of now, we have a bunch of unrelated, badly documented R-scripts in a folder of 600 MB, that are hard to read and even harder to understand. Where is the vector estimation happening, what is the k you choose for the k-means languages that you select as neighbors, what is the impact of k on your results, etc. It makes me extremely nervous to see such a huge bunch of barely commented R-scripts that often do the same, but bear another name of another language family. This is definitely not how you make a new method successful. The least we would expect is a package in R with a tutorial that runs us through your code, for one language family, and then an extended tutorial with all four language families.

Third, speaking of four, I hate to say this, but I was reviewing this study before, not negatively, but pointing to the code, and to other issues. Interestingly, the number of language families has now dropped from 7 to 4. How the hack did that happen? How do the authors explain that they discard three language families now? I know having the same reviewers for the same paper across journals is annoying, but please, good scientific practice requires you to be transparent and tell us what happened here. Did you discard them, because they did not bring the results you hoped for?

Fourth, the claim of the method not using phylogenetic information is a bit exaggerated: we know geography correlates often with language relatedness (see for example here: https://doi.org/10.1371/journal.pone.0265460), so if geography explains the tree, you cannot say you do not use the tree if you use geography as a proxy for the construction of your vectors.

Fifth, the question of homeland has always been problematic, but if you already use data by Wichmann and Rama, you should also check the much simpler baseline published in Glottolog by now (www.pyglottolog. readthedocs.io/en/latest/homelands.html#module-pyglottolog.homelands). This method seems to work as well as the one by Wichmann and Rama, but it is even simpler, so I would say there's one more baseline to be tested. And when speaking of testing: why restrict your study to four datasets (or seven), if there are many more available in terms of phylogenies now, which are all with nicely coded cognate sets in standardized data formats (see e.g., https://doi.org/10.1038/s41597-022-01432-0 for a very large collection of standardized data)? It seems the data has been cherry-picked to yield good results. Taking ten of the datasets in the Lexibank collection should not be difficult and would tell us much more clearly where we are with this new method.

Sixth, the method has the rather infelicitous name "language velocity field estimation", and I could not find any explanation why the authors chose to call it like that, since the name is very confusion and difficult to parse, and it does not really help to understand what the method could be about. I think in general it would be useful to 1) change the name to something that explains the method in a better way (dynamic trait vectors? I am not sure) and 2) to explain the method in much, much more detail. For this, figures would be needed that show how vectors for some of the traits are estimated, and the authors would need to also check the resulting vectors on an individual basis in order to see if they make sense.

Seventh, the authors praise their method for not needing trees, but at the same time, they do not tell the readers why trees are so useful: they tell us various scenarios of character evolution in a very transparent way, in which we have scenario and can plot how the trait evolved. Of course, this is not always done, but they should tell the readers to which the method they propose allows us to get some insights into the black box, since a simple black box, even if it works, is not satisfying from a scientific viewpoint, and we talk about scientific approaches here.

Eighth, and final point, the paper is not nice to read, the authors should check their wordings, which are often hard to follow, at times with flaws in grammar, and it would really profit from a complete overhaul and a thorough checking by a proof reader.

Due to all these reservations, I recommend that the paper be rejected, but I emphasize that it is not for poor quality, but for lack of fit. I look forward to see a new methods paper emerging from this, in which the authors work hard to share a useful new approach with the scientific world that they also evaluate rigorously against existing approaches. I am convinced they have the potential to turn their paper into such a study, and I am also very confident that this would be the right way to go, instead of trying to sell this as some study with new insights, or a study with a method that beats all existing approaches, since this is obviously not the case.

# Response Letter to Reviewers

**Replies to Reviewer 1:**

*Q1: The authors propose a vectorial framework to reconstruct the spatial dispersal of four language families around the world. The authors use a very wide range of methods that are borrowed from data science, physics and others from linguistics. I do not have the expertise to cover all of these methods, however the authors could help the reader understand if these methods are clustering algorithms, prediction methods, accuracy tests, etc. Some methods are called in the main text without further description, while some others are wrongly described, e.g. PCA is described in the main text as a similarity or clustering algorithm, actually PCA helps filtering out the least important features in order to describe a target variable in a space defined by superposition of few important features.*

**Replies to Q1:**

We sincerely appreciate the invaluable suggestions provided by the reviewer. Our computational approach can be characterized as a kind of spatial reconstruction method that primarily encompasses other three distinct methods. The first one is the Principal Component Analysis (PCA) which is an unsupervised dimensionality reduction technique for rearranging linguistic traits into fewer more important new traits. The second one is the dynamic model consisting of ordinary differential equations for reconstructing the past states of linguistic traits. The third one is the geographic projection technique utilized for mapping the velocity vectors from the PC space into the geographic space. In the revised manuscript, we have modified unclear and problematic descriptions of our approaches and added more corresponding comprehensive explanations (*Lines 110-151* of the revised main text).

The reviewer has pointed out: "*PCA is described in the main text as a similarity or clustering algorithm, actually PCA helps filtering out the least important features in order to describe a target variable in a space defined by superposition of few important features*". We are sorry for the imprecise descriptions of the PCA algorithm in the previous version of our manuscript. In this study, the PCA algorithm is not implemented to cluster language samples. Instead, it is used to reduce the dimension of linguistic traits by reassembling them into two important new traits (i.e., PC1 and PC2). Accordingly, each language sample can be visualized in the two-dimensional

PC space based on its PC1 and PC2 values. The Euclidean distances among pair-wise language samples in the PC space (i.e., PCA-based distance) represent their linguistic relatedness with each other. To be specific, the language samples sharing closer linguistic relatedness tend to distribute closer in the PC space. Therefore, the linguistic relatedness can be shown through the Euclidean distances among the language samples in the PC space.

It is noted that utilizing the PCA-based distance metric to assess sample relatedness is a prevailing practice in many studies within the fields of genetics and linguistics [1-3]. Accordingly, we employ the PCA-based distance to quantify the linguistic relatedness among language samples in this study. Following the reviewer's comments, we have revised all the contents related to the PCA algorithm (*Lines 114-122* of the main text).

**Reference**

[1] Wang, Chuan-Chao, et al. "Genomic insights into the formation of human populations in East Asia." Nature 591.7850 (2021): 413-419.

[2] Haak, Wolfgang, et al. "Massive migration from the steppe was a source for Indo-European languages in Europe." Nature 522.7555 (2015): 207-211.

[3] Norvik, Miina, et al. "Uralic typology in the light of a new comprehensive dataset." Journal of Uralic Linguistics 1.1 (2022): 4-42.
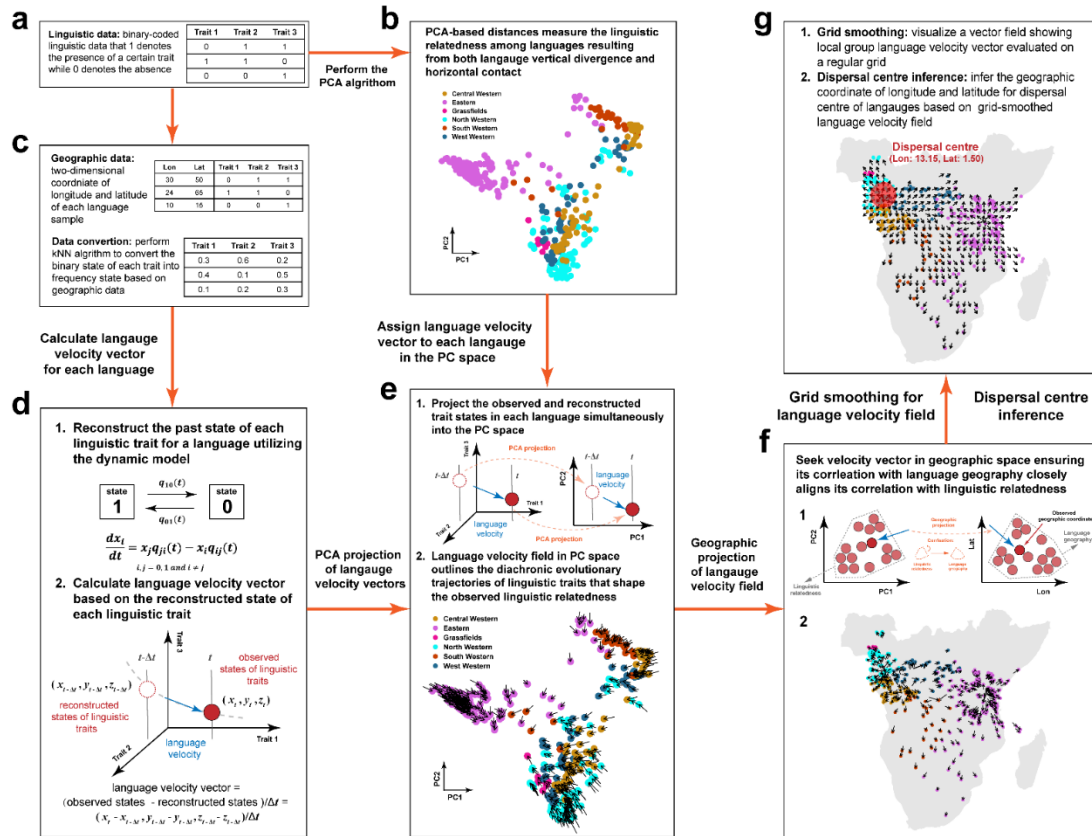
*Q2: The methods section is a repetition of the vague description of the tools made in the main text and no further information is provided. The reader needs to get to Supplementary Information #3 to finally get a technical description of the methods that should actually appear in the Methods section. However, here the technical details are not clearly expressed and the physical meaning of the vector is unclear. Due to this, all the following results are unclear.*

**Replies to Q2:**

We appreciate these comments. Following the reviewer's comments, we have rephrased some vague descriptions of our approach and added more technical

descriptions and key mathematical formulas in the Materials and Methods section. Considering the readability of the manuscript, detailed mathematical derivations and professional mathematical terminology descriptions have still been retained in Supplementary Note 3. Moreover, we have also provided a new schematic diagram (Figure 1 in the revised manuscript) to illustrate the rationale and procedure of our approach comprehensively. For the convenience of the reviewer, this figure is attached below namely Figure to Q2.



**Figure to Q2. Schematic diagram of language velocity field estimation (LVF) for inferring the dispersal trajectories and centers of languages.** The computational procedures of the LVF comprise two major steps. Subfigures (a) to (e) illustrate the first step which is to estimate a velocity field on the PC space to outline the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. Subfigures (f) to (g) illustrate the second step, which is to project the velocity field from PC space into geographic space. Within the velocity field in geographic space, the directions of the velocity vectors compose a set of continuously changing trajectories that delineate from where these languages diffuse to their current locations. These procedures are exemplified using the Bantu language family.

80  Comprehensive insights into the underlying principles and computational steps can be
81  found in the Materials and Methods section, as well as Supplementary Note 1.

82

83  *Q3：The text is hard to read, mostly due the presence of many typos and other*
84  *grammar issues. Long sentences are used for speculative purposes, while key*
85  *methodological descriptions are narrowed down to few vague sentences.*

86  **Replies to Q3:**

87  In the revised manuscript, we corrected the typos and grammar errors and
88  modified several long and vague sentences. Furthermore, we engaged the AJE
89  language editing service to thoroughly polish our manuscript (ID: Q2K9ZRSF). To
90  make our methodological description clearer, we rephrased some vague sentences and
91  added detailed mathematical formulas and explanations for our approach in the
92  Materials and Methods section.

93

94  *Q4：I realize that the authors did a very hard work and that the storytelling is not*
95  *easy to unroll in a linear way. Still, I feel that the authors should make an effort to*
96  *simplify, correct and make the text clearer in order to be readable by an*
97  *interdisciplinary audience.*

98  **Replies to Q4:**

99  We sincerely appreciate the reviewer's comments. Considering the readability of
100  the interdisciplinary audience, we have rephrased the sentences in the manuscript to
101  enhance the clarity and comprehensibility of the narrative. Moreover, we have
102  rearranged the structure of our whole manuscript to improve its clarity and readability.

103

104 *Q5：typo in the abstract, the sentence "And its effectiveness and robustness have*
105 *been carefully verified by both simulated and empirical validations" starts with*
106 *"and".*

107 **Replies to Q5:**

108     We greatly thank the reviewer for pointing this out. We have corrected this typo
109 in the abstract as shown in the *Line 35* of the revised main text.

110

111 *Q6: line 87: again the sentence starts with "and"*

112 **Replies to Q6**

113     This typo has been corrected in the revision.

114

115 *Q7: "And such relatedness could vary with time when languages continuously*
116 *dispersal into new regions." sentence starts with "and" + dispersal is a noun, the*
117 *verb is disperse. The same is repeated in many other sentences, please correct.*

118 **Replies to Q7**

119     These grammatical errors have been corrected in the revision.

120

121 *Q8: line 106: "The Principal Component Analysis (PCA) is implemented to exhibit*
122 *the linguistic relatedness of present languages." it is not clear on what variables the*
123 *PCA is implemented. PCA identifies the most important variables to explain the*
124 *variance of a target variable (in this case, I guess, the target variable is the*
125 *languages relatedness?). Clustering classification is a forthcoming step.*

126 **Replies to Q8:**

127     We appreciate these important comments. In our study, Principal Component
128 Analysis (PCA) has been applied to the binary-coded lexical trait, where the value 1

129 indicates the presence of the lexical trait in a language, while 0 signifies its absence.
130 More specifically, our dataset is organized in the form of a matrix comprising binary
131 values. The rows of this matrix correspond to diverse language samples, while the
132 columns denote distinct binary-coded lexical traits, as illustrated in Table to Q8. In
133 this study, both the empirical and simulated datasets adhere to this form.

134 The target variables derived from the PCA process are not referred to as linguistic
135 relatedness. Instead, linguistic relatedness among language samples is represented by
136 their Euclidean distances within the PC space. Specifically, in this study, PCA is
137 employed to linearly transform lexical traits into two critical variables designated as
138 PC1 and PC2. These PC1 and PC2 variables are the target variables extracted by the
139 PCA algorithm. They represent the two most significant dimensions capable of
140 capturing the primary variations within the original linguistic traits. Consequently, we
141 can visually represent language samples based on their coordinates (PC1, PC2) within
142 a two-dimensional PC space. In this space, language samples with closer linguistic
143 relatedness are naturally distributed together. In such instances, the Euclidean
144 distances among language samples within the PC space serve as a manifestation of
145 their linguistic relatedness.

146 **Table to Q8.** The format of the linguistic dataset utilized in this study.

|  | **Trait 1** | **Trait 2** | **Trait 3** | **…** | **Trait $k$** |
|---|---|---|---|---|---|
| **Language 1** | 0 | 1 | 0 | … | 1 |
| **Language 2** | 1 | 0 | 1 | … | 1 |
| **…** | … | … | … | … | … |
| **Language $n$** | 1 | 1 | 0 | … | 0 |

147

148 *Q9: subsection "Simulated validations for language velocity field estimation". I*
149 *really struggle here to understand what data did the authors use to validate their*
150 *results. The dataset that is supposedly used as ground truth is also simulated by a*
151 *phylogeographic algorithm. The authors claimed in the introduction that this*
152 *method only captures vertical dependency of languages and not horizontal contacts*
153 *and borrowings. I am confused about what is the contribution of this validation.*
154 *Maybe the authors could add this discussion in the limitations of the study.*

155 **Replies to Q9:**

156    We are grateful for these comments. The reasons for utilizing the simulated
157 datasets in this study are given below:

## 1. Simulated datasets with known dispersal centers can be used for model validations

160    The optimal validation for our methodology should be implemented relying on
161 benchmark datasets where the actual language dispersal centers are already
162 documented. These datasets enable us to validate our approach by comparing the
163 estimated dispersal center locations with the documented ones. Since empirical
164 datasets often lack precise information on the actual dispersal center locations,
165 validating our approach using empirical datasets is challenging due to the credibility
166 of the estimated dispersal center is hard to verify. Fortunately, a viable solution is
167 provided by simulated datasets from Wichmann et al. (2021) [1]. These simulated
168 datasets include known locations of true language dispersal centers, as they are
169 generated through a random walk model applied to a phylogenetic tree assigned with
170 given dispersal centers. Given the locations of the language dispersal centers are
171 known in these simulated datasets, they can serve as robust benchmarks for validating
172 our approach. In the previous manuscript, we extensively demonstrated the
173 effectiveness and robustness of our approach based on these simulated datasets.

## 2. Simulated datasets are not generated by the phylogeographic approach

175    We would like to clarify that the simulated datasets are not generated through the
176 phylogeographic approach but the random walk model. We understand that the
177 unclear descriptions in the previous manuscript may have led the reviewer to consider
178 these two approaches are the same. However, the phylogeographic approach is just a
179 specific application of the random walk model in the phylogenetic domain [2-3]. The
180 phylogeographic approach aims to backwardly reconstruct the language dispersal
181 center based on the locations of observed language samples assigned to a
182 phylogenetic tree. In contrast, the random walk model utilized in Wichmann et al.
183 (2021) [1] is employed to forwardly generate the locations of observed language
184 samples based on a phylogenetic tree assigned with a given language dispersal center.
185 As mentioned in Wichmann et al. (2021), the generation of the simulated datasets
186 follows below procedures:

187    "…*The simulation process can be summarized as follows. Movements are*

*constrained to any populated place on Earth, i.e. a place included in the geonames.org database. A starting point is found by randomly choosing from this set of populated places. At each time step there is a preset probability of moving to a new place within a square containing at least ch populated places…….The kind of movement we simulate here may be called a semi-random walk, since it is a kind of random walk constrained to populated places…….Maps of all 1000 cases, showing the homeland, intermediate stations, locations of current languages, and inferred homelands similarly to Figure 2 below, as well as the script that produced the maps, are provided in the electronic supplementary material (SI-11)….*"*

Therefore, it is important to note that the simulated datasets are not produced through the phylogeographic approach, even though the simulation process incorporates the phylogenetic tree and random walk model.

## 3. Simulated datasets as benchmarks for model comparisons

**(i) Our approach and the phylogeographic approach share a common theoretical foundation but employ distinct implementation strategies.** Both our approach and phylogeographic approach involve two key steps in inferring language dispersal through the diachronic evolution of linguistic traits (Figure 1 to Q9). The first step entails delineating the diachronic evolutionary trajectories of linguistic traits that contribute to linguistic relatedness among observed language samples. The second step involves transforming these trajectories into language dispersal trajectories based on the correlation between linguistic relatedness and language geography [2, 4].

However, these two approaches differ in their detailed strategies for implementing these steps (Figure 1 to Q9). The primary distinctions revolve around how linguistic relatedness is represented. Specifically, in the phylogeographic approach, linguistic relatedness is represented by the phylogenetic tree, which captures only vertical language divergence. In contrast, our approach measures linguistic relatedness through the Euclidean distances among language samples in the two-dimensional PC space (PCA-based distance). This method can capture both vertical divergence and horizontal contact. We anticipate that our approach would perform similarly to the phylogeographic approach when linguistic relatedness can be explained by the tree model (Table to Q9). However, when linguistic relatedness cannot be fully explained by the tree model, there is a notable difference between the two approaches (Table to

221      To illustrate this, we conducted comprehensive simulated and empirical
222 comparisons between our approach and the phylogeographic approach. The results of
223 the comparisons are summarized in Figure 3 in the revised main text. For the
224 reviewer's convenience, we have attached this figure to this reply as Figure 2 to Q9. It
225 is important to note that Figure 2 to Q9 outlines the comparison results not only
226 between our approach and the phylogeographic approach but also against four other
227 spatial reconstruction approaches: the diversity approach, the minimal distance
228 approach, and the centroid approach. However, in this response, we focus solely on
229 the comparison between our approach and the phylogeographic approach to highlight
230 their similarities and differences (Figure 1 to Q9).

231 **(ii) Simulated comparisons when linguistic relatedness can be explained by the**
232 **tree model.** The simulated datasets can serve as benchmarks to compare the
233 performance between our approach and the phylogeographic approach when the
234 linguistic relatedness can be explained by the tree model. Due to simulated datasets
235 being generated based on a specific phylogenetic tree, the linguistic relatedness of the
236 simulated language samples is solely raised by the vertical divergence. In other words,
237 the linguistic relatedness among simulated language samples can be well captured by
238 the tree model. Therefore, based on the simulated datasets, the dispersal centers
239 inferred by the phylogeographic approach and our approach should be the same as
240 each other.

241      Fortunately, the simulated results indeed showed the same performance between
242 the phylogeographic approach and our approach ($p$-value $> 0.05$; Figures 2b1 to Q9).
243 More importantly, under the circumstance of the linguistic relatedness being solely
244 raised by vertical divergence, the phylogenetic tree and PCA-based distance
245 estimation can both adequately explain the linguistic relatedness ($p$-value $< 0.05$;
246 Figure 2b6 to Q9). It evidences that our approach and phylogeographic approach
247 indeed share the same theoretical foundation but with different implementations.

248 **(iii) Empirical comparisons using simulated results as baselines when linguistic**
249 **relatedness cannot be explained by the tree model.** The four empirical datasets can
250 be utilized for comparisons between our approach and the phylogeographic approach
251 when the linguistic relatedness cannot be explained by the tree model. Based on the

252    phylogenetic topology of simulated language samples as baseline (Figure 2b2 to Q9),

253    the phylogenetic topology of language samples in four empirical datasets utilized in

254    this study significantly deviates from the tree topology in this study ($p$-value < 0.05;

255    Figure 2b6 to Q9). It indicates that both vertical divergence and horizontal contact

256    could have contributed to the linguistic relatedness among these empirical language

257    samples. Accordingly, the phylogenetic tree cannot be able to adequately interpret the

258    linguistic relatedness within these four empirical cases. Under this circumstance, we

259    would anticipate different dispersal centers estimated by our approach and the

260    phylogeographic approach in empirical applications.

261         With the estimated difference in simulated comparisons as the baseline, the

262    empirical comparisons demonstrated a significant difference in performances between

263    our approach and the phylogeographic approach in Sino-Tibetan and Arawak ($p$-value

264    < 0.05; Figure 2a to Q9) languages. However, such difference was not observed in the

265    Bantu and Indo-European languages ($p$-value > 0.05; Figure 2a to Q9). The reason is

266    that for Bantu and Indo-European languages, PCA-based distance and phylogenetic

267    tree can both explain the linguistic relatedness among language samples ($p$-value <

268    0.05; Figure 2b6 to Q9). It indicates that the phylogenetic tree can explain the

269    linguistic relatedness under the influence of a certain degree of horizontal contact. In

270    contrast to Bantu and Indo-European languages, the comparison results showed that

271    PCA-based distance (Sino-Tibetan: $p$-value < 0.05; Arawak: $p$-value < 0.05; Figure

272    2b6 to Q9) could well explain the linguistic relatedness of Sino-Tibetan and Arawak

273    languages, while the phylogenetic tree cannot (Sino-Tibetan: $p$-value = 0.115; Arawak:

274    $p$-value = 0.121; Figure 2b6 to Q9). These empirical comparisons confirm that the

275    difference between our approach and the phylogeographic approach can be attributed

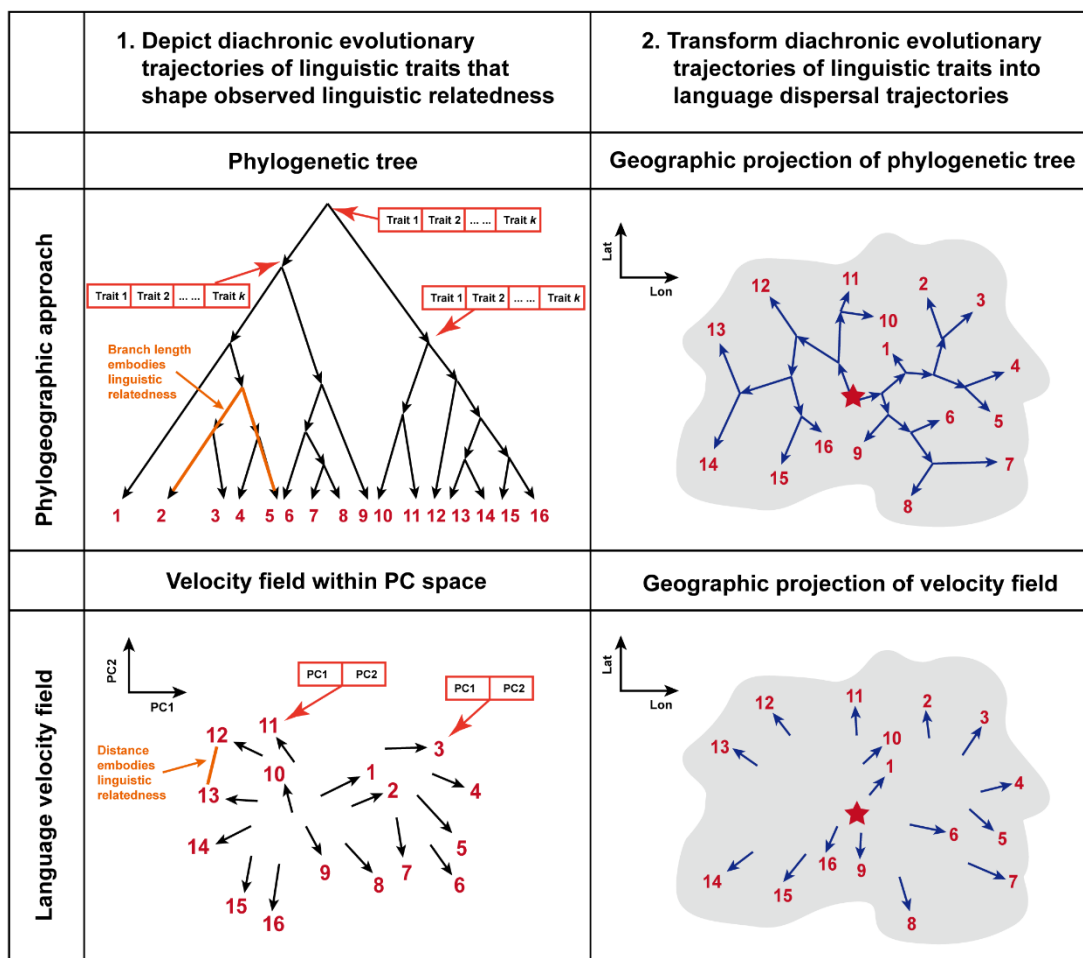276    to the distinct strategies for representing linguistic relatedness.

277    **In summary, the simulated and empirical comparisons confirm that the**

278    **distinction between our approach and the phylogeographic approach is raised by**

279    **their different explanatory power for linguistic relatedness. To be specific, when**

280    **linguistic relatedness can be explained by the family-tree model, the performance**

281    **between the phylogeographic approach and our approach is identical. However,**

282    **when linguistic relatedness cannot be explained by the family-tree model, a**

283    **notable distinction would emerge between the phylogeographic approach and**

284    **our approach.** In the revision, all the aforementioned contents have been added to the

285    revised main text as shown in the *Lines 153-172 and Lines 210-303*.

286 **Table to Q9.** Expected performance between the phylogeographic approach and our
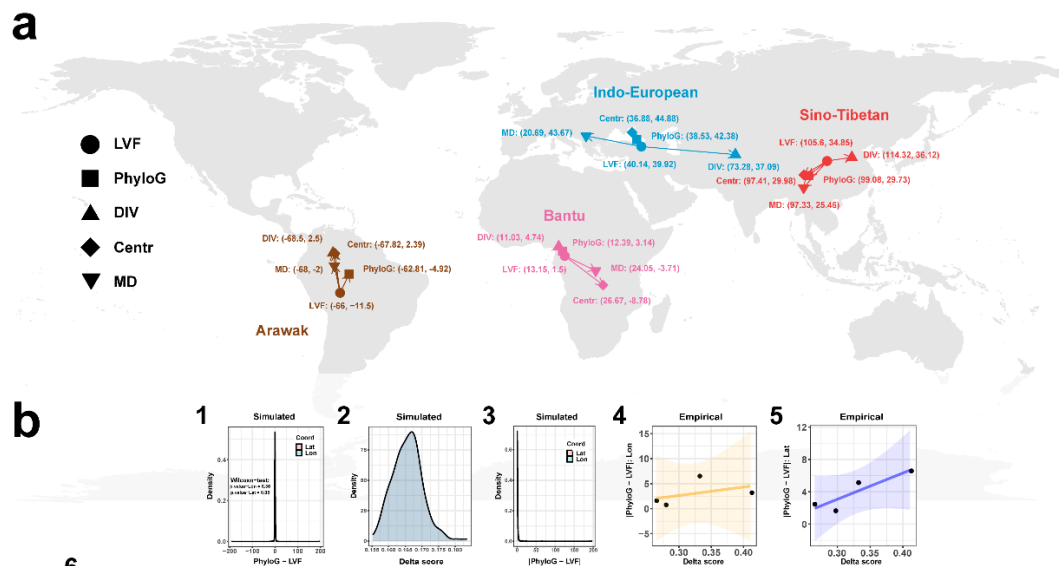287 approach utilizing simulated and empirical datasets.

| | | Simulated dataset | Empirical dataset | |
|---|---|---|---|---|
| **Linguistic relatedness attribution** | | Vertical divergence | Vertical divergence | Horizontal contact |
| **Whether the approaches can capture the divergence or contact** | **Phylogeographic approach** | √ | √ | ✕ |
| | **Language velocity field** | √ | √ | √ |
| **Equality of two approaches** | | = | ≠ | |

288



289

**Figure 1 to Q9. Language velocity field estimation (LVF) shares the same foundation as the phylogeographic approach but with different implementation strategies.** Both LVF and phylogeographic approach entails two major steps to infer language dispersal pattern. The first is to depict the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. The second is to transform these diachronic evolutionary trajectories of linguistic traits into language dispersal trajectories. In the phylogenetic tree, each language is determined by $k$ linguistic traits. In the velocity field within PC space, each language is determined by PC1 and PC2 which are rearranged from the $k$ linguistic traits through the PCA algorithm. The red number denotes a language. The black arrow signifies the evolutionary direction of linguistic traits in a language. The blue arrow represents the dispersal direction of a language. The red star denotes the estimated dispersal center.

| Tree-likeness | | Simulation | | Indo-European | | Sino-Tibetan | | Bantu | | Arawak | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Delta-score | 0.1657 | | 0.2656 | | 0.3324 | | 0.2976 | | 0.4129 | |
| | *p*-value | - | | < 0.01 | | < 0.01 | | < 0.01 | | < 0.01 | |
| | | Lon | Lat | Lon | Lat | Lon | Lat | Lon | Lat | Lon | Lat |
| \|PhyloG-LVF\| | Difference | 1.55 | 0.94 | 1.61 | 2.46 | 6.52 | 5.12 | 0.76 | 1.64 | 3.19 | 6.58 |
| | *p*-value | - | - | 0.158 | 0.069 | 0.020 | 0.012 | 0.443 | 0.142 | 0.058 | 0.007 |
| Linguistic relatedness explanatory power | | $R^2$ | *p*-val | $R^2$ | *p*-val | $R^2$ | *p*-val | $R^2$ | *p*-val | $R^2$ | *p*-val |
| | PCA-based distance | 0.90 | < 0.01 | 0.37 | < 0.01 | 0.44 | < 0.01 | 0.68 | < 0.01 | 0.53 | < 0.01 |
| | Phylogenetic tree | 0.93 | < 0.01 | 0.39 | < 0.01 | 0.05 | 0.146 | 0.31 | < 0.01 | 0.09 | 0.075 |

**Figure 2 to Q9. The comparison between LVF and other spatial reconstruction approaches.** a) The dispersal centres of four empirical language families and groups inferred by five different approaches: language velocity field estimation (LVF),

phylogeographic approach (PhyloG), diversity approach (DIV), centroid approach (Centr), and minimal distance approach (MD). b1) The density plot for the distribution of differences between the coordinates of dispersal centres in the aspects of longitude and latitude inferred from LVF and PhyloG based on 1,000 simulated datasets. The p-value is calculated based on the Wilcoxon rank-sum test, where < 0.05 indicates that the difference between the inferred coordinates is significantly different from zero. b2) The density plot for the average delta score of the languages whose linguistic relatedness can be well-explained by the tree model. It was estimated from 200 bootstrap replicates on the simulated languages. b3) The density plot for the distribution of the absolute differences in the aspects of longitude and latitude between the coordinates of dispersal centres inferred from LVF and PhyloG based on 1,000 simulated datasets. b4) The linear relation between the average delta score and the absolute difference of the longitude estimated from LVF and PhyloG. The orange ribbon denotes the 95% confidence interval. b5) The linear relation between the average delta score and the absolute difference of the latitude estimated from LVF and PhyloG. The blue ribbon denotes the 95% confidence interval. b6) The table of the delta score, estimated difference between LVF and PhyloG, and linguistic relatedness explanatory power of PCA-based distance estimation and phylogenetic tree. The p-value is calculated by the Wilcoxon rank-sum test where < 0.05 indicates the significance of the delta score, estimated difference, and linguistic relatedness explanatory power.

**Reference**

[1] Wichmann, Søren, and Taraka Rama. "Testing methods of linguistic homeland detection using synthetic data." Philosophical Transactions of the Royal Society B 376.1824 (2021): 20200202.

[2] Bouckaert, Remco, et al. "Mapping the origins and expansion of the Indo-European language family." Science 337.6097 (2012): 957-960.

[3] Grollemund, Rebecca, et al. "Bantu expansion shows that habitat alters the route and pace of human dispersals." Proceedings of the National Academy of Sciences 112.43 (2015): 13296-13301.

[4] Koile, Ezequiel, et al. "Geography and language divergence: The case of Andic languages." Plos one 17.5 (2022): e0265460.

339

**Replies to Q10:**

We appreciate this comment. As elucidated in our **Replies to Q9**, the most significant characteristic of the simulated datasets is that they are generated based on the given language dispersal centers. In other words, the actual locations of the dispersal centers are already known within the simulated datasets. Following the reviewer's suggestion, we have corrected the word "*true*" as "*given*" in the revised manuscript.

*Q11: what is the delta score of tree-likeness?*

**Replies to Q11:**

In the revision, we have added a comprehensive explanation of the delta score in *Lines 253-257* of the revised main text. Here, we provide a brief description. The delta score, denoted as $\delta$ score, serves as a widely used metric for quantifying the likeness between the language phylogenetic topology and the tree topology in the phylo-linguistics [1-3]. In other words, the delta score quantifies the degree of linguistic relatedness of languages that can be explained by the tree model. The delta score is calculated based on the distance among the languages, with a value ranging from 0 to 1. A larger value of the delta score denotes that the language phylogenetic topology is more compatible with the tree topology [4]. In other words, a larger value of the delta score signifies that the linguistic relatedness is less affected by the horizontal contacts and can be better explained by the tree model.

**Reference**

[1] Greenhill, Simon J., et al. "Evolutionary dynamics of language systems." Proceedings of the National Academy of Sciences 114.42 (2017): E8822-E8829.

[2] Kolipakam, Vishnupriya, et al. "A Bayesian phylogenetic study of the Dravidian

367    language family." Royal Society open science 5.3 (2018): 171504.

368    [3] Birchall, Joshua, Michael Dunn, and Simon J. Greenhill. "A combined
369    comparative and phylogenetic analysis of the Chapacuran language family."
370    International Journal of American Linguistics 82.3 (2016): 255-284.

371    [4] Holland, Barbara R., et al. "δ plots: a tool for analyzing phylogenetic distance
372    data." Molecular biology and evolution 19.12 (2002): 2051-2059.

373

374    *Q12: the authors do not describe the data they used accurately. For instance, what*
375    *is a trait? What is a cognate? It is never stated.*

376    **Replies to Q12:**

377    We thank the reviewer for pointing this out. In this study, our datasets contain the
378    Indo-European, Sino-Tibetan, Bantu, and Arawak lexical cognate datasets derived
379    from the previous publications respectively [1-4]. These datasets contain several
380    lexical words following a specific wordlist such as Swadesh 100 or 200 wordlist.
381    Each word (item) contains different lexical cognates identified by linguistic experts,
382    which manifest the same meaning and similar sounds. Furthermore, each cognate has
383    been transformed into a binary-coded lexical trait where the value of 1 denotes the
384    presence of this cognate in the language, while 0 indicates its absence (an example of
385    cognate coding is shown in Table to Q12). Accordingly, the Indo-European dataset
386    contains 5,995 binary lexical cognates across 103 language samples; the Sino-Tibetan
387    dataset encompasses 949 binary lexical cognates across 109 Sino-Tibetan language
388    samples; the Bantu dataset comprises 3,859 binary lexical cognates across 420
389    language samples; Arawak dataset involves 694 binary lexical cognates across 60
390    language samples. The detailed cognate coding process for each case is described as
391    follows.

392    For the Indo-European lexical dataset, Bouckaert et al. compiled 207 lexical
393    items [1]. According to these lexical items, they identified 5,995 cognates across 103
394    Indo-European languages, which were further recoded as 5,995 binary-coded lexical
395    traits. Bouckaert et al. described their cognate coding process as follows: "*We*
396    *recorded word forms and cognacy judgments across 207 meanings in 103*

15

397 *contemporary and ancient languages…. Cognate data were coded as binary*
398 *characters showing the presence or absence of a cognate set in a language. There*
399 *were 5995 cognate sets in total, with most meanings represented by several different*
400 *cognate sets. All cognate coding decisions were checked with published historical*
401 *linguistic sources (Table S1). The database contained 25908 cognate coded lexemes.*
402 *Of these, 67% came originally from ref. (17 ), 14% from ref. (16 ), and 19% were*
403 *newly compiled from published sources. Ref. (17 ) required considerable correction,*
404 *and changes were made to approximately 26% of coding decisions on individual*
405 *lexemes. Ref. (16 ) required corrections to only 0.5% of lexemes.*".

406    For the Sino-Tibetan lexical dataset, Zhang et al. compiled 90 lexical items from
407 the *Sino-Tibetan Etymological Dictionary and Thesaurus* (STEDT) project [5]. These
408 lexical items also appear in *Swadesh's 100-word list* [6]. These selected lexical items
409 facilitated the identification of 949 cognates across 109 Sino-Tibetan languages,
410 which were then encoded as 949 binary-coded lexical traits. Zhang et al. described
411 their cognate coding process as below: "*The lexical root-meanings used in this study*
412 *came from the Sino-Tibetan Etymological Dictionary and Thesaurus (STEDT)*
413 *project1, which was developed by a number of experienced historical linguists led by*
414 *James A. Matisoff over a 30-year period (URL: http://stedt.berkeley.edu/)......To*
415 *minimize the word lateral transfers, in this study we chose only the words with*
416 *meaning inside the Swadesh 100-word list, since they are relatively resistant to*
417 *borrowing2……In order to make sure that all the languages were comparable to each*
418 *other, we filtered only those languages with at least 90 lexical meanings of Swadesh*
419 *100-word list recorded (no matter whether an RM exists) and 30 – 120*
420 *RMs……Finally, we retained 109 ST language samples with 949 binary-coded lexical*
421 *RMs for further phylogenetic analyses.*"

422    For the Bantu lexical dataset, Grollemund et al. compiled 100 lexical items from
423 the *Atlas Linguistique du GABon list* [7], of which 68 lexical items overlap with
424 *Swadesh's 100-word list*. According to these lexical items, they recognized 3,859
425 cognates across 420 Bantu languages. These 3,859 cognates were further transformed
426 into 3,859 binary-coded lexical traits. Grollemund described their cognate coding
427 process as: "*For phylogenetic inference, we used a selection of 100 meanings*
428 *comprising a modified version of the Atlas Linguistique du GABon list (52). The Atlas*
429 *includes 159 meanings, and our sample of 100 meanings are those that are best*
430 *documented for the languages we studied……We identified 3,859 cognate sets across*

431    *the n = 100 meanings. These were coded as binary characters for purposes of*
432    *phylogenetic analysis.*"

433    For the Arawak lexical dataset, Walker et al. compiled *Swadesh's 100-word list*
434    and identified 694 cognates across 60 Arawak languages. Subsequently, these
435    cognates were then recoded as 694 binary-coded lexical traits. Walker et al. described
436    their cognate coding process as below: "*We compiled Swadesh [20] lists of 100*
437    *common vocabulary items and scored cognate sets across 60 Arawak languages and*
438    *dialects representing all the major branches of the Arawak language family (see*
439    *electronic supplementary material, table S1)……We transformed coded cognates into*
440    *binary codes for each variant with sites representing whether any particular cognate*
441    *set is present ('1') or absent ('0') in that language...... The method yields 694 sites of*
442    *which 88 per cent are complete.*"

443    According to the reviewer's suggestions, we have incorporated the
444    aforementioned contents about the cognate and binary-coded lexical trait in *Lines*
445    *373-382* of the revised main text.

446    **Tabel to Q12.** Example of cognate coding using two lexical items (Mouth and Bone)
447    for four languages: Apurina, Bare, Yavitero, and Palikur. Lexical lists (left table) are
448    transformed into binary codes for each cognate variant with sites representing whether
449    any particular cognate is present ("1") or absent ("0") in that language (right table).

|  | Lexical item | |
|---|---|---|
|  | **Mouth** | **Bone** |
| **Apurina** | *nama* | ***api*** |
| **Bare** | *numa* | *bani* |
| **Yavitero** | *numa* | *ihiu* |
| **Palikur** | *by* | *api* |

**Transform data into binary codings** →

|  | Mouth | | Bone | |
|---|---|---|---|---|
| **Lexical trait** | **A** | **B** | **A** | **B** |
| **Apurina** | 1 | 0 | 1 | 0 |
| **Bare** | 1 | 0 | 1 | 0 |
| **Yavitero** | 1 | 0 | 0 | 1 |
| **Palikur** | 0 | 1 | 1 | 0 |

450

451    **Reference**

[1] Bouckaert, Remco, et al. "Mapping the origins and expansion of the Indo-European language family." Science 337.6097 (2012): 957-960.

[2] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

[3] Grollemund, Rebecca, et al. "Bantu expansion shows that habitat alters the route and pace of human dispersals." Proceedings of the National Academy of Sciences 112.43 (2015): 13296-13301.

[4] Walker, Robert S., and Lincoln A. Ribeiro. "Bayesian phylogeography of the Arawak expansion in lowland South America." Proceedings of the Royal Society B: Biological Sciences 278.1718 (2011): 2562-2567.

[5] Matisoff, James A. "Sino-Tibetan etymological dictionary and thesaurus (STEDT)." Berkeley: Sino-Tibetan Etymological Dictionary and Thesaurus Project.(stedt. berkeley. edu/dissemination/STEDT. pdf)[accessed on18 October 2020] (2015).

[6] Swadesh, Morris. "Towards greater accuracy in lexicostatistic dating." International journal of American linguistics 21.2 (1955): 121-137.

[7] Hombert, Jean-Marie. "Atlas linguistique du Gabon." Revue gabonaise des Sciences de l'homme 2 (1990): 37-42.

*Q13: "Third, the changes in the state frequencies of linguistic traits are proportional to their sociolinguistic prestige in a certain area.". I don't get the logic of this sentence. What is the meaning of prestige here? The definition of prestige is expressed only in the next paragraph, it should be introduced before going into interpretations.*

**Replies to Q13:**

We thank the reviewer for pointing this out. This prestige parameter reflects the social opportunities or convenience for individuals who speak a specific language containing a particular trait state [1]. States of linguistic traits with higher prestige

480 would be more prevalent in future generations, while those with lower prestige would
481 be less prevalent. Accordingly, the prestige of a specific state in a linguistic trait can
482 be mathematically defined as the probability of this linguistic trait remaining in that
483 state after a unit of time. According to the reviewer's comment, we have modified the
484 corresponding section and rearranged the sequence of the paragraph related to the
485 prestige parameter as shown in *Lines 408-412* of the revised main text.

486 **Reference:**

487 [1] Abrams, Daniel M., and Steven H. Strogatz. "Modelling the dynamics of
488 language death." Nature 424.6951 (2003): 900-900.

489

490 *Q14: "It is noted that the larger length of the velocity vector of a language denotes*
491 *the more rapid change of this language during its evolution". The reader is*
492 *provided with no tools to understand this sentence. A schematic representation of a*
493 *vector could really help. E.g. what are the elements of a vector?*

494 **Replies to Q14:**

495 We appreciate your comment. In the revision, we have added a more
496 comprehensive schematic representation for the velocity vector in Figure 1d of the
497 revised main text. We also attach this subfigure related to the calculation of the
498 velocity vector at the end of this **Replies to Q14** (Figure to Q14).

499 As shown in Figure to Q14, we can see that each velocity vector contains two
500 aspects: direction and length. Each vector is calculated as the difference between the
501 past reconstructed and current trait states divided by the reconstruction time.
502 Accordingly, the direction of each vector signifies the direction of the diachronic
503 change of the linguistic traits in each language in the high-dimensional space and
504 low-dimensional PC space (i.e., 2-D PC plot). In short, the direction of each vector
505 depicts how the linguistic traits evolve into their current states. Moreover, when the
506 linguistic traits of a language undergo rapid evolution, its trait states should change
507 significantly over a given time period. Such change can be represented by the length
508 of the velocity vector visualized as an arrow in the high-dimensional space and
509 low-dimensional PC space. However, our study exclusively concentrates on the

510  language dispersal pattern which can be reflected solely by the directions of the
511  velocity vectors. Accordingly, the lengths of the velocity vectors are actually not
512  utilized in this study. Noting these, we have removed the descriptions about the
513  lengths of velocity vectors in the revised manuscript.



514

515  **Figure to Q14.** Schematic diagram of the calculation of velocity vector.

516

517  *Q15: Does the PCA find only two components, or the authors found that more*
518  *components did not lead to more variance explainability? Again, PCA here is*
519  *presented as a tool to find similarities among datapoints, actually it is a*
520  *rearrangement of the predictors of the model that tells what are the most important*
521  *features in the model. The authors say nothing about all this. Projecting the points*
522  *in to the PC space allows to visualize clusters, but actual clustering is performed by*
523  *other tools, such as k-nearest-neighbors.*

524  **Replies to Q15:**

525      Thank you for your comments. We have three specific reasons for selecting only
526  two principal components in this study which are explained below:

527     Firstly, visualizing the samples in a two-dimensional plane using two principal
528     components is a common and effective practice [1-3]. It enables a clear visualization
529     of the distribution of the data points in a two-dimensional plane. Accordingly, we also
530     selected two principal components and visualized the language samples in the
531     two-dimensional space.

532     Secondly, in the subsequent step of our approach, the language velocity field will
533     be projected from the PC space into the two-dimensional (i.e., longitude and latitude)
534     geographic space. By selecting the first two principal components, we ensure that the
535     PC space and the geographic space share an identical dimension, thereby preventing
536     the loss of information during geographic mapping. For instance, attempting to map a
537     three-dimensional language velocity field to a two-dimensional geographic space
538     would result in the loss of one crucial dimension of information regarding the
539     language velocity field. Nevertheless, the reviewer provided a novel insight into our
540     approach. It is that when the geographic coordinates of language samples have a
541     higher dimension, it would be prudent to retain more principal components for the
542     geographic mapping of the language velocity field.

543     Thirdly, according to the simulated validations, we found that relying on two
544     principal components was sufficient to estimate a reliable language velocity field in
545     the geographic map. Based on this language velocity field, we could accurately reflect
546     the language dispersal trajectories and centers. Consequently, we only selected two
547     principal components for the construction of the velocity field in this study.

548     Fourthly, we do not conduct the PCA algorithm to cluster or find similarities
549     among language samples. Actually, in this study, the PCA algorithm is only conducted
550     to recombine the original traits into two important traits. we plot each language
551     sample according its coordinate (PC1, PC2) in the 2-dimensional PC space. The
552     shorter Euclidean distances among language samples in PC space embody their higher
553     linguistic relatedness. However, if we aim to further identify which language samples
554     should be clustered together, we will need to employ other clustering approaches.
555     According to the reviewer's comments, we have revised the descriptions about the
556     PCA algorithm as shown in the *Lines 114-122* of the revised main text.

557     **Reference**

558     [1] Wang, Chuan-Chao, et al. "Genomic insights into the formation of human

559    populations in East Asia." Nature 591.7850 (2021): 413-419.

560    [2] Haak, Wolfgang, et al. "Massive migration from the steppe was a source for

561    Indo-European languages in Europe." Nature 522.7555 (2015): 207-211.

562    [3] Norvik, Miina, et al. "Uralic typology in the light of a new comprehensive

563    dataset." Journal of Uralic Linguistics 1.1 (2022): 4-42.

564

565    *Q16: it is not clear how the vectors are formed in the PC space. Up to my*

566    *understanding the PCA describes the datapoint with two components, hence I*

567    *expect to observe a single point with coordinates (PC1,PC2) in the PC space. By the*

568    *way, we cannot build a vector with one point. I understand from SI-3 that the*

569    *vectors are computed as the difference in the PC space of X(0) − X(-m), where t=0*

570    *represents now and t=-m represents a moment in the past. What is this moment in*

571    *the past? Then I read "Therefore, Vl describes the change of the state frequencies*

572    *of language l in a unit of time.". what is the unit of time? Years, centuries?*

573    **Replies to Q16:**

574    We sincerely thank the reviewer for bringing up these important points. We

575    address the reviewer's concern as below.

576    **1. The derivation of the velocity vectors in PC space**

577    We agree with the reviewer that PCA can describe each current language sample

578    with two components PC1 and PC2. The PC1 and PC2 are derived by applying a

579    matrix $A_{2 \times n}$ (2 rows and $n$ columns) to each current language sample $l_{current} = [x_1, \ldots,$

580    $x_n]^T$ ($n$ linguistic traits): $[PC1_{current}, PC2_{current}]^T = Al_{current} = A[x_1, \ldots, x_n]^T$. It can be

581    regarded as projecting a $n$-dimensional vector into a 2-dimensional PC space as a

582    2-dimensional vector. Therefore, we can only observe a single language point with a

583    coordinate ($PC1_{current}$, $PC2_{current}$) in the PC space.

584    However, given a dynamic model [1-3], our approach can reconstruct the past

585    trait states for each language sample according to its current observed trait states

586    noted as $l_{past} = [y_1, \ldots, y_n]^T$. When projecting current trait states for each language

587    sample into the PC space, we simultaneously project its past trait states into this PC

588      space as well: $[PC1_{past}, PC2_{past}]^T = Al_{past} = A[y_1, …, y_n]^T$. Therefore, we can observe

589      two points noted as $(PC1_{current}, PC2_{current})$ and $(PC1_{past}, PC2_{past})$ in the PC space,

590      where one represents the current trait states of this language sample, and another

591      represents its past trait states. By taking the difference between these two points

592      divided by the reconstruction time, we can derive a vector that describes the rate and

593      direction of the changes in the trait states of this language sample.

594      According to the reviewer's suggestions, the calculations of the vectors are

595      illustrated by a schematic diagram in Figure 1 in the revised main text. For the

596      convenience of the reviewer, we have attached the subfigure of Figure 1 related to the

597      calculation of velocity vectors below as Figure 1 to Q16.

598      **2.    The definition of a unit of time**

599      **(i) The definition of a unit of time is identical to the one in the phylogenetic study**.

600      The velocity vector is calculated as the difference in the PC space of $X(0) – X(-m)$

601      divided by reconstruction time $m$, where $t = 0$ represents the present time, and $t = -m$

602      represents a moment in the past. Here, $m$ denotes $m$ units of times, and $-m$ thus

603      represents $m$ units of times before the present time. Given that we often have limited

604      knowledge regarding the precise origin time of past languages, we thus define a unit

605      of time as one generation. It serves as a dimensionless time indicator representing the

606      period during which the linguistic traits in language accumulate one mutation. This

607      definition of the unit of time in our study is identical to the definition in the

608      phylogenetic tree where no exact time calibrations have been made (hereafter

609      non-time-calibrated phylogenetic tree).

610      To be specific, in a non-time-calibrated phylogenetic tree, the branch length

611      between a parent node and a child node (where the language is referred to as a node

612      for convenience hereafter) represents the time during which the child language has

613      evolved from its parent language. This branch length is typically represented by the

614      number of mutations that occurred in linguistic traits during the evolution of the child

615      language from its parent language. Because the longer evolutionary time of a

616      language results in more mutations being accumulated in linguistic traits (see Figure 2

617      to Q16 attached below) [4-5]. Under this circumstance, a unit of time is defined as the

618      period in which the linguistic traits of language undergo one mutation.

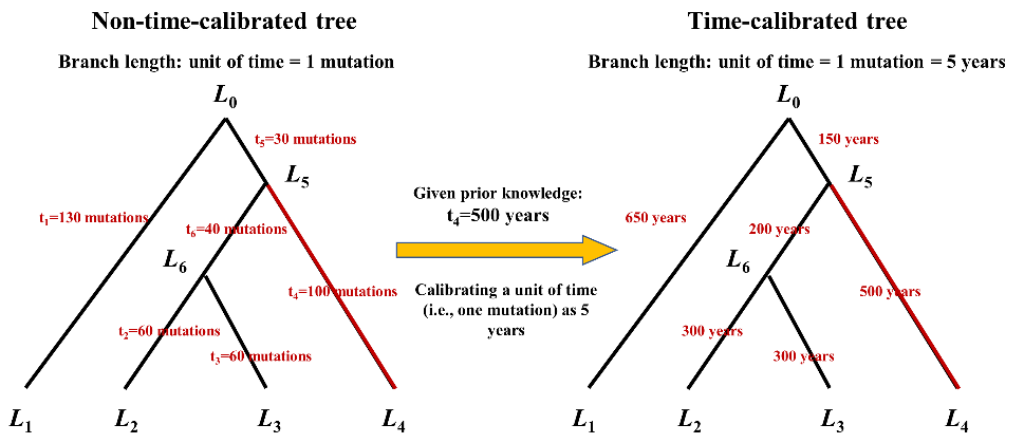619      **(ii) A unit of time can be calibrated based on prior origin time.** This dimensionless

620    unit of time can be further converted into the exact period once given the precise

621    origin time of the parent language (see Figure 2 to Q16 attached below). For instance,

622    we assume that one branch length between a parent language *L*5 and a child language

623    *L*4 within a non-time-calibrated phylogenetic tree corresponds to 100 mutations (see

624    Figure 2 to Q16 attached below). Moreover, we assume that we also possess prior

625    knowledge about the precise origin time of that parent language, said 500 years ago.

626    Accordingly, we can calibrate the unit of time as 500/100 = 5 years using the

627    commonly utilized strict molecular clock model in linguistics which assumes the

628    mutation rate is constant [3, 6]. According to this unit of time with exact time

629    calibration, we can calibrate all the branch lengths with exact periods in the

630    non-time-calibrated phylogenetic tree according to the times of mutations (see Figure

631    2 to Q16 attached below).

632        Similarly, the unit of time defined in our approach can also be converted to an

633    exact period in our approach, once we have prior knowledge about the precise origin

634    times of the past language samples. Nevertheless, the calibration of the unit of time in

635    our approach is not essential, since our approach is not designed to estimate the

636    divergence time of languages. It is just like the application of the phylogeographic

637    approach to a non-time-calibrated phylogenetic tree to solely infer the geographical

638    dispersal center of languages [7]. We have added the definition of unit of time into the

639    *Lines 441-443* of the revised main text.

640

**Figure 1 to Q16.** The calculation of velocity vectors in the PC space.



642

**Figure 2 to Q16.** Calibrating each branch length of the non-time-calibrated tree based on the mutation times and prior knowledge about language divergence times.

**Reference**

[1] Yang, Ziheng. "Maximum-likelihood estimation of phylogeny from DNA sequences when substitution rates differ over sites." Molecular biology and evolution 10.6 (1993): 1396-1401.

649    [2] Penny, David, et al. "Mathematical elegance with biochemical realism: the
650    covarion model of molecular evolution." Journal of Molecular Evolution 53 (2001):
651    711-723.

652    [3] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in
653    northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

654    [4] Choudhuri, Supratim. Bioinformatics for beginners: genes, genomes, molecular
655    evolution, databases and analytical tools. Elsevier, 2014.

656    [5] Lewis, Paul O. "A genetic algorithm for maximum-likelihood phylogeny
657    inference using nucleotide sequence data." Molecular biology and evolution 15.3
658    (1998): 277-283.

659    [6] Chang, Will, et al. "Ancestry-constrained phylogenetic analysis supports the
660    Indo-European steppe hypothesis." Language (2015): 194-244.

661    [7] Walker, Robert S., and Lincoln A. Ribeiro. "Bayesian phylogeography of the
662    Arawak expansion in lowland South America." Proceedings of the Royal Society B:
663    Biological Sciences 278.1718 (2011): 2562-2567.

664

665    *Q17: what is the delta score and how is it computed? It is never stated in the text,*
666    *nor in the SI*

667    **Replies to Q17:**

668    We thank the reviewer for pointing this out. The rationale for the delta score has
669    been introduced in the **Replies to Q11**. Here, we offer a brief description of its
670    calculation procedure.

671    For any quarter of four elements *x*, *y*, *u*, and *v*, we denote $d_{xy|uv} = d_{xy} - d_{uv}$.
672    Then, the delta score is defined as the ratio $\delta_q = \frac{d_{xv|yu} - d_{xu|yv}}{d_{xv|yu} - d_{xy|uv}}$ [1]. This ratio
673    measures the tree-likeness of the quartet *q* that $\delta_q = 0$ if $d_{xv|yu} = d_{xu|yv} = d_{xy|uv}$
674    hold. The larger the value of $\delta_q$ indicates the less treelike of *q*. The average value of
675    $\delta_q$ of the all-possible quarter of the language samples thus can serve as the metric to

676 quantify the overall tree-likeness of the language topology. In this study, the delta
677 score is calculated using the "*delta.plot*" function of the "*ape*" package [2]. The
678 corresponding contents have been included in *Lines 579-580* of the Materials and
679 Method section of the revised main text.

680 **Reference**

681 [1] Holland, Barbara R., et al. "δ plots: a tool for analyzing phylogenetic distance
682 data." Molecular biology and evolution 19.12 (2002): 2051-2059.

683 [2] Paradis, Emmanuel, and Klaus Schliep. "ape 5.0: an environment for modern
684 phylogenetics and evolutionary analyses in R." Bioinformatics 35.3 (2019): 526-528.

685

686 *Q18: Later on I read "In this study, we set m = 1.", but no reason is given, nor the*
687 *unit of time is stated. One year? One century? Again, this is very opaque.*

688 **Replies to Q18:**

689 We appreciate the reviewer for pointing these out. As mentioned in the **Replies to**
690 **Q16**, a unit of time in this study is defined as one generation, which serves as a
691 dimensionless time indicator representing a period during which the linguistic traits in
692 language accumulate one mutation. This dimensionless unit of time can be converted
693 into an exact time once the precise divergence time of the past language sample is
694 given. However, the exact time calibration of the unit of time is not necessary in our
695 approach, since our approach is designed to infer the dispersal pattern of languages
696 rather than their origin time.

697 In this study, the setting of $m = 1$ is chosen based on the results of both empirical
698 and simulated validations. To be specific, in simulated validations, we demonstrated
699 that relying on the setting $m = 1$ could estimate a reliable language velocity field in
700 the geographic space. Based on this language velocity field, the estimated language
701 dispersal center shows no significant difference from the prior given dispersal center
702 (Figure 1 to Q18).

703 Without a loss of generality, we also tested the robustness of the language
704 velocity field estimated through different settings of $m$ in simulated validations. The

705 results indicate that there are no significant differences among the language velocity
706 fields estimated through different settings of *m* (Figure 2 to Q18). These results
707 indicate that the rate of change of linguistic traits can remain relatively constant
708 during different evolutionary periods. It is compatible with the rate assumption of the
709 widely-used molecular clock model in linguistics that postulates the evolutionary rate
710 of linguistic traits is constant [1-2]. In other words, the velocity vector is almost
711 unchanged either setting *m* = 1 or setting *m* as other different reconstruction times.
712 Therefore, it is feasible to estimate the velocity vector for representing the diachronic
713 change in linguistic traits by setting *m* = 1.

714    According to the simulated validations, we further set *m* = 1 in the empirical
715 applications. Without a loss of generality, we also tried different parametric settings of
716 *m* in the empirical applications. The results also suggested that the language velocity
717 field was robust under different settings of *m* (Figure 3 to Q18), and all could identify
718 the language dispersal centers that can be supported by genetic and archaeological
719 evidence. Based on all these empirical and simulated validations, we ultimately set *m*
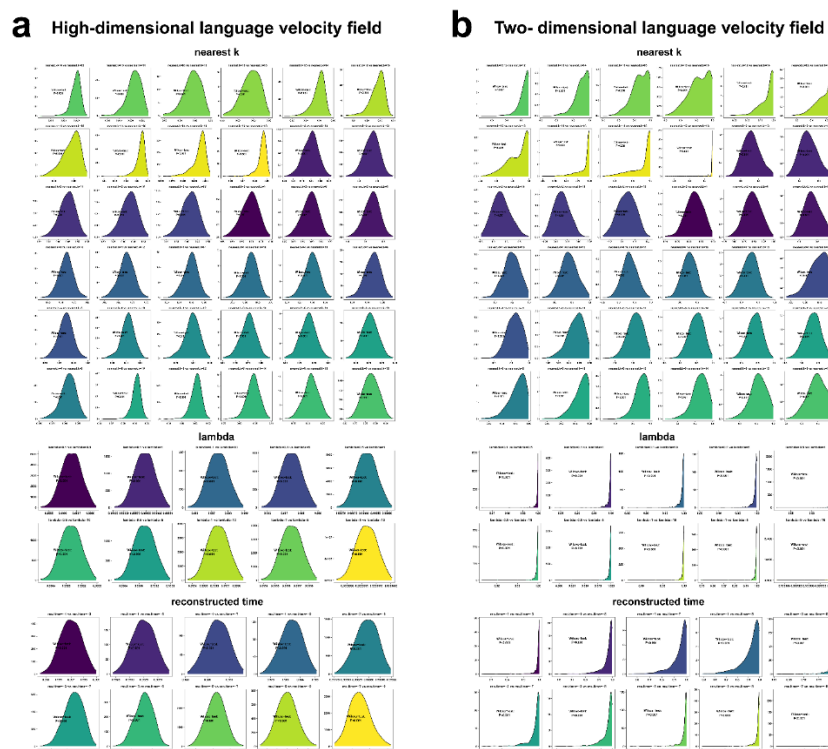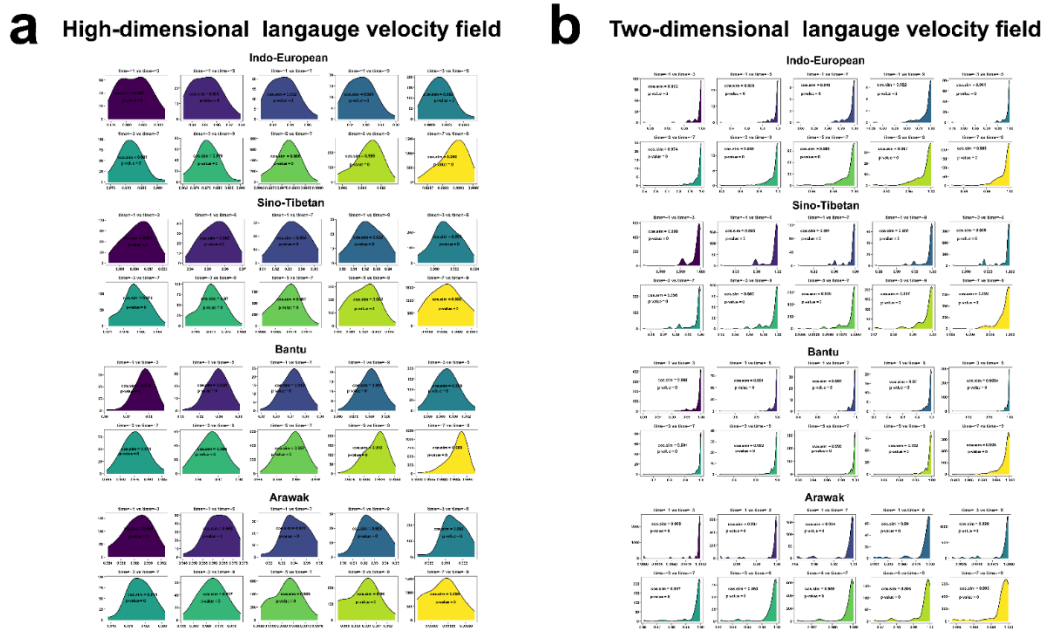720 = 1 as the default parametric value in our approach.

721


722 **Figure 1 to Q18. The simulated validation for the effectiveness of the language**
723 **velocity field estimation (LVF) under different parametric settings.** The
724 probability density plot demonstrates the distributions of the errors of the longitude
725 and latitude respectively between the true and inferred language dispersal center
726 estimated from 1,000 simulated datasets under different parametric settings. These

28

727      parameters are the number of the grid points *n.grid* (*n.grid* = 50, 100, 200, 300, 400,

728      and 500); the number of the nearest neighbors *k* (*k* = 2, 4, 6, …, and 18); mutation rate

729      of Poisson process $\lambda$ ($\lambda$ = 0.1, 0.5, 1, 5, and 10); reconstruction time *m* (*m* = 1, 3, 5, 7,

730      and 9). We set the default parametric values as *n.grid* = 300, *k* = 4, $\lambda$ = 1, and *m* = 1

731      when varying across the settings of these parameters respectively. The black texts are

732      the *p*-value of the statistical significance of the error derived from the Wilcoxon

733      rank-sum test. *p*-value > 0.05 denotes the statistical non-significance of the error

734      (significantly equal to 0).

735



736

737      **Figure 2 to Q18. The simulated validation for the robustness of the language**

738      **velocity field estimation (LVF) under different parametric settings.** The

739      probability density plot demonstrates the distribution of the average cosine similarity

740      between language velocity fields estimated from 1,000 simulated datasets under

741      different parametric settings. The parameters are the number of the nearest neighbors

742      *k* (*k* = 2, 4, 6, …, and 18); mutation rate of Poisson process $\lambda$ ($\lambda$ = 0.1, 0.5, 1, 5, and

743      10); reconstruction time *m* (*m* = 1, 3, 5, 7, and 9). We set the default parametric values

744      as *k* = 4, $\lambda$ = 1, and *m* = 1 when varying across the settings of these parameters

745      respectively. The black texts are the *p*-value of the statistical significance of this

748



749

750 **Figure 3 to Q18. The empirical validation for the robustness of the language**
751 **velocity field estimation (LVF) against the setting of the reconstruction time.** The
752 probability density plot demonstrates the distribution of the cosine similarity among
753 the language velocity vectors calculated under different settings of reconstruction time
754 *m* (*m* = 1, 3, 5, 7, and 9) before the current time in four language families and groups.
755 We set the default parametric values as $k = 10$ and $\lambda = 1$ when varying across the
756 settings of *m*. The black texts are the average similarity of the distribution of
757 similarity and the *p*-value of the statistical significance of this average similarity
758 derived from the permutation test (Permutation Times = 500). The average similarity
759 ranges from 0 to 1, where 1 denotes that these two velocity fields are most similar and
760 0 is dissimilar. *p*-value < 0.05 denotes the statistical significance of the average
761 similarity.

762 **Reference**

763 [1] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in
764 northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

765     [2] Chang, Will, et al. "Ancestry-constrained phylogenetic analysis supports the
766     Indo-European steppe hypothesis." Language (2015): 194-244.

767

768     *Q19. I do not understand the physical meaning of this vectorial framework because*
769     *no clear explanation is provided.*
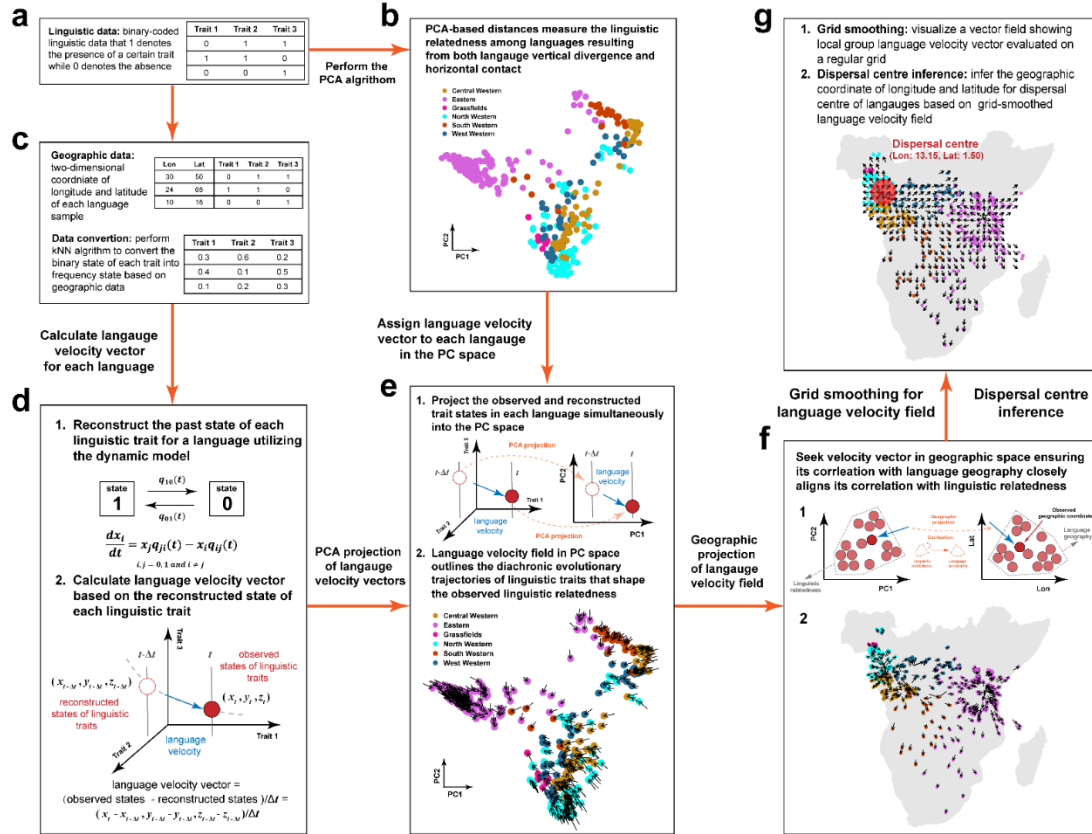
770     **Replies to Q19:**

771        To provide a clearer explanation of our approach, we have added more detailed
772     explanations of our approach in the revised main text (*Lines 109-151*). Moreover, we
773     have also added detailed mathematical formulas of our approach in the Materials and
774     Methods section. As the supplementary, we have also redrawn the schematic diagram
775     presented as Figure 1 to Q19 (also referred to as Figure 1 in the revised main text) to
776     visually elucidate the rationale and calculation procedure of our approach.

777        **Our approach shares the same theoretical foundation as the phylogeographic**
778     **approach but with different implementation strategies.** As the most prevailing
779     approach, the phylogeographic approach performs two major steps to infer language
780     dispersal patterns. The first is to obtain a phylogenetic tree to delineate the
781     evolutionary trajectories of linguistic traits that shape the observed linguistic
782     relatedness (Figure 2 to Q19) [1-3]. The second is to project the phylogenetic tree into
783     the geographic space based on the correlation between linguistic relatedness and
784     language geography (Figure 2 to Q19) [1-4]. With the projection, evolutionary
785     trajectories of linguistic traits can be transformed into language dispersal trajectories.
786     Our approach shares the similar two major steps as the phylogeographic approach that
787     infers language dispersal through the diachronic evolution of linguistic traits (Figure 2
788     to Q19). However, our approach employs different strategies to carry out these two
789     steps compared to the phylogeographic approach.

790        **The velocity field in PC space delineates the diachronic evolutionary**
791     **trajectories of linguistic traits that shape the observed linguistic relatedness.** Our
792     approach conducts the PCA-based distance rather than a phylogenetic tree to
793     represent linguistic relatedness. Specifically, the PCA algorithm is conducted to
794     rearrange the lexical traits into two principal components namely PC1 and PC2.
795     According to PC1 and PC2, the distribution of language samples can be visualized in
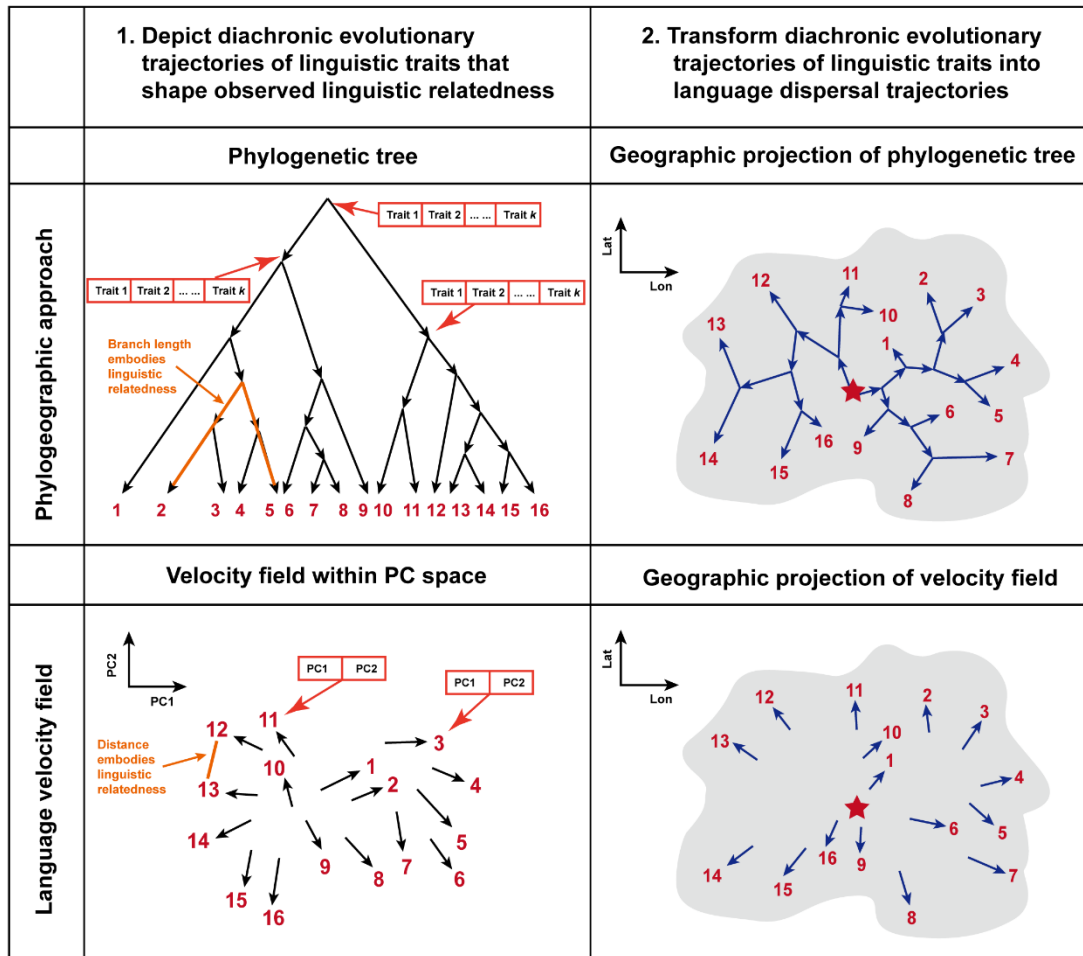
796  the PC space. The shorter distances among language samples in the PC space imply
797  their higher linguistic relatedness. In parallel, the language velocity vector is
798  estimated to demonstrate the direction of the average change of trait states for each
799  language sample in a unit of time. With the past trait states reconstructed by the
800  dynamic model, the velocity vector can be calculated by dividing the diachronic
801  changes in trait states of each language sample by the *m* unit of time. This velocity
802  vector depicts how the linguistic traits in a language sample evolve into their current
803  states. By mapping these velocity vectors into the PC space, a language velocity field
804  can be derived on the PC space to delineate the diachronic evolutionary trajectories of
805  linguistic traits that shape the observed linguistic relatedness. This velocity field in PC
806  space functions similarly to the phylogenetic tree in the phylogeographic approach.

807  **Projecting velocity field into geographic space to transform the evolutionary**
808  **trajectories of linguistic traits into the language dispersal trajectories**. Based on
809  the correlation between observed linguistic relatedness and language geography, we
810  further project each velocity vector from PC space into geographic space utilizing
811  kernel projection [3]. The rationale of this projection is to search for the velocity
812  vector in the geographic space ensuring that its correlation with language geography
813  closely matches with its correlation with linguistic relatedness. With the kernel
814  projection, the vector directions in the geographic space, which compose a set of
815  trajectories, render from where the observed language samples diffuse into their
816  current locations. This geographic projection of the velocity field is similar to the
817  projection of the phylogenetic tree into the geographic space to outline the dispersal
818  trajectories in the phylogeographic approach.

819

**Figure 1 to Q19. Schematic overview of the language velocity field estimation (LVF) for inferring the dispersal trajectories and centers of languages.** The computational procedures of the LVF comprise two major steps. Subfigures (a) to (e) illustrate the first step which is to estimate a velocity field on the PC space to outline the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. Subfigures (f) to (g) illustrate the second step, which is to project the velocity field from PC space into geographic space. Within the velocity field in geographic space, the directions of the velocity vectors compose a set of continuously changing trajectories that delineate from where these languages diffuse to their current locations. These procedures are exemplified using the Bantu language family. Comprehensive insights into the underlying principles and computational steps can be found in the Materials and Methods section, as well as Supplementary Note 1.

833

**Figure 2 to Q19. Language velocity field estimation (LVF) shares the same foundation as the phylogeographic approach but with different implementation strategies.** Both LVF and phylogeographic approach entails two major steps to infer language dispersal pattern. The first is to depict the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. The second is to transform these diachronic evolutionary trajectories of linguistic traits into language dispersal trajectories. In the phylogenetic tree, each language is determined by $k$ linguistic traits. In the velocity field within PC space, each language is determined by PC1 and PC2 which are rearranged from the $k$ linguistic traits through the PCA algorithm. The red number denotes a language. The black arrow signifies the evolutionary direction of linguistic traits in a language. The blue arrow represents the dispersal direction of a language. The red star denotes the estimated dispersal center.

**Reference**

34

848 [1] Bouckaert, Remco, et al. "Mapping the origins and expansion of the
849 Indo-European language family." Science 337.6097 (2012): 957-960.

850 [2] Grollemund, Rebecca, et al. "Bantu expansion shows that habitat alters the route
851 and pace of human dispersals." Proceedings of the National Academy of Sciences
852 112.43 (2015): 13296-13301.

853 [3] Currie, Thomas E., et al. "Cultural phylogeography of the Bantu Languages of
854 sub-Saharan Africa." Proceedings of the Royal Society B: Biological Sciences
855 280.1762 (2013): 20130695.

856 [4] Koile, Ezequiel, et al. "Geography and language divergence: The case of Andic
857 languages." Plos one 17.5 (2022): e0265460.

858 [5] La Manno, Gioele, et al. "RNA velocity of single cells." Nature 560.7719 (2018):
859 494-498.

860

861 *Q20: the authors said that they study the spatial dispersal of languages along*
862 *10,000 years, to my understanding the vector field describes the change of the*
863 *language between one exact moment of the past and t=0, which is supposed to be*
864 *today.*

865 **Replies to Q20:**

866     We are grateful for the reviewer's comments. The reason why we mentioned
867 10,000 years in the main text is that all four language families and groups utilized in
868 this study originated within the last 10,000 years. For the Indo-European languages,
869 different phylogenetic studies have reported that their origin time could be either
870 approximately 8,000 to 9,500 years ago [1] or approximately 6,000 years ago [2]. For
871 the Sino-Tibetan languages, its initial divergence has been estimated to occur between
872 4,000 to 8,000 years ago [3-5]. The origin of the Bantu languages has been traced
873 back to roughly 5,000 years ago [6]. Although the detailed origin time of the Arawak
874 languages remains unclear, its origin is interlinked with the agricultural advancement
875 in lowland South America around 5,000 years ago [7-8]. Consequently, the origin of
876 Arawak languages should have dated at most 5,000 years ago. Overall, 10,000 years

877     is the upper limit of the origin time for these four language families and groups.

878     **Reference**

879     [1] Bouckaert, Remco, et al. "Mapping the origins and expansion of the
880     Indo-European language family." Science 337.6097 (2012): 957-960.

881     [2] Chang, Will, et al. "Ancestry-constrained phylogenetic analysis supports the
882     Indo-European steppe hypothesis." Language (2015): 194-244.

883     [3] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in
884     northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

885     [4] Zhang, Hanzhi, et al. "Dated phylogeny suggests early Neolithic origin of
886     Sino-Tibetan languages." Scientific Reports 10.1 (2020): 20792.

887     [5] Sagart, Laurent, et al. "Dated language phylogenies shed light on the ancestry of
888     Sino-Tibetan." Proceedings of the National Academy of Sciences 116.21 (2019):
889     10317-10322.

890     [6] Grollemund, Rebecca, et al. "Bantu expansion shows that habitat alters the route
891     and pace of human dispersals." Proceedings of the National Academy of Sciences
892     112.43 (2015): 13296-13301.

893     [7] Diamond, Jared, and Peter Bellwood. "Farmers and their languages: the first
894     expansions." Science 300.5619 (2003): 597-603.

895     [8] Clement, Charles Roland, et al. "Crop domestication in the upper Madeira River
896     basin." Boletim do Museu Paraense Emílio Goeldi. Ciências Humanas 11 (2016):
897     193-205.

898 **Replies to Reviewer 2:**

899 *Q1: As I stated in my previous reviews of this paper, it is interesting, convincing,*
900 *and historically significant in its conclusions. I am pleased to see that the authors*
901 *have cut down the paper to deal with the four clearest examples, these being*
902 *Indo-European, Sino-Tibetan, Bantu, and Arawak. The more troublesome*
903 *Austroasiatic, Japonic and Oceanic examples have been removed, and I think this*
904 *decision has added greatly to the clarity of the paper. It deserves to be published in*
905 *Nature Communications. My first comment is that the paper still needs a light level*
906 *of English editing. I do not have time to do this on behalf of the authors, but*
907 *perhaps I can use the abstract as an example of how some light editing might*
908 *increase its clarity:*

909 **Replies to Q1:**

910 We are deeply grateful for the reviewer's great support and affirmation of our
911 work. Moreover, we also would like to express our sincere appreciation for the
912 reviewer personally revising our abstract. According to this valuable example of
913 revision, we have carefully revised our manuscript. This revision involves correcting
914 many typos and grammatical errors, and rephrasing some lengthy and vague sentences.
915 Moreover, we also engaged the AJE language editing service to thoroughly edit the
916 language of our manuscript (ID: Q2K9ZRSF). We expect that our revisions could
917 enhance the readability and clarity of our manuscript for native English speakers.

918

*Q2: Figure 2 shows the proposed agricultural homeland in northern Amazonia for Arawak. This conflicts with text lines 184-186, where it is stated that " In addition, the language velocity field posited the dispersal of Arawak languages originated from the border of Peru, Brazil, and Bolivia in Western Amazonia, which was geographically close to the known ancient agricultural homeland of South America in the Andes". This statement implies a homeland much further to the south than shown on the map, which is what the archaeology would suggest. The map shows an area too far north. I note in Supplementary Notes 1 Table S2 that the Arawak homeland is put in the northern lowlands of Bolivia (upper Madeira River), which is precisely where I would expect it to be!*

**Replies to Q2:**

We are sincerely grateful for the reviewer to point these out. According to the reviewer's suggestions, we found inaccuracies in our descriptions regarding the origin of Arawak languages near the Andes, since their estimated dispersal center was indeed located too far from the Andes foothills. As mentioned by the reviewer, the dispersal center of Arawak languages estimated by our approach is located in the upper Madeira River basin within the northern lowlands of Bolivia. Accordingly, we further made some literature investigations about the upper Madeira River basin.

To our knowledge, the Madeira River rises from the Andes and flows through a larger part of the Southwestern Amazonian [1]. The upper Madeira River basin, which has raised numerous complex Neolithic Societies, has long been regarded as an important homeland of ancient agriculture in lowland South America [2]. In this area, plenty of crops have been domesticated, such as manioc, peanuts, peach palms, coca, and tobacco. It is noted that the estimated dispersal center of the Arawak language is located in the upper Madeira River basin. This estimation implies that the Arawak language origin is associated with the agricultural origin in Southwestern Amazonian. Accordingly, we revise the sentences of *Lines 184-186* in the original main text into: "***In addition, the LVF showed the dispersal of Arawak languages originating from the northern lowlands of Bolivia in the upper Madeira River basin, which is an important homeland of ancient agriculture in lowland South America.***" as shown in *Lines 202-204* of the revised main text.

**Reference**

951  [1] Clement, Charles Roland, et al. "Crop domestication in the upper Madeira River
952  basin." Boletim do Museu Paraense Emílio Goeldi. Ciências Humanas 11 (2016):
953  193-205.

954  [2] Piperno, Dolores R. "The origins of plant cultivation and domestication in the
955  New World tropics: patterns, process, and new developments." Current anthropology
956  52.S4 (2011): S453-S470.

957

958  *Q3: Likewise, lines 187-189 state " Moreover, in the case of Sino-Tibetan*
959  *languages, their dispersal center was inferred in the Gansu province of China*
960  *(Figure 2b). It was approximate to the geographic ranges of the Yangshao*
961  *(7,000-5,000 years BP) and/or Majiayao (5,500-4,000 years BP) Neolithic cultures,*
962  *although it was far from the ancient agricultural homelands known in the Yangzi*
963  *and Yellow River Basins of China." Surely, Yangshao and Majiayao were centrally*
964  *located in the Yellow River homeland of millet and pig agriculture? I cannot*
965  *understand what is meant here, although, of course, the Yangzi is a different matter.*

966  **Replies to Q3:**

967  We greatly appreciate the reviewer for bringing these points out. The original
968  intention of our statement was to express that the dispersal and origin of Sino-Tibetan
969  languages appear to have stronger connections with the agriculture that originated in
970  the Yellow River basin rather than the Yangzi River basin.

971  Early farming in China can be divided into two distinct attributes. One originated
972  in the Yellow River basin with a focus on millet cultivation, while another one was
973  developed in the Yangzi River basin with a focus on rice cultivation [1].
974  Geographically located in the center of the Yellow River basin, Yangshao, and
975  Majiayao Neolithic cultures were predominantly engaged in millet cultivation, as
976  evidenced by the archaeological materials [2-3]. Therefore, the estimated
977  Sino-Tibetan language dispersal center located in the geographic ranges of Yangshao
978  and Majiayao Neolithic cultures indicates that the Sino-Tibetan languages could have
979  dispersed with the spread of millet from the Yellow River basin rather than the Yangzi
980  River basin.

However, according to the reviewer's suggestion, we think that it is not necessary to mention the agriculture in the Yangzi River basin in this study which is not relevant to the case of Sino-Tibetan languages. The agriculture in Yangzi River should be another story in another research. Accordingly, we have revised the sentences in *Lines 187-189* of the original main text as: "***Moreover, in the case of Sino-Tibetan languages, their dispersal centre was inferred to be located in the Gansu Province of China (Figure 2b). This centre is situated within the geographic ranges of the Yangshao (7,000-5,000 years BP) and/or Majiayao (5,500-4,000 years BP) Neolithic cultures 6 in the ancient agricultural homeland of China, the Yellow River plains.***" in the *Lines 195-198* of the revised main text.

**Reference**

[1] Deng, Zhenhua, et al. "From early domesticated rice of the middle Yangtze Basin to millet, rice and wheat agriculture: Archaeobotanical macro-remains from Baligang, Nanyang Basin, Central China (6700–500 BC)." PLoS One 10.10 (2015): e0139885.

[2] Sagart, Laurent, et al. "Dated language phylogenies shed light on the ancestry of Sino-Tibetan." Proceedings of the National Academy of Sciences 116.21 (2019): 10317-10322.

[3] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

1001  *Q4: The discussion from lines 197 to 298 is highly technical, and I have no*
1002  *observations on it. Much the same applies to the materials and methods section. I*
1003  *can understand from lines 301-9 that the basic data come from a geographical*
1004  *plotting of cognate presences and absences, but I was puzzled by the statement*
1005  *(lines 304-6) "Lexical cognates of these language samples in each language family*
1006  *or group were binary-coded traits..." This sentence seems to confuse the concepts of*
1007  *cognate and language. How many cognate terms were used in the analysis, and*
1008  *from which proto-language levels were these cognates derived? In other words, how*
1009  *was a cognate defined? This might be explained in the supplementary data, but I*
1010  *think it should be clearer here in the main text.*

1011  **Replies to Q4:**

1012  We appreciate these valuable comments. In this study, we have used four lexical
1013  datasets encompassing 103 Indo-European, 109 Sino-Tibetan, 420 Bantu, and 60
1014  Arawak languages, respectively, which were derived from previously published works
1015  [1-4]. These lexical datasets are constructed upon the foundation of cognates (also
1016  referred to as cognate sets) which are varied word expressions for a particular lexical
1017  item (meaning) across diverse languages. These linguistic expressions (cognates) for
1018  the same lexical item have been identified as being inherited from a common ancestor.
1019  Within each lexical dataset, every linguistic expression (cognate) has been recorded as
1020  a binary lexical trait, where a value of 1 indicates its presence in a language, while 0
1021  indicates its absence.

1022  To be specific, for the Indo-European lexical dataset, Bouckaert et al. compiled
1023  207 lexical items [1] which facilitated the identification of 5,995 lexical cognates
1024  across 103 Indo-European languages. These cognates were further recoded into 5,995
1025  binary-coded lexical traits. Bouckaert et al. described their cognate coding process as
1026  follows: "*We recorded word forms and cognacy judgments across 207 meanings in*
1027  *103 contemporary and ancient languages…. Cognate data were coded as binary*
1028  *characters showing the presence or absence of a cognate set in a language. There*
1029  *were 5995 cognate sets in total, with most meanings represented by several different*
1030  *cognate sets. All cognate coding decisions were checked with published historical*
1031  *linguistic sources (Table S1). The database contained 25908 cognate-coded lexemes.*
1032  *Of these, 67% came originally from ref. (17 ), 14% from ref. (16 ), and 19% were*
1033  *newly compiled from published sources. Ref. (17 ) required considerable correction,*

1034 *and changes were made to approximately 26% of coding decisions on individual*
1035 *lexemes. Ref. (16 ) required corrections to only 0.5% of lexemes.*".

1036     For the Sino-Tibetan lexical dataset, Zhang et al. compiled 90 lexical items from
1037 the *Sino-Tibetan Etymological Dictionary and Thesaurus* (STEDT) project [5]. These
1038 lexical items can be also found in *Swadesh's 100-word list* [6]. These chosen lexical
1039 items led to the detection of 949 cognates across 109 Sino-Tibetan languages, which
1040 were then encoded as 949 binary-coded lexical traits. Zhang et al. described their
1041 cognate coding process as below: "*The lexical root-meanings used in this study came*
1042 *from the Sino-Tibetan Etymological Dictionary and Thesaurus (STEDT) project1,*
1043 *which was developed by a number of experienced historical linguists led by James A.*
1044 *Matisoff over a 30-year period (URL: http://stedt.berkeley.edu/)......To minimize the*
1045 *word lateral transfers, in this study we chose only the words with meaning inside the*
1046 *Swadesh 100-word list since they are relatively resistant to borrowing2……In order*
1047 *to make sure that all the languages were comparable to each other, we filtered only*
1048 *those languages with at least 90 lexical meanings of the Swadesh 100-word list*
1049 *recorded (no matter whether an RM exists) and 30 – 120 RMs……Finally, we*
1050 *retained 109 ST language samples with 949 binary-coded lexical RMs for further*
1051 *phylogenetic analyses.*"

1052     For the Bantu lexical dataset, Grollemund et al. selected 100 lexical items from
1053 the *Atlas Linguistique du GABon list* [7], of which 68 lexical items overlap with
1054 *Swadesh's 100-word list*. According to these lexical items, they recognized 3,859
1055 cognates across 420 Bantu languages. These cognates were further transformed into
1056 3,859 binary-coded lexical traits. Grollemund described their cognate coding process
1057 as: "*For phylogenetic inference, we used a selection of 100 meanings comprising a*
1058 *modified version of the Atlas Linguistique du GABon list (52). The Atlas includes 159*
1059 *meanings, and our sample of 100 meanings are those that are best documented for the*
1060 *languages we studied……We identified 3,859 cognate sets across the n = 100*
1061 *meanings. These were coded as binary characters for purposes of phylogenetic*
1062 *analysis.*"

1063     For the Arawak lexical dataset, Walker et al. compiled *Swadesh's 100-word list*
1064 and identified 694 cognates across 60 Arawak languages. Subsequently, these
1065 cognates were then recoded as 694 binary-coded lexical traits. Walker et al. described
1066 their cognate coding process as below: "*We compiled Swadesh [20] lists of 100*

1067 *common vocabulary items and scored cognate sets across 60 Arawak languages and*
1068 *dialects representing all the major branches of the Arawak language family (see*
1069 *electronic supplementary material, table S1)……We transformed coded cognates into*
1070 *binary codes for each variant with sites representing whether any particular cognate*
1071 *set is present ('1') or absent ('0') in that language...... The method yields 694 sites of*
1072 *which 88 per cent are complete.*"

1073     According to the reviewer's suggestions, we have revised the corresponding
1074 contents as shown in the *Lines 373-382* of the revised main text.

1075 **Reference**

1076 [1] Bouckaert, Remco, et al. "Mapping the origins and expansion of the
1077 Indo-European language family." Science 337.6097 (2012): 957-960.

1078 [2] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in
1079 northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

1080 [3] Grollemund, Rebecca, et al. "Bantu expansion shows that habitat alters the route
1081 and pace of human dispersals." Proceedings of the National Academy of Sciences
1082 112.43 (2015): 13296-13301.

1083 [4] Walker, Robert S., and Lincoln A. Ribeiro. "Bayesian phylogeography of the
1084 Arawak expansion in lowland South America." Proceedings of the Royal Society B:
1085 Biological Sciences 278.1718 (2011): 2562-2567.

1086 [5] Matisoff, James A. "Sino-Tibetan etymological dictionary and thesaurus
1087 (STEDT)." Berkeley: Sino-Tibetan Etymological Dictionary and Thesaurus
1088 Project.(stedt. berkeley. edu/dissemination/STEDT. pdf)[accessed on18 October 2020]
1089 (2015).

1090 [6] Swadesh, Morris. "Towards greater accuracy in lexicostatistic dating."
1091 International journal of American linguistics 21.2 (1955): 121-137.

1092 [7] Hombert, Jean-Marie. "Atlas linguistique du Gabon." Revue gabonaise des
1093 Sciences de l'homme 2 (1990): 37-42.

1094

*Q5: Lines 449-40 state: "The diversity approach is an alternative phylogenetic tree-free approach and simply infers the location of the language homeland to the areas with the highest linguistic diversity." What is meant here by linguistic diversity? Does it relate to relative times of splitting from an inferred phylogenetic family tree? (i.e., deeper-splitting subgroups are older)? I presume it is not simply related to number of languages.*

**Replies to Q5:**

We are sorry for the lack of clarity regarding the definition of linguistic diversity. As the reviewer correctly mentioned, linguistic diversity is not determined solely by the number of languages. As described by Wichmann and Sapir [1-3], the level of linguistic diversity is determined by the degree of differentiation among languages within a specific geographical area. Higher linguistic diversity indicates greater dissimilarities among the languages within that region. Consequently, even if there is a large number of languages in a particular geographic area, the linguistic diversity might still be low if those languages do not exhibit significant distinctions with each other.

The traditional diversity approach does not directly involve the divergence time provided by the phylogenetic tree for calculation. It simply measures the degree of distinctions among the observed languages (i.e., linguistic diversity) and assumes that the homeland of languages should be located in the area possessing the largest linguistic diversities [3]. Nevertheless, the theoretical foundation of this approach is somewhat related to the divergence time as the reviewer mentioned. In short, the diversity approach assumes that early divergence exhibits a higher divergence rate, which subsequently leads to the birth of an extraordinary number of distinct languages around the language homeland [3]. However, this theoretical underpinning has always been criticized because no solid evidence has been proposed to link divergence rate and homeland location. Additionally, other population activities, such as the migration of native speakers out of their original homeland, could also alter the linguistic diversity of the language homeland [4].

Following the reviewer's suggestions, we have added more detailed descriptions of the linguistic diversity approach in *Lines 308-311* of the revised main text. Moreover, a more comprehensive discussion of the linguistic diversity approach can

1127     be found in Supplementary Note 2: Section 2.

1128     **Reference**

1129     [1] Sapir, Edward. Time perspective in aboriginal American culture: a study in
1130     method. No. 13. Government Printing Bureau, 1916.

1131     [2] Wichmann, Søren, and Taraka Rama. "Testing methods of linguistic homeland
1132     detection using synthetic data." Philosophical Transactions of the Royal Society B
1133     376.1824 (2021): 20200202.

1134     [3] Wichmann, Søren, André Müller, and Viveka Velupillai. "Homelands of the
1135     world's language families: A quantitative approach." Diachronica 27.2 (2010):
1136     247-276.

1137     [4] Neureiter, Nico, et al. "Can Bayesian phylogeography reconstruct migrations and
1138     expansions in linguistic evolution?." Royal Society open science 8.1 (2021): 201079.

1139

1140     *Q6: I noticed in Supplementary Note 1 that phylogenetic discussions of*
1141     *Austroasiatic, Japonic and Oceanic are still mentioned, even through these*
1142     *groupings are no longer discussed in the main text.*

1143     **Replies to Q6:**

1144     We express our appreciation to the reviewer for bringing these points to our
1145     attention. In the revision, we have deleted the discussions related to the Austroasiatic,
1146     Japonic, and Oceanic languages in Supplementary Note 1.

1147

1148     *Q7: Supplementary Notes 2: it is not clear to me that Supplementary sections 2 and*
1149     *3 are really necessary (The interdisciplinary alignment of Genetics, Archaeology,*
1150     *and Linguistics; The Age-Area Hypothesis for inferring the language homeland). I*
1151     *think the observations made in this paper can stand quite well without them.*

1152     **Replies to Q7:**

1153    We sincerely appreciate the reviewer's suggestions. Following these suggestions,
1154    we have made several revisions to the main text and Supplementary Note 2. To be
1155    specific, we have excluded Section 3 (i.e., The Interdisciplinary Alignment of
1156    Genetics, Archaeology, and Linguistics) from Supplementary Note 2. After careful
1157    consideration, we have decided to retain Section 2 within Supplementary Note 2. This
1158    decision is motivated by the fact that the diversity approach is another famous
1159    phylogeny-free approach for identifying the language dispersal center. In our study,
1160    we have undertaken empirical comparisons between our approach and this
1161    methodology. As a result, Section 2 of Supplementary Note 2 offers an invaluable
1162    complement to the main text, providing readers with a more comprehensive grasp of
1163    the underlying rationale and limitations of the diversity approach.

1164

1165

**Replies to Reviewer 3:**

*Q1: I find this study generally quite interesting, since the authors claim that they have developed a new method that allows to represent historical dynamics of individual languages in comparison with neighboring languages by multidimensional vectors, which can then be projected in lower-dimensional space in order to even infer the original locations from which the language family as a whole dispersed. While interesting, I see some general problems with the study, mainly its fit with the journal where it was submitted to, and as a result, I recommend it to be rejected -- not because it is too low in quality, but rather because it is not a good fit with the journal, as I'll explain below. Apart from this, I see some major and minor flaws, which I'll discuss below. First, regarding the fit of the approach: What the authors propose is a methodological study, a new methodology of which they claim it outperforms established -- albeit controversial -- methods. In such a case, the journal where they submitted their study to, does not really qualify as a good fit, since we do not deal with new findings (they cannot be made until the method has been thoroughly evaluated) but rather with a new method that needs to be shown to work. For this reason, I think some journal like "Nature Methods" would be a much better fit here.*

**Replies to Q1:**

We are genuinely grateful for the reviewer's recommendation regarding the potential fit of our manuscript with *Nature Methods*, which is another outstanding Nature-branded journal renowned for its specialization in novel methods. Nonetheless, we firmly maintain our conviction that our work is ideally suited for *Nature Communications*.

Firstly, *Nature Communications* stands as a top-rank multidisciplinary journal that is devoted to publishing high-quality research in all interdisciplinary areas. Apart from reporting novel discoveries, it also has published many papers that propose novel methods to address interesting scientific questions. Specifically, diverse velocity field-based methods applicable to various research fields have been published in *Nature Communications*. These velocity fields have contributed to inferring the trajectories of dynamic changes in natural and social systems such as single-cell differentiation [1-2], human mobility [3], and atmospheric circulation [4].

1198 Accordingly, we think that our paper, which proposes a novel velocity field-based
1199 method to infer the language dispersal trajectory, is also suitable to the aim and scope
1200 of *Nature Communications*.

1201 Secondly, although our paper presents a new computational approach, its essence
1202 remains firmly rooted in multidisciplinary exploration. Our study seeks to investigate
1203 the spatial alignment of linguistic, genetic, and archaeological evidence in
1204 reconstructing prehistoric population activities worldwide. We believe that this topic
1205 could spark broad interest among researchers devoted to the interdisciplinary studies
1206 of human prehistory. It should also meet the aim and scope of *Nature*
1207 *Communications*.

1208 **Reference:**

1209 [1] Gao, Mingze, Chen Qiao, and Yuanhua Huang. "UniTVelo: temporally unified
1210 RNA velocity reinforces single-cell trajectory inference." Nature Communications
1211 13.1 (2022): 6586.

1212 [2] Riba, Andrea, et al. "Cell cycle gene regulation dynamics revealed by RNA
1213 velocity and deep-learning." Nature Communications 13.1 (2022): 2865.

1214 [3] Mazzoli, Mattia, et al. "Field theory for recurrent mobility." Nature
1215 communications 10.1 (2019): 3895.

1216 [4] Sohn, Byung-Ju, et al. "Regulation of atmospheric circulation controlling the
1217 tropical Pacific precipitation change in response to $CO_2$ increases." Nature
1218 communications 10.1 (2019): 1108.

1219

**Replies to Q2:**

We are grateful for the reviewer's suggestion. It greatly enhances the readability of our R codes and the convenience of the replications and utilizations of our approach by other users. Following the reviewer's suggestions, we have built an R package and provided some detailed tutorials on this package. Please see https://github.com/Stan-Sizhe-Yang/Language-velocity-field-estimation-for-language-dispersal-pattern-inference.

*Q3: Third, speaking of four, I hate to say this, but I was reviewing this study before, not negatively, but pointing to the code, and to other issues. Interestingly, the number of language families has now dropped from 7 to 4. How the hack did that happen? How do the authors explain that they discard three language families now? I know having the same reviewers for the same paper across journals is annoying, but please, good scientific practice requires you to be transparent and tell us what happened here. Did you discard them, because they did not bring the results you hoped for?*

**Replies to Q3:**

We appreciate the reviewer for pointing this out. Moreover, we are deeply

49

1251 grateful for the reviewer to dedicate valuable personal time to review our manuscript
1252 again. As mentioned by the reviewer, the previous version of our manuscript
1253 contained seven cases: Sino-Tibetan, Indo-European, Bantu, Arawak, Japonic,
1254 Austroasiatic, and Oceanic languages. However, for this version submitted to *Nature*
1255 *Communications*, we have excluded three language cases: Japonic, Austroasiatic, and
1256 Oceanic languages.

1257 **The primary reason for dropping three language cases**. We dropped these
1258 three language cases due to the lack of language samples around their suggested
1259 language homelands. To be specific, the proposed homelands of Indo-European,
1260 Sino-Tibetan, Bantu, and Arawak languages are situated in geographic ranges where
1261 sufficient language samples can be found [1-4]. However, there lack of sufficient
1262 language samples within the geographic areas covering the suggested homelands of
1263 Japonic (West Liao River of China [5-7]), Oceanic (Taiwan of China [6-8]), and
1264 Austroasiatic (Southern China [6]) languages respectively. Due to the lack of
1265 available language samples, it is nearly possible to determine the homelands of these
1266 three language cases in China solely based on the geographic coordinates of their
1267 language samples observed today. Accordingly, we can solely reconstruct the parts of
1268 their complete dispersal histories. The estimated results of these three language cases
1269 are described as follows.

1270 **The estimated results of three dropped language cases**. (i) The Japonic
1271 languages are regarded as the branch of the Trans-Eurasian languages [5]. Our
1272 approach traced their dispersal originating from the Honshu, followed by spread
1273 northward and southward across Japan. This dispersal pattern is in accordance with
1274 the expansion of the Trans-Eurasian languages from the Korean peninsula into Japan
1275 archipelago [5-7]. (ii) The Oceanic languages are a branch of the Austronesian
1276 languages [8]. We estimated their dispersal from the region near Southern Halmahera
1277 Island with subsequent eastward expansion across the Pacific settlement. The
1278 Southern Halmahera Island region is located at the easternmost edge of the
1279 geographic range of Oceanic language samples. Therefore, the estimated Oceanic
1280 dispersal pattern is compatible with the expansion of the Oceanic branch of the
1281 Austronesian language in the Pacific settlement [6-8]. (iii) For the Austroasiatic
1282 languages, our approach inferred their dispersal from the Mekong River region (one
1283 of the agricultural homelands in Mainland Southeast Asia), with subsequent
1284 expansion throughout Mainland Southeast Asia. This result favors the "Riverine

1285 hypothesis" proposed by Sidwell [9].

1286 Overall, due to a lack of sufficient language samples, the inferred dispersal
1287 patterns of Japonic, Oceanic, and Austroasiatic languages can only reflect a portion of
1288 their complete dispersal histories respectively. And, the estimated dispersal centers of
1289 these languages may be the secondary centers that are formed after they diffused into
1290 their current observed geographic ranges. Therefore, these three cases are unable to
1291 depict the full picture of their corresponding language dispersal patterns and illustrate
1292 the full power of our approach. More importantly, retaining these three language cases
1293 in our manuscript would make our narrative less clear which would potentially
1294 confuse the readers. In the version submitted to *Nature Communications*, we therefore
1295 decided to drop these three more troublesome cases.

1296 **Reference:**

1297 [1] Bouckaert, Remco, et al. "Mapping the origins and expansion of the
1298 Indo-European language family." Science 337.6097 (2012): 957-960.

1299 [2] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in
1300 northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

1301 [3] Grollemund, Rebecca, et al. "Bantu expansion shows that habitat alters the route
1302 and pace of human dispersals." Proceedings of the National Academy of Sciences
1303 112.43 (2015): 13296-13301.

1304 [4] Walker, Robert S., and Lincoln A. Ribeiro. "Bayesian phylogeography of the
1305 Arawak expansion in lowland South America." Proceedings of the Royal Society B:
1306 Biological Sciences 278.1718 (2011): 2562-2567.

1307 [5] Robbeets, Martine, et al. "Triangulation supports agricultural spread of the
1308 Transeurasian languages." Nature 599.7886 (2021): 616-621.

1309 [6] Diamond, Jared, and Peter Bellwood. "Farmers and their languages: the first
1310 expansions." science 300.5619 (2003): 597-603.

1311 [7] Skoglund, Pontus, and Iain Mathieson. "Ancient genomics of modern humans:
1312 the first decade." Annual review of genomics and human genetics 19 (2018): 381-404.

1313 [8] Gray, Russell D., Alexei J. Drummond, and Simon J. Greenhill. "Language
1314 phylogenies reveal expansion pulses and pauses in Pacific settlement." science
1315 323.5913 (2009): 479-483.

1316 [9] Paul, Sidwell. "The Austroasiatic central riverine hypothesis." Вопросы
1317 языкового родства 16 (59) (2010): 117-134.

1318

*Q4: Fourth, the claim of the method not using phylogenetic information is a bit exaggerated: we know geography correlates often with language relatedness (see for example here: https://doi.org/10.1371/journal.pone.0265460), so if geography explains the tree, you cannot say you do not use the tree if you use geography as a proxy for the construction of your vectors.*

**Replies to Q4:**

Thank you for your comments. At first, we fully agree with the reviewer that language geography usually strongly correlates with linguistic relatedness. The languages with closer geographic locations often possess higher relatedness due to either vertical divergence or horizontal contact. This connection guarantees the viability of various methods to reconstruct the dispersal pattern of languages based on linguistic relatedness, such as the phylogeographic approach [1] and our language velocity field estimation approach. To be specific, both the phylogeographic approach and our approach initially delineate the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. Subsequently, based on the correlation between linguistic relatedness and language geography, these evolutionary diachronic evolutionary trajectories are transformed into language dispersal trajectories.

Secondly, we would like to emphasize that our approach necessitates the phylogenetic information, but this phylogenetic information is not represented by the phylogenetic tree. To be specific, phylogenetic information or linguistic relatedness is not identical to the phylogenetic tree. It is noted that linguistic relatedness can be shaped by both vertical divergence and horizontal contact. The phylogenetic tree is just one of the models utilized to extract and represent the part of the linguistic relatedness of languages solely resulting from vertical divergence [2]. In our approach,

1344 we do not utilize the phylogenetic tree but a more general approach—the PCA
1345 algorithm to measure the linguistic relatedness through the distances among languages
1346 in a two-dimensional PC space (PCA-based distance). In the PC space, the languages
1347 exhibiting more linguistic relatedness resulting from either vertical divergence or
1348 horizontal contacts are intended to be distributed closer. Accordingly, if the linguistic
1349 relatedness is solely attributed to vertical divergence, the PCA-based distance should
1350 be able to capture the phylogenetic information similar to that of the phylogeographic
1351 tree.

1352 Thirdly, in our approach, we first depict the diachronic evolutionary trajectories
1353 of linguistic traits that shape the observed linguistic relatedness within the PC space.
1354 Based on the correlation between linguistic relatedness and language geography, we
1355 subsequently transform these diachronic evolutionary trajectories into language
1356 dispersal trajectories. Accordingly, we actually utilize language geography to
1357 approximate the linguistic relatedness for constructing the velocity field. Although the
1358 linguistic relatedness can be partially captured by the phylogenetic tree, it does not
1359 mean that our approach adopts the topological structure of the phylogenetic tree as
1360 input data used in our computational approach. **However, if the linguistic**
1361 **relatedness can be adequately captured by the phylogenetic tree, the**
1362 **phylogenetic information distilled by our approach should be similar to that**
1363 **distilled by the phylogenetic tree. Under this circumstance, our approach can be**
1364 **somehow regarded as utilizing the phylogenetic tree as well. In contrast, if**
1365 **linguistic relatedness bears more influence from horizontal contacts, our**
1366 **approach cannot be regarded as utilizing the phylogenetic tree. This conclusion**
1367 **has been verified in the revised main text (*Lines 210-303*).**

1368 **Reference**

1369 [1] Bouckaert, Remco, et al. "Mapping the origins and expansion of the
1370 Indo-European language family." Science 337.6097 (2012): 957-960.

1371 [2] François, Alexandre. "Trees, waves and linkages: Models of language
1372 diversification." The Routledge handbook of historical linguistics. Routledge, 2015.
1373 161-189.

1374

<i>Q5: Fifth, the question of homeland has always been problematic, but if you already use data by Wichmann and Rama, you should also check the much simpler baseline published in Glottolog by now (www.pyglottolog.readthedocs.io/en/latest/homelands.html#module-pyglottolog.homelands). This method seems to work as well as the one by Wichmann and Rama, but it is even simpler, so I would say there's one more baseline to be tested.</i>

**Replies to Q5:**

We highly value these insightful suggestions. Therefore, we compared our approach—language velocity field estimation (LVF) to two other baseline approaches suggested by the reviewer. These comparisons were achieved based on 1,000 simulated datasets and 4 empirical datasets. These two baseline approaches are referred to as "centroid (Centr)" and "minimal distance (MD)" approaches. The Centr approach postulates that the center of the polygon formed by the extension of current language geographic locations should be the dispersal center. The MD approach posits that the location of the language that exhibits the smallest average geographic distance to the other languages should be the dispersal center.

**1. Simulated validations for baseline approaches.**

It is noted that the simulated datasets are generated by applying a random walk model to the phylogenetic tree given a set of predefined dispersal centers. Accordingly, we have already known the true dispersal centers in these simulated datasets. Utilizing these simulated datasets provided by Wichmann et al., we first verified whether Centr and MD approaches can effectively estimate the predefined dispersal center. By applying Centr and MD approaches to the simulated datasets, we computed the errors in terms of longitude and latitude respectively between the true and estimated dispersal centers (Figure 1a to Q5). For either Centr or MD approaches, the outcomes of the Wilcoxon rank-sum test demonstrated that the errors between true and estimated dispersal centers were not significantly different from zero in both terms of longitude and latitude ($p$-value > 0.05; Figure 1a to Q5). It indicates that there is no difference between the dispersal centers estimated by either Centr or MD approaches and the true ones, thus affirming the high effectiveness of both Centr and MD approaches.

**2. Simulated comparisons between LVF and baseline approaches.**

54

1407      After justifying the effectiveness of the Centr and MD approach, we further
1408 compared the performance of LVF within these two approaches respectively based on
1409 1,000 simulated datasets. It is noted that the effectiveness of the LVF had already been
1410 verified using these simulated datasets in our previous manuscript. Therefore, we
1411 anticipated that LVF should exhibit the same performance as the Centr and MD
1412 approaches in simulated applications. Noting these, we calculated the differences in
1413 terms of longitude and latitude between the dispersal centers estimated by LVF and
1414 these two approaches respectively (Figure 1b to Q5). According to the Wilcoxon
1415 rank-sum test, we indeed found no significant differences in terms of longitude and
1416 latitude between the dispersal center estimated by LVF and those estimated by these
1417 two approaches respectively ($p$-value > 0.05; Figure 1b to Q5). This result confirms
1418 that LVF exhibits identical performance as these two baseline approaches in simulated
1419 applications.

1420 **3.    Empirical comparisons between LVF and baseline approaches.**

1421      We proceeded to compare the performance between LVF and baseline approaches
1422 in empirical applications. However, we found significant differences between the
1423 dispersal centers estimated by LVF and those estimated by these two baseline
1424 approaches (Figure 2 to Q5). Moreover, it appeared that the estimated dispersal
1425 centers of Centr and MD approaches seemed to lack support from the genetic and
1426 archeological evidence and were well less aligned with linguistics' conventional
1427 intuitions. In contrast, the estimated results of LVF can be more favored by the
1428 archaeological and genetic evidence, implying the better performance of LVF in
1429 empirical applications as compared to Centr and MD approaches.

1430 **4.    The possible reasons why two baseline approaches are useful in simulated**
1431 **validations but not in empirical applications.**

1432      Given the distinctions between the theoretical foundations of LVF and these two
1433 baseline approaches (i.e., Centr and MD), it is not surprising to see such obvious
1434 differences between the estimated result of LVF and those of the two baseline
1435 approaches in empirical applications. The LVF reconstructs language dispersal by
1436 transforming the diachronic evolutionary trajectories of linguistic traits that shape the
1437 observed linguistic relatedness into the language dispersal trajectories. In contrast,
1438 these two baseline approaches rely solely on the geographic locations of language

1439 samples, making their estimated results more susceptible to the biased geographic
1440 distribution of language samples. Nevertheless, these two baseline approaches exhibit
1441 high effectiveness in simulated validations probably owing to that simulated datasets
1442 are generated by the random walk model. The random walk model simulates that
1443 languages diffuse evenly as an outward radiating pattern from a given center.
1444 Accordingly, such simulation may display two characteristics:

1445 (a) The simulated language samples tend to be evenly distributed around this
1446    given dispersal center in the geographic space.
1447 (b) Due to (a), the simulated language samples located closer to the center of
1448    their geographic distribution would have a shorter average geographic
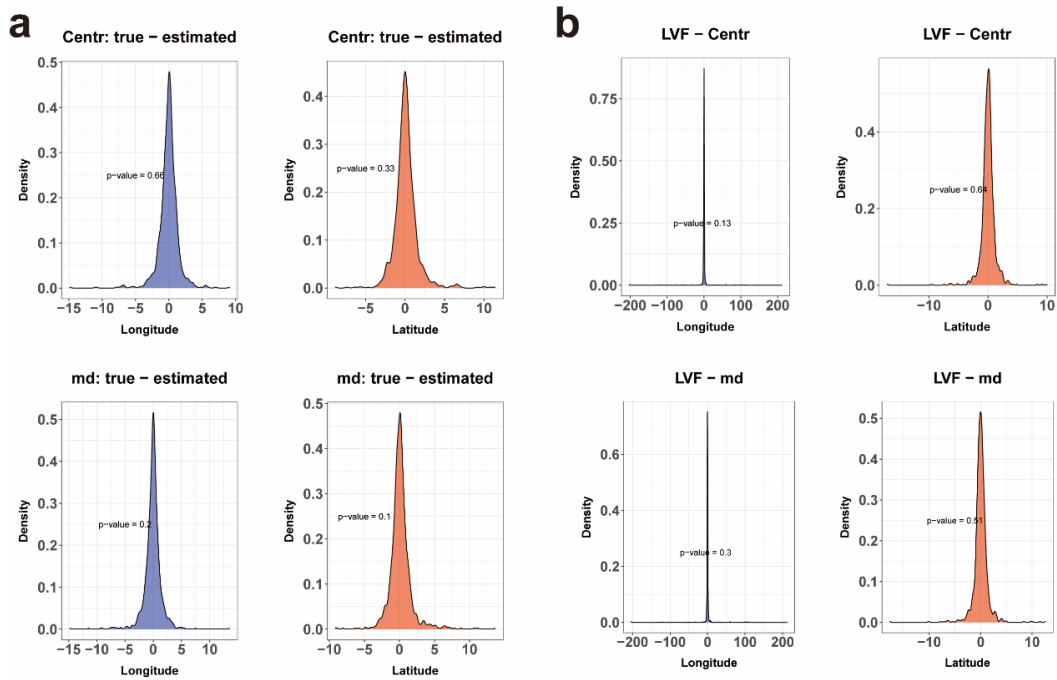1449    distance to other languages.

1450 Due to these two characteristics, both Centr and MD approaches can exhibit good
1451 performance in identifying the language dispersal center within simulated applications.
1452 Nevertheless, the empirical language samples may be not geographically distributed
1453 around the dispersal center uniformly, due to numerous reasons such as sampling bias,
1454 environmental constraints (i.e., mountain, desert, and river), and population
1455 movement (carrying languages out of the dispersal center) [1-2]. Consequently, Centr
1456 and MD approaches solely relying on the geographic locations of language samples
1457 may not perform as effectively in empirical applications.

1458 **Reference**

1459 [1] Grollemund, Rebecca, et al. "Bantu expansion shows that habitat alters the route
1460 and pace of human dispersals." Proceedings of the National Academy of Sciences
1461 112.43 (2015): 13296-13301.

1462 [2] Neureiter, Nico, et al. "Can Bayesian phylogeography reconstruct migrations and
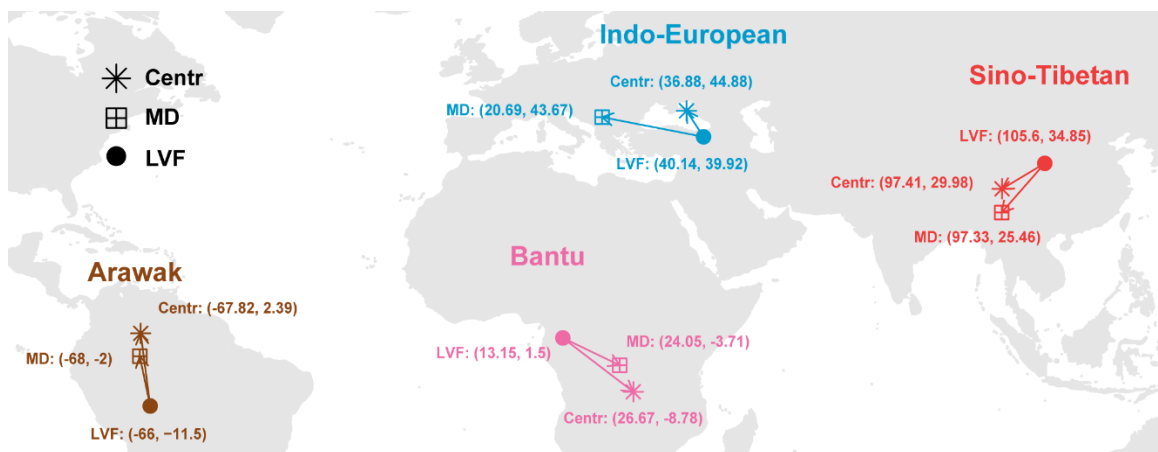1463 expansions in linguistic evolution?" Royal Society open science 8.1 (2021): 201079.

1464

1465

1466 **Figure 1 to Q5. Simulated validations of two baseline approaches and simulated**
1467 **comparisons between LVF and baseline approaches.** a) density plot shows the
1468 distribution of the error between the true and estimated dispersal center in terms of
1469 longitude and latitude. The p-value is calculated based on the Wilcoxon rank-sum test.
1470 b) density plot shows the distribution of the difference between the dispersal center
1471 estimated by LVF and baseline approaches in terms of longitude and latitude. The
1472 p-value is calculated based on the Wilcoxon rank-sum test.

1473



1474

1475 **Figure 2 to Q5. The dispersal centers estimated by LVF, Centr, and MD**

1476 **approaches for four language families and groups.**

1477

1478 *Q6: And when speaking of testing: why restrict your study to four datasets (or*
1479 *seven), if there are many more available in terms of phylogenies now, which are all*
1480 *with nicely coded cognate sets in standardized data formats (see e.g.,*
1481 *https://doi.org/10.1038/s41597-022-01432-0 for a very large collection of*
1482 *standardized data)? It seems the data has been cherry-picked to yield good results.*
1483 *Taking ten of the datasets in the Lexibank collection should not be difficult and*
1484 *would tell us much more clearly where we are with this new method.*

1485 **Replies to Q6:**

1486    We express our sincere gratitude to the reviewer for introducing the *Lexibank*
1487 which is an important lexical dataset to us. The *Lexibank* covers nearly 3,000
1488 language samples of around 300 language families and groups around the world. This
1489 lexical dataset could provide comprehensive insights into the origins and dispersals of
1490 various language families and groups around the world.

1491    The primary objective of our paper is to examine the alignment of language
1492 dispersal, demic diffusion, and Neolithic/Agricultural cultures spread in human
1493 prehistory. Therefore, the language cases utilized in our paper are expected to fulfill
1494 the following criteria. Firstly, the language case should have a possible association
1495 with the origin and development of ancient agriculture. Secondly, the demic or
1496 cultural diffusions in the specific geographic areas where these languages are spoken
1497 should be supported by corresponding genetic or archaeological evidence. Thirdly, the
1498 language cases are preferably renowned cases with sufficient language samples that
1499 have been rigorously investigated in previous phylogenetic research. More
1500 importantly, the lexical items in these language cases should have been carefully
1501 collated and well coded into cognate sets that meet the standard of computational
1502 linguistics. With these criteria, we hope that the empirical cases can better serve our
1503 paper's primary objective and make our estimated results more acceptable to the
1504 broad range of audiences.

1505    According to these criteria, four language cases which are Indo-European,
1506 Sino-Tibetan, Bantu, and Arawak languages are included in our study. These

1507 languages are hypothesized to be closely associated with agricultural development in
1508 this area [1-2]. Moreover, they are widely spoken in their corresponding geographic
1509 area and have all been rigorously studied by former phylogenetic studies [3-6]. More
1510 importantly, the lexical items utilized in these four cases have undergone careful
1511 selections and validations. In the geographic areas where the languages are spoken,
1512 the demic diffusion and cultural spread have been delineated based on sufficient
1513 genetic or archaeological evidence [1-2].

1514 Following these criteria, within the *Lexibank*, we filtered out the language cases
1515 with a sample size lower than 20, ultimately leaving us with 17 language cases. These
1516 cases are *Afro-Asiatic, Arawak, Atlantic-Congo, Austroasiatic, Austronesian,*
1517 *Hmong-Mien, Indo-European, Nuclear Trans New Guinea, Pama-Nyungan,*
1518 *Quechuan, Sino-Tibetan, Dravidian, Tucanoan, Tupian, Turkic, Uralic, and*
1519 *Uto-Aztecan languages*. Among them, the *Indo-European, Sino-Tibetan, Austroasiatic,*
1520 *and Arawak languages* have been incorporated into our study and *Afro-Asiatic and*
1521 *Pama-Nyungan languages* are the hunter-gatherer languages. Additionally, there lack
1522 of sufficient Austronesian language samples within their suggested homeland in China.
1523 Therefore, we ultimately selected 10 language cases: *Uralic, Trans-New-Guinea,*
1524 *Quechuan, Turkic, Tukanoan, Tupian, Uto-Aztecan, Hmong-Mien, Atlantic-Congo,*
1525 *and Dravidian languages*.

1526 However, either the evolution or dispersals of these 10 language cases has not
1527 been well investigated and remains highly controversial in the previous computational
1528 linguistic studies. Therefore, investigating their dispersal patterns seems worthy of
1529 being pursued as separate research endeavors for publication. Moreover,
1530 corresponding genetic and archaeological evidence is also hard to find to support the
1531 demic diffusion and cultural spread within the area where these languages are spoken.
1532 Given these constraints, we hold the view that including these 10 language cases in
1533 our study may not align with our primary research objective and make the narrative of
1534 our manuscript less clear. Therefore, we still hope to retain the original well-attested
1535 four language cases (i.e., Indo-European, Sino-Tibetan, Bantu, and Arawak languages)
1536 in our manuscript.

1537 Although we have decided not to include these language cases in our revision, we
1538 still have applied our approach—language velocity field estimation (LVF) to these
1539 language cases to infer their dispersal patterns. In this reply, we present the results

1540 regarding the dispersal patterns of these 10 language cases to the reviewer below

1541 (Table 1 to Q6 and Figure 1 to Q6). The datasets of these 10 language cases and the R

1542 codes for replicating the results of these 10 language cases can be downloaded from

1543 https://github.com/Stan-Sizhe-Yang/Language-velocity-field-estimation-for-language

1544 -dispersal-pattern-inference.

## 1. Uralic languages

1546 Uralic languages are widely distributed across northeastern Europe and Northern

1547 Asia. The lexical dataset of Uralic languages was sourced from Honkola et al. (2013)

1548 [7]. The LVF inferred that the dispersal center of Uralic languages is situated in the

1549 steppe region in the southeast of the Ural Mountains (Lon: 64.6, Lat: 54.9) (Figure 1b

1550 to Q6). From this dispersal center, Uralic languages dispersed westward crossing the

1551 Ural Mountains into Europe and eastward into the Far East region. It advocates the

1552 "east of the southern Urals origin hypothesis" of Uralic languages, which is proposed

1553 according to the historical contact between Uralic and Indo-Iranian languages [8].

## 2. Trans-New-Guinea languages

1555 Trans–New Guinea languages are widely spoken on the island of New Guinea

1556 and neighboring islands. The Trans–New Guinea lexical dataset was obtained from

1557 Greenhill (2015) [9]. The LVF depicted the dispersal of Trans-New-Guinea languages

1558 originating from the center in central Papua New Guinea (Lon: 144.3, Lat: -6.4),

1559 which used to be the ancient agricultural homeland of New Guinea island (Figure 1c

1560 to Q6). This result is compatible with the conclusion drawn from recent linguistic

1561 studies and corroborated by the archaeological evidence [10-11]. It suggests that the

1562 Trans–New Guinea dispersal could be closely associated with the development and

1563 spread of agriculture across the New Guinea island.

## 3. Quechuan languages

1565 The Quechuan languages are widely spoken by the native peoples in South

1566 America. We collected the Quechuan lexical dataset from the Blum et al. (2023) [12].

1567 The dispersal center of Quechuan languages (Lon: -75.5, Lat: -9.8) was inferred more

1568 adjacent to the Lima near the Andes which is the ancient agricultural homeland in

1569 South America [13] (Figure 1d to Q6). From this dispersal center, Quechuan

1570 languages spread northward and southward along the Andes. These results are

1571     compatible with the evidence drawn from the Quechua dialectology [14].

1572     **4. Turkic languages**

1573     Turkic languages span the vast expanse of the Eurasian continent, stretching from
1574     the northwest of China to the west of Eastern Europe, and from the north of Siberia to
1575     the south of Iran. The precise homeland of Turkic languages remains a subject of
1576     intense debate. The expansive geographic area encompassing the Transcaspian steppe
1577     to the far northeastern reaches of Manchuria in Asia is regarded as a potential
1578     homeland for these languages [15]. We applied LVF to the Turkic lexical dataset
1579     structured by Savelyev et al. (2020) [16]. The spatial reconstruction showed that
1580     Turkic languages spread westward into Europe and eastward into the Far East region
1581     from the dispersal center inferred in Kazakhstan near Mongolia and Southern Siberia
1582     (Lon: 77.1, Lat: 54.4) (Figure 1e to Q6). This result can be advocated by the genetic
1583     evidence that suggests the potential origin of Turkic-speaking populations in the area
1584     near Mongolia and Southern Siberia [15]. However, we noticed that the Turkic
1585     language samples manifested an exceedingly sparse geographic distribution across the
1586     Eurasian continent. Such sparse geographic distribution may introduce more
1587     uncertainties into the LVF estimation. Therefore, collecting more Turkic language
1588     samples may enable LVF to yield a more precise depiction of the Turkic dispersal
1589     pattern.

1590     **5. Tukanoan languages**

1591     Tukanoan, also referred to as Tucanoan, is a language family of Colombia, Brazil,
1592     Ecuador, and Peru in South America. We applied the LVF to the Tucanoan dataset
1593     derived from Chacon et al. (2017) [17]. The dispersal center of Tucanoan languages
1594     was inferred in the region of the Japurá River (Lon: -70.0, Lat: -0.9) (Figure 1f to Q6).
1595     The location of this dispersal center is compatible with the conclusion drawn from
1596     previous linguistic studies and can be advocated by the archaeological evidence
1597     [17-18].

1598     **6. Tupian languages**

1599     The Tupian language family is one of the largest linguistic groups in South
1600     America. The dataset of the Tupian language was sourced from Galucio et al. (2015).
1601     We applied LVF to this dataset for inferring the dispersal pattern of Tupian languages.

1602 The result showed that Tupian languages dispersed from the center located in the
1603 regions of Rondônia in Brazil within the Madeira River basin (Lon: -62.3, Lat: -11.6)
1604 across South America. This result is compatible with previous linguistic studies [19]
1605 (Figure 1g to Q6).

1606 **7. Uto-Aztecan languages**

1607 The Uto-Aztecan languages are the mother tongue of native Americans, which are
1608 primarily spoken in the Great Basin region, including states such as California,
1609 Nevada, and Arizona, and extending into Mexico. The Uto-Aztecan lexical dataset
1610 was derived from the Greenhill (2023) [20]. The LVF identified the dispersal center of
1611 Uto-Aztecan languages in Southern Arizona (Lon: -113.5, Lat: 33.9) near the border
1612 between Arizona and Mexico (Figure 1h to Q6). This location was compatible with
1613 the one inferred by the phylogeographic approach as reported in Greenhill (2023)
1614 (Lon: -116.7, Lat: 34.8). From this dispersal center, the Uto-Aztecan languages spread
1615 southeastward and northwestward along the coastline, and northeast into South
1616 America. These results favor the "Northern origin hypothesis" supported by the
1617 reconstruction of flora and fauna terms [20-21]. This hypothesis postulates that
1618 Uto-Aztecan languages originated in the area between Southern California's Mojave
1619 Desert and the Sonoran and Chihuahuan desert regions of Arizona and northern
1620 Mexico.

1621 **8. Hmong-Mien languages**

1622 The Hmong-Mien languages are primarily spoken by various ethnic groups in
1623 southern China, northern Vietnam, Laos, Thailand, and Myanmar. Linguistic
1624 reconstructions focusing on ancient terminology related to flora and fauna have
1625 suggested that the origins of Hmong-Mien languages might be found in the provinces
1626 to the south of the Yangzi River [22]. In our investigation, we applied the LVF to the
1627 Hmong-Mien lexical dataset derived from Chen (2013) [23] (Figure 1i to Q6). The
1628 results consistently indicated that the dispersal center of Hmong-Mien languages is
1629 indeed located within Guizhou province, situated to the south of the Yangzi River
1630 (Longitude: 107.7, Latitude: 27.0).

1631 **9. Atlantic-Congo languages**

1632 The Atlantic-Congo languages, which constitute a prominent subgroup of the

Niger-Congo language family, have a significant presence across the African continent. The Atlantic-Congo lexical dataset was collected from the public dataset compiled by Koelle (1853) [24]. Utilizing the LVF, we traced the dispersal of Atlantic-Congo languages initiating from Nigeria near Cameroon (Lon: 5.6, Lat: 6.4), which used to be the ancient agricultural homeland in Africa [25] (Figure 1j to Q6). It suggests that the Atlantic-Congo dispersal could be associated with agricultural expansion in Africa.
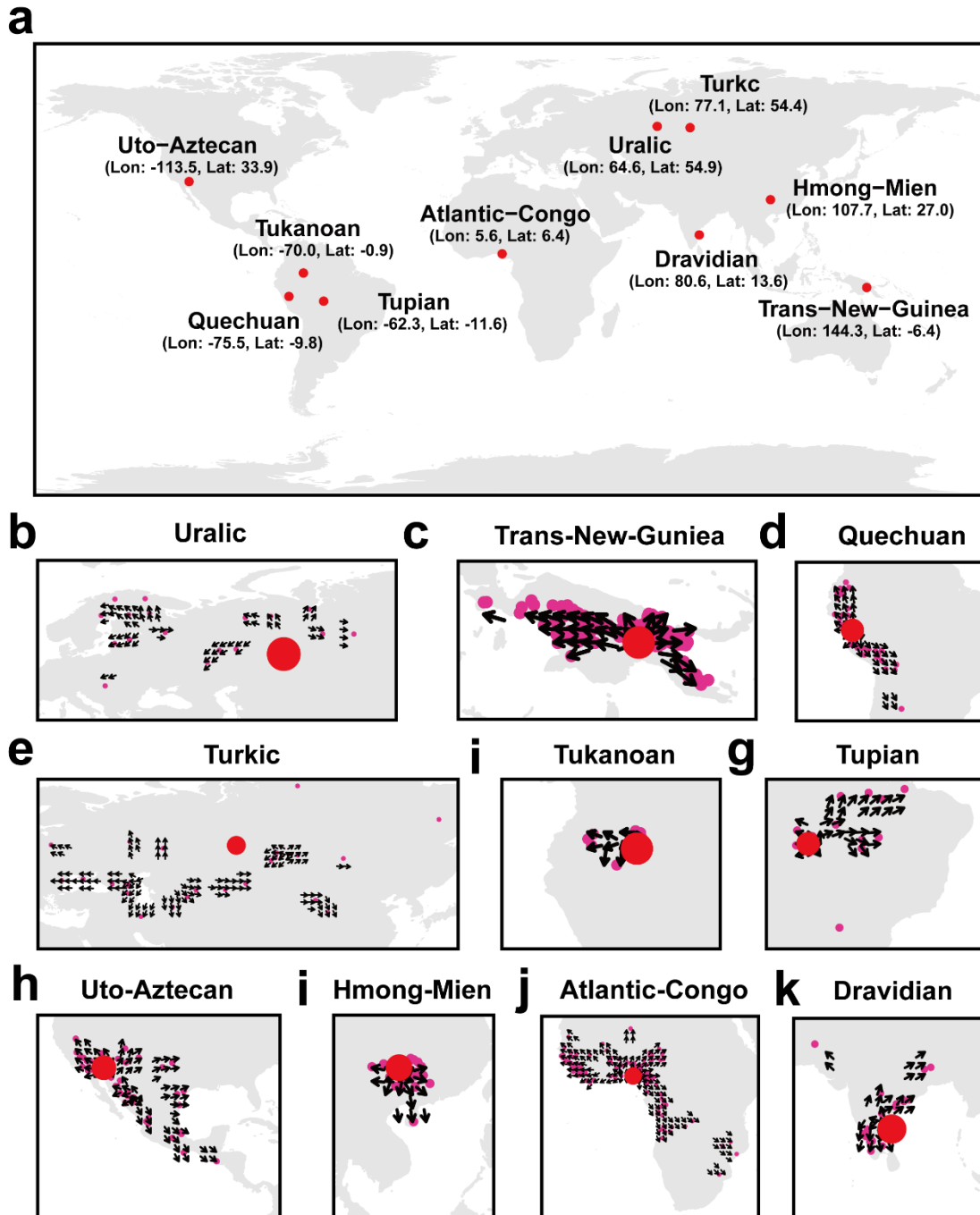
**10. Dravidian languages**

The Dravidian languages are widely scattered across southern and central India and surrounding countries. The dispersal of Dravidian languages has been a long-standing debate. The genetic evidence indicates the potential origin of Dravidian languages in the Indus Valley, with subsequent southward and eastward expansion across the Indian subcontinent [26]. The linguistic evidence drawn from the term reconstruction suggests that Dravidian languages might originate somewhere in South India (i.e., Peninsular India) [26]. Archaeological evidence yields the connection between the origin of the Dravidian language and the development of the Southern Neolithic complex in Karnataka and Andhra Pradesh [27, 28]. Based on the Dravidian lexical dataset derived from Kolipakam et al. (2018) [29], LVF inferred the dispersal of Dravidian languages originating from the center located in the range of Andhra Pradesh (Lon: 80.6, Lat: 13.6) (Figure 1k to Q6). This result can be supported by the archaeological evidence that implies the close association between Dravidian dispersal and Neolithic culture spread in India.

**Table and Figure**

**Table 1 to Q6.** The coordinates of dispersal centers inferred by LVF for ten language families and groups.

| Language | Longitude | Latitude |
|---|---|---|
| Uralic | 64.6 | 54.9 |
| Trans-New-Guinea | 144.3 | -6.4 |
| Quechuan | -75.5 | -9.8 |
| Turkic | 77.1 | 54.4 |
| Tukanoan | -70.0 | -0.9 |

| | | |
|---|---|---|
| Tupian | -62.3 | -11.6 |
| Uto-Aztecan | -113.5 | 33.9 |
| Hmong-Mien | 107.7 | 27.0 |
| Atlantic-Congo | 5.6 | 6.4 |
| Dravidian | 80.6 | 13.6 |

**a**



**b** Uralic  **c** Trans-New-Guniea  **d** Quechuan

**e** Turkic  **i** Tukanoan  **g** Tupian

**h** Uto-Aztecan  **i** Hmong-Mien  **j** Atlantic-Congo  **k** Dravidian

1658

**Figure 1 to Q6.** The Language velocity fields reveal the dispersal patterns of 10

language families and groups worldwide. The red dot denotes the dispersal center inferred by LVF. The pink dot signifies the language sample. The black arrow represents the grid-smoothed velocity vector.

**Reference**

[1] Diamond, Jared, and Peter Bellwood. "Farmers and their languages: the first expansions." science 300.5619 (2003): 597-603.

[2] Skoglund, Pontus, and Iain Mathieson. "Ancient genomics of modern humans: the first decade." Annual review of genomics and human genetics 19 (2018): 381-404.

[3] Bouckaert, Remco, et al. "Mapping the origins and expansion of the Indo-European language family." Science 337.6097 (2012): 957-960.

[4] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

[5] Grollemund, Rebecca, et al. "Bantu expansion shows that habitat alters the route and pace of human dispersals." Proceedings of the National Academy of Sciences 112.43 (2015): 13296-13301.

[6] Walker, Robert S., and Lincoln A. Ribeiro. "Bayesian phylogeography of the Arawak expansion in lowland South America." Proceedings of the Royal Society B: Biological Sciences 278.1718 (2011): 2562-2567.

[7] Honkola, Terhi, et al. "Cultural and climatic changes shape the evolutionary history of the Uralic languages." Journal of Evolutionary Biology 26.6 (2013): 1244-1253.

[8] Nichols, Johanna. "The origin and dispersal of Uralic: Distributional typological view." Annual Review of Linguistics 7 (2021): 351-369.

[9] Greenhill, Simon J. "TransNewGuinea. org: An online database of New Guinea languages." PLoS One 10.10 (2015): e0141563.

[10] Diamond, Jared, and Peter Bellwood. "Farmers and their languages: the first expansions." science 300.5619 (2003): 597-603.

1687 [11]Schapper, Antoinette. "Farming and the Trans-New Guinea family." Language
1688 dispersal beyond farming (2017): 155-181.

1689 [12]Blum, Frederic, et al. "A phylolinguistic classification of the Quechua language
1690 family." (2023).

1691 [13]Diamond, Jared, and Peter Bellwood. "Farmers and their languages: the first
1692 expansions." science 300.5619 (2003): 597-603.

1693 [14]King, Kendall A., and Nancy H. Hornberger. "Quechua as a lingua franca."
1694 Annual Review of Applied Linguistics 26 (2006): 177-196.

1695 [15]Yunusbayev, Bayazit, et al. "The genetic legacy of the expansion of
1696 Turkic-speaking nomads across Eurasia." PLoS Genetics 11.4 (2015): e1005068.

1697 [16]Savelyev, Alexander, and Martine Robbeets. "Bayesian phylolinguistics infers the
1698 internal structure and the time-depth of the Turkic language family." Journal of
1699 Language Evolution 5.1 (2020): 39-53.

1700 [17]Chacon, Thiago. "Arawakan and Tukanoan contacts in Northwest Amazonia
1701 prehistory." PAPIA Rev. Bras. Estud. Crioulos E Similares 27 (2017): 237-265.

1702 [18]Chacon, Thiago. "On Proto-Languages and Archaeological Cultures: pre-history
1703 and material culture in the Tukanoan Family." Revista Brasileira de Linguística
1704 Antropológica 5.1 (2013): 217-245.

1705 [19]Galucio, Ana Vilacy, et al. "Genealogical relations and lexical distances within
1706 the Tupian linguistic family." Boletim do Museu Paraense Emílio Goeldi. Ciências
1707 Humanas 10 (2015): 229-274.

1708 [20]Greenhill, Simon J., et al. "A recent northern origin for the Uto-Aztecan family."
1709 Language (2023).

1710 [21]Campbell, Lyle. "What drives linguistic diversification and language spread."
1711 Examining the farming/language dispersal hypothesis (2002): 49-63.

1712 [22]Ratliff, Martha. Hmong-Mien language history. Pacific Linguistics, Research
1713 School of Pacific and Asian Studies, The Australian National University, 2010.

[23] Qiguang, Chen. "Miao-Yao yuwen [Miao and Yao Language]." Beijing: Zhongyang Minzu Daxue Chubanshe (2013).

[24] Koelle, Sigismund Wilhelm. "Polyglotta Africana; or a comparative vocabulary of nearly three hundred words and phrases in more than one hundred distinct African languages." (No Title) (1853).

[25] Diamond, Jared, and Peter Bellwood. "Farmers and their languages: the first expansions." science 300.5619 (2003): 597-603.

[26] Narasimhan, Vagheesh M., et al. "The formation of human populations in South and Central Asia." Science 365.6457 (2019): eaat7487.

[27] Southworth, Franklin. Linguistic archaeology of south Asia. Routledge, 2004.

[28] Kolipakam, Vishnupriya, et al. "A Bayesian phylogenetic study of the Dravidian language family." Royal Society open science 5.3 (2018): 171504.

[29] Kolipakam, Vishnupriya, et al. "DravLex: A Dravidian lexical database:(Version v1. 0.0)[Data set]." (2018).

*Q7: Sixth, the method has the rather infelicitous name "language velocity field estimation", and I could not find any explanation why the authors chose to call it like that, since the name is very confusion and difficult to parse, and it does not really help to understand what the method could be about. I think in general it would be useful to 1) change the name to something that explains the method in a better way (dynamic trait vectors? I am not sure) and 2) to explain the method in much, much more detail. For this, figures would be needed that show how vectors for some of the traits are estimated, and the authors would need to also check the resulting vectors on an individual basis in order to see if they make sense.*

**Replies to Q7:**

We are sorry for not being clear about the rationale of our approach. After careful consideration, we have decided to retain the original name "language velocity field" of our approach. Because this name can intuitively reflect the characteristics of our

1742 approach. Following the valuable *suggestion 2)* offered by the reviewer, we have
1743 redrawn our original schematic diagrams for the rationale and calculation procedure
1744 of our approach with greater detail and accuracy as shown in Figure 1 of the revised
1745 main text. For the convenience of the reviewer, we attach Figure 1 of the revised main
1746 text to the end of this reply as Figure 1 to Q7. Additionally, we have added more
1747 detailed descriptions of our approach into the *Lines 109-151* of the revised main text.
1748 Considering the word limit in the main text, more detailed explanations of our
1749 approach can be found in Supplementary Note 1. Here, we provide a concise
1750 explanation of the rationale of our approach.

1751     **The inspiration for proposing language velocity field estimation.** The velocity
1752 field can be visualized as a collection of arrows with given magnitudes and directions
1753 estimated by a specific dynamic model, which demonstrates the directions of the
1754 spatiotemporal changes of individuals [1]. The directions of the vectors in the velocity
1755 field compose sets of continuously changing paths that visualize the dynamic
1756 trajectories of natural phenomena such as atmospheric circulation [2] (e.g., water
1757 vapor transport), and cell differentiation [3] (e.g., RNA transcription). Furthermore,
1758 this approach has now extended to infer the trajectories of the spatial-temporal
1759 changes of social phenomena such as demic diffusion [4] (e.g., human mobility), and
1760 cultural spread [5] (e.g., Neolithic culture propagation). Given that humans are the
1761 carriers of languages which are also the carriers of cultures, we believe that the
1762 velocity field could also contribute to the inference of the language dispersal.
1763 Accordingly, our approach is designed to establish a language velocity field on the
1764 geographic map to depict language dispersal patterns. By visualizing the language
1765 velocity field on the geographic map, the directions of velocity vectors can intuitively
1766 show how and from where (i.e., dispersal trajectory and center) these languages have
1767 dispersed into their current locations.

1768     **Our approach shares the same theoretical foundation as the phylogeographic**
1769 **approach but with different implementation strategies.** As the most prevailing
1770 approach, the phylogeographic approach implements two major steps to infer
1771 language dispersal from the diachronic evolution of linguistic traits [6]. The first is to
1772 establish a phylogenetic tree to depict the diachronic evolutionary trajectories of
1773 linguistic traits that shape the observed linguistic relatedness (Figure 2 to Q7). The
1774 second is to project the phylogenetic tree into the geographic space to transform these
1775 diachronic evolutionary trajectories into dispersal trajectories, based on the correlation

68

between linguistic relatedness and geography (Figure 2 to Q7). Akin to the phylogeographic approach, our approach also infers language dispersal through the diachronic evolution of linguistic traits with two major steps (Figure 2 to Q7). The first is to establish a velocity field to depict the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. The second is to project this velocity field into the geographic space to outline the language dispersal trajectories. These two steps are described as follows.

**The velocity field in PC space delineates diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness.** Our approach conducts the PCA-based distance rather than a phylogenetic tree to represent linguistic relatedness. To be specific, we employ the PCA algorithm to extract two optimal principal components (i.e., PC1 and PC2) from the linguistic traits. According to PC1 and PC2, we represent the linguistic relatedness among language samples as the distances among them in the PC space that can be shaped by both divergence and contact (Figure 1b to Q7). In parallel, we use a dynamic model, similar to the widely-used covarion model for linguistic trait evolution [7-9], to reconstruct the past states of linguistic traits for each language sample (Figure 1d1 to Q7). Given the differences between the past and current trait states of each language sample, we can obtain a velocity vector that reflects the direction of diachronic changes in its linguistic traits (Figure 1d2 to Q7). In other words, the velocity vector depicts how the linguistic traits in each language sample evolve into their current states. Finally, we project this language velocity field into the PC space formed by the aforementioned two principal components (Figure 1e to Q7). For convenience, we can interpret the language velocity field in the PC space as the collection of arrows connecting the past and current states of linguistic traits within language samples in the PC space (Figure 1e1 to Q7). Accordingly, the past and current states of linguistic traits within language samples can simultaneously be visualized in the PC space. Each arrow connecting the past and current states of linguistic traits for each language sample outlines the diachronic change of the linguistic traits in this language. Therefore, the arrows in the PC space compose a set of trajectories to depict the diachronic evolution of the linguistic traits that shape the observed linguistic relatedness (Figure 1e2 to Q7).

**Transforming the diachronic evolutionary trajectories of the linguistic traits into language dispersal trajectories.** We project the language velocity field from the PC space to the geographic space based on the correlation between linguistic
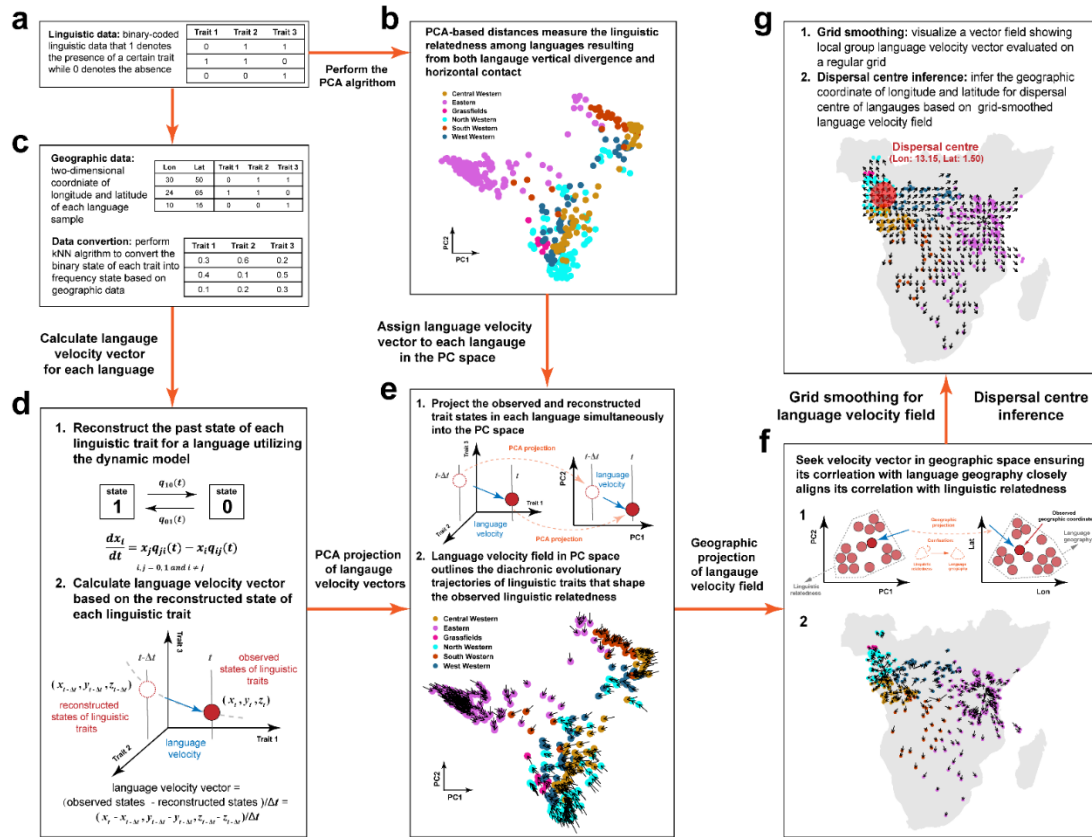
1810    relatedness and geography [6-8] (Figure 1f to Q7). To achieve this, we utilize the
1811    kernel projection approach proposed by La Manno et al. [3] to project the language
1812    velocity field from the PC space into the two-dimensional geographic space. The
1813    rationale behind this kernel projection is to estimate the velocity vectors of language
1814    samples in the geographic space, ensuring that their correlation with language
1815    distributions in the PC space can be best preserved within the geographic space
1816    (Figure 1f1 to Q7). This projection is similar to the projection of the phylogenetic tree
1817    to the geographic space in the phylogeographic approach. Accordingly, the directions
1818    of these vectors compose a set of trajectories that depict from where the observed
1819    language samples have diffused into their current locations (Figure 1f2 to Q7). We
1820    hope these contents supplemented by Figure 1 to Q7 can provide the reviewer with a
1821    clearer understanding of our approach.

1822    **Validation of velocity field.** The direction of the ultimate velocity vector of a
1823    language sample we estimated within the geographic space manifests the direction
1824    from where this language sample diffuses into its current locations. However, it is
1825    important to highlight that the power of any spatial reconstruction method is
1826    inevitably affected by the heterogeneity of the spatial distribution of samples.
1827    Therefore, each estimated velocity vector cannot signify exactly the diffusion
1828    direction of each language sample. However, our approach aims to reconstruct the
1829    general dispersal pattern of the entire language family or group rather than the exact
1830    dispersal direction of just one language sample. Moreover, relying solely on a single
1831    velocity vector is insufficient to ascertain the dispersal pattern of the entire language
1832    family. And, the overall dispersal pattern of the entire language family is deduced by
1833    the continuously changing trajectories formed by a collection of velocity vectors.
1834    Consequently, it appears less critical to validate the effectiveness of a solitary velocity
1835    vector on the individual level. Accordingly, we consider that the effectiveness of our
1836    approach should be validated on the global level of the language velocity field rather
1837    than the individual level of a single language velocity vector. Under this circumstance,
1838    simulated validations of our approach have confirmed its ability to reconstruct
1839    accurate language dispersal patterns based on the language velocity field in our
1840    previous manuscript. Therefore, with these simulated validations, we believe that the
1841    velocity vectors can indeed contribute to reconstructing the language dispersal
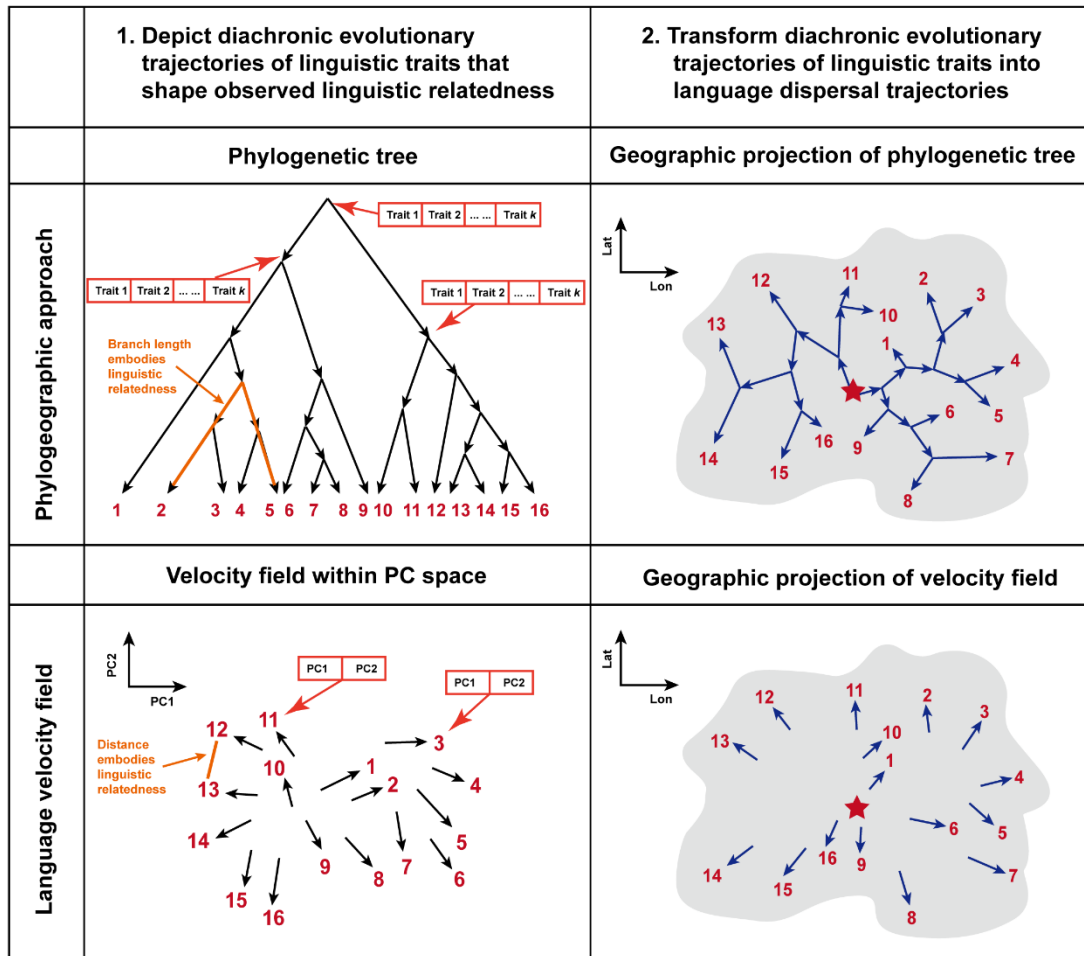1842    pattern.

1843    **Reference**

1844 [1] Galbis, Antonio, and Manuel Maestre. Vector analysis versus vector calculus.
1845 Springer Science & Business Media, 2012.

1846 [2] Sohn, Byung-Ju, et al. "Regulation of atmospheric circulation controlling the
1847 tropical Pacific precipitation change in response to CO2 increases." Nature
1848 communications 10.1 (2019): 1108.

1849 [3] La Manno, Gioele, et al. "RNA velocity of single cells." Nature 560.7719 (2018):
1850 494-498.

1851 [4] Mazzoli, Mattia, et al. "Field theory for recurrent mobility." Nature
1852 communications 10.1 (2019): 3895.

1853 [5] Fort, Joaquim. "Demic and cultural diffusion propagated the Neolithic transition
1854 across different regions of Europe." Journal of the Royal Society interface 12.106
1855 (2015): 20150166.

1856 [6] Bouckaert, Remco, et al. "Mapping the origins and expansion of the
1857 Indo-European language family." Science 337.6097 (2012): 957-960.

1858 [7] Yang, Ziheng. "Maximum-likelihood estimation of phylogeny from DNA
1859 sequences when substitution rates differ over sites." Molecular biology and evolution
1860 10.6 (1993): 1396-1401.

1861 [8] Penny, David, et al. "Mathematical elegance with biochemical realism: the
1862 covarion model of molecular evolution." Journal of Molecular Evolution 53 (2001):
1863 711-723.

1864 [9] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in
1865 northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

1866 **Figure**

1867

**Figure 1 to Q7. Schematic diagram of language velocity field estimation (LVF) for inferring the dispersal trajectories and centers of languages.** The computational procedures of the LVF comprise two major steps. Subfigures (a) to (e) illustrate the first step which is to estimate a velocity field on the PC space to outline the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. Subfigures (f) to (g) illustrate the second step, which is to project the velocity field from PC space into geographic space. Within the velocity field in geographic space, the directions of the velocity vectors compose a set of continuously changing trajectories that delineate from where these languages diffuse to their current locations. These procedures are exemplified using the Bantu language family. Comprehensive insights into the underlying principles and computational steps can be found in the Materials and Methods section, as well as Supplementary Note 1.

1881

**Figure 2 to Q7. Language velocity field estimation (LVF) shares the same foundation as the phylogeographic approach but with different implementation strategies.** Both LVF and phylogeographic approach entails two major steps to infer language dispersal pattern. The first is to depict the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. The second is to transform these diachronic evolutionary trajectories of linguistic traits into language dispersal trajectories. In the phylogenetic tree, each language is determined by $k$ linguistic traits. In the velocity field within PC space, each language is determined by PC1 and PC2 which are rearranged from the $k$ linguistic traits through the PCA algorithm. The red number denotes a language. The black arrow signifies the evolutionary direction of linguistic traits in a language. The blue arrow represents the dispersal direction of a language. The red star denotes the estimated dispersal center.

1895 *Q8: Seventh, the authors praise their method for not needing trees, but at the same*
1896 *time, they do not tell the readers why trees are so useful: they tell us various*
1897 *scenarios of character evolution in a very transparent way, in which we have*
1898 *scenario and can plot how the trait evolved. Of course, this is not always done, but*
1899 *they should tell the readers to which the method they propose allows us to get some*
1900 *insights into the black box, since a simple black box, even if it works, is not*
1901 *satisfying from a scientific viewpoint, and we talk about scientific approaches here.*

1902 **Replies to Q8:**

1903    We really appreciate the reviewer for raising this crucial point. To improve the
1904 credibility and interpretability of our approach, we have added more comprehensive
1905 descriptions and explanations of our approach to the revised main text (*Line*s
1906 109-151). Here, we offer a brief answer.

1907 **1. The phylogenetic tree visualizes the diachronic evolutionary trajectories of**
1908    **the linguistic traits that shape the observed linguistic relatedness.**

1909    The phylogeographic approach infers the language dispersal through the
1910 diachronic evolution of linguistic traits. As the reviewer mentioned, the phylogenetic
1911 tree plays an important role in the phylogeographic approach. To be specific, the
1912 phylogenetic tree is a power representation for the diachronic evolutionary trajectories
1913 of the linguistic traits that shape the observed linguistic relatedness (Figure 1 to Q8).
1914 This representation relies on the branching pattern within the phylogenetic tree. This
1915 branching pattern visualizes the diachronic evolution of linguistic traits in languages
1916 after diverging from their ancestors [1]. The shorter branch linking two languages
1917 indicates fewer diachronic changes occurring between their traits, resulting in a higher
1918 linguistic relatedness between them. This phylogenetic tree can be projected into the
1919 geographic space based on the correlation between linguistic relatedness and language
1920 geography (Figure 1 to Q8) [1-2]. To be specific, each branch within the phylogenetic
1921 tree, that has been projected into the geographic space, is regarded as a segment of the
1922 dispersal trajectories (Figure 1 to Q8). With this projection, the evolutionary
1923 trajectories of linguistic traits can thus be transformed into language dispersal
1924 trajectories.

1925 **2. The theoretical foundation and interpretability of our approach.**

74

1926      Akin to the phylogeographic approach, our approach also aims to reconstruct the
1927 language dispersal pattern through the diachronic evolution of linguistic traits. Our
1928 approach and phylogeographic approach actually share the same theoretical
1929 foundation but with different implementation strategies (Figure 1 to Q8).

1930      **The velocity field in PC space depicts the diachronic evolutionary**
1931 **trajectories of the linguistic traits that shape the observed linguistic relatedness.**
1932 Our approach represents the linguistic relatedness of observed language samples
1933 through the distances among them in a two-dimensional PC space instead of a
1934 phylogenetic tree. This PC space is determined by two optimal axes (PC1 and PC2)
1935 estimated through the PCA algorithm (Figure 2b to Q7). In this PC space, the
1936 language samples with higher relatedness, due to both divergence and contact, would
1937 be distributed closer. In parallel, we reconstruct the past states of linguistic traits for
1938 each language sample using a dynamic model that is derived from the widely-used
1939 covarion model for linguistic trait evolution [3-5] (Figure 2d to Q7). Subsequently, we
1940 also project these past trait states onto the PC space. Accordingly, both past and
1941 current states of linguistic traits for each language sample can be visualized in the PC
1942 space. By computing the differences between the current and past trait states divided
1943 by the reconstruction time for each language sample in the PC space, we can derive a
1944 velocity vector representing the diachronic changes of its linguistic traits (Figure 2e1
1945 to Q7). In other words, this velocity vector illustrates how the linguistic traits in this
1946 language sample evolve into their current states. Accordingly, these velocity vectors
1947 consist of a velocity field in the PC space. And, this velocity field outlines a set of
1948 trajectories that represent the diachronic change of linguistic traits that shape the
1949 observed linguistic relatedness (Figure 2e2 to Q7).

1950      **Transforming the diachronic evolutionary trajectories of the linguistic traits**
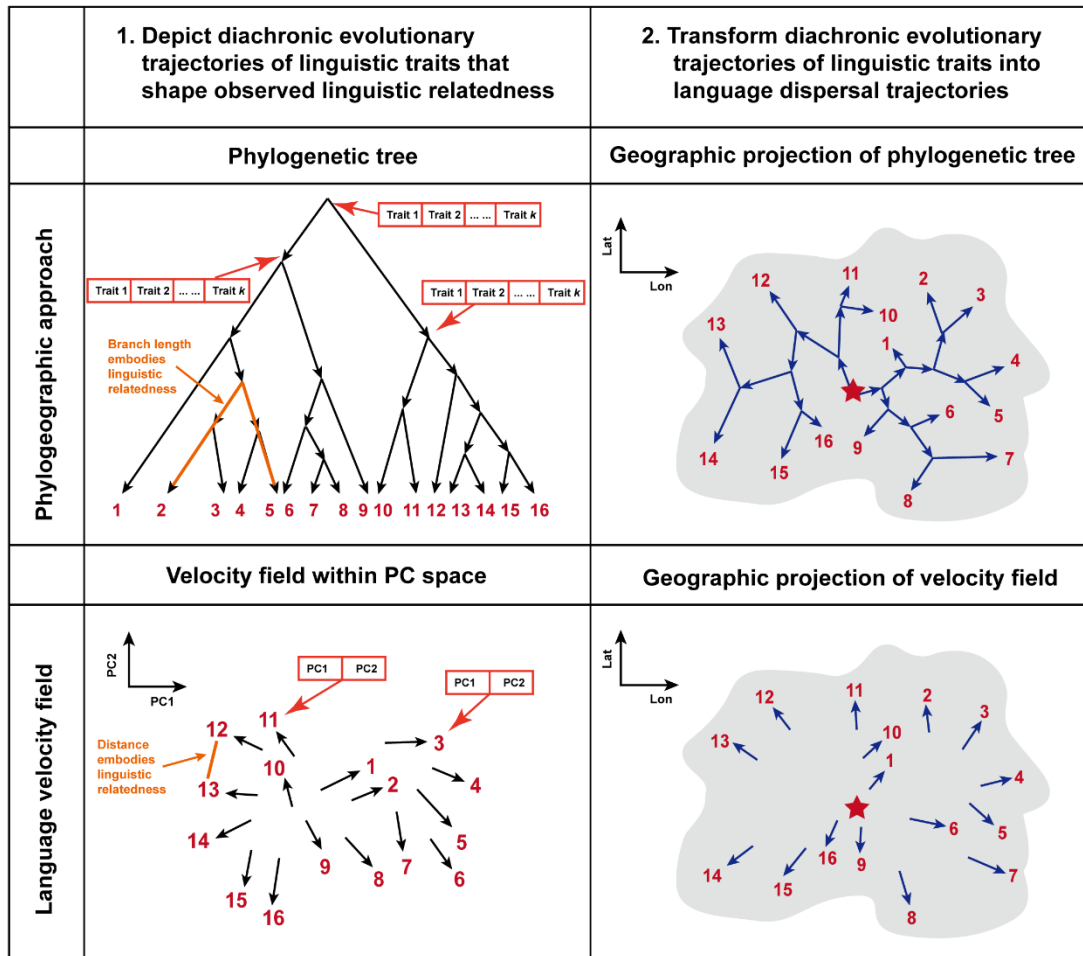1951 **into language dispersal trajectories**. Subsequently, we adopt the kernel projection
1952 proposed by La Manno et al. to map the velocity field from PC space into the
1953 geographic space. This projection seeks the velocity vector in the geographic space
1954 ensuring that its correlation with language geography aligns closely with its
1955 correlation with linguistic relatedness (Figure 2f1 to Q7). This projection is similar to
1956 the projection of each branch within the phylogenetic into the geographic space as a
1957 segment of dispersal trajectories (Figure 1 to Q8). With the kernel projection, the
1958 velocity vectors compose a set of trajectories in geographic space that depict from
1959 where the observed language samples have diffused into their current locations

1960     (Figure 2f2 to Q7).

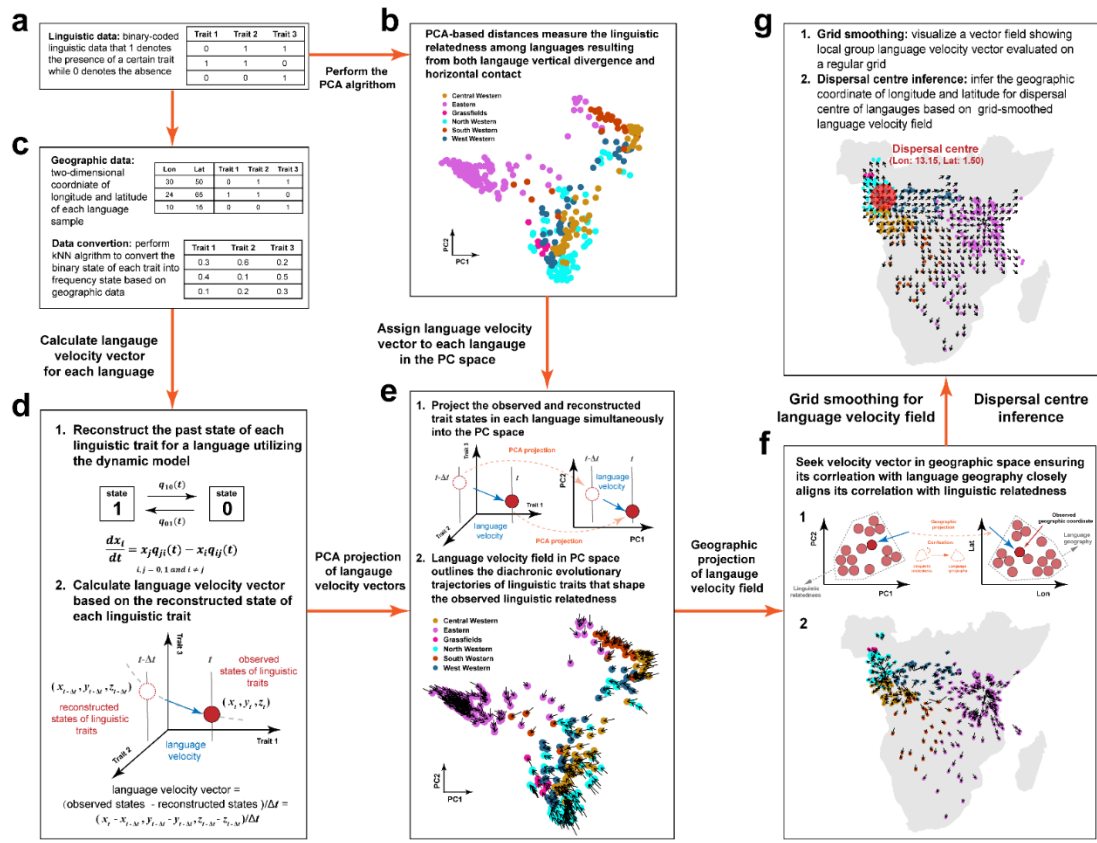1961     **The relationship between the phylogeographic approach and our approach.**
1962 It is noted that if linguistic relatedness can be adequately demonstrated by the
1963 phylogenetic tree, our approach and phylogenetic tree can capture similar linguistic
1964 relatedness. Accordingly, our approach and phylogeographic approach would exhibit
1965 the same performance. In contrast, if linguistic relatedness cannot be adequately
1966 demonstrated by the phylogenetic tree, our approach can capture additional
1967 phylogenetic information from linguistic relatedness due to horizontal contacts as
1968 compared to the phylogeographic approach. Accordingly, our approach may derive a
1969 more reliable result than the phylogeographic approach. In summary, our approach
1970 can be seen as an extension of the phylogeographic approach by relaxing its tree
1971 topology assumption of the phylogeographic approach. This conclusion has been
1972 verified in the revised main text (*Lines 210-303*). Therefore, our approach does not
1973 stand as the opposite of the phylogeographic approach but as its extension.

1974     **Figure**

1975

**Figure 1 to Q8. Language velocity field estimation (LVF) shares the same foundation as the phylogeographic approach but with different implementation strategies.** Both LVF and phylogeographic approach entails two major steps to infer language dispersal pattern. The first is to depict the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. The second is to transform these diachronic evolutionary trajectories of linguistic traits into language dispersal trajectories. In the phylogenetic tree, each language is determined by $k$ linguistic traits. In the velocity field within PC space, each language is determined by PC1 and PC2 which are rearranged from the $k$ linguistic traits through the PCA algorithm. The red number denotes a language. The black arrow signifies the evolutionary direction of linguistic traits in a language. The blue arrow represents the dispersal direction of a language. The red star denotes the estimated dispersal center.

1989

**Figure 2 to Q8. Schematic diagram of language velocity field estimation (LVF) for inferring the dispersal trajectories and centers of languages.** The computational procedures of the LVF comprise two major steps. Subfigures (a) to (e) illustrate the first step which is to estimate a velocity field on the PC space to outline the diachronic evolutionary trajectories of linguistic traits that shape the observed linguistic relatedness. Subfigures (f) to (g) illustrate the second step, which is to project the velocity field from PC space into geographic space. Within the velocity field in geographic space, the directions of the velocity vectors compose a set of continuously changing trajectories that delineate from where these languages diffuse to their current locations. These procedures are exemplified using the Bantu language family. Comprehensive insights into the underlying principles and computational steps can be found in the Materials and Methods section, as well as Supplementary Note 1.

**Reference**

[1] Bouckaert, Remco, et al. "Mapping the origins and expansion of the Indo-European language family." Science 337.6097 (2012): 957-960.

2006     [2] Grollemund, Rebecca, et al. "Bantu expansion shows that habitat alters the route
2007     and pace of human dispersals." Proceedings of the National Academy of Sciences
2008     112.43 (2015): 13296-13301.

2009     [3] Yang, Ziheng. "Maximum-likelihood estimation of phylogeny from DNA
2010     sequences when substitution rates differ over sites." Molecular biology and evolution
2011     10.6 (1993): 1396-1401.

2012     [4] Penny, David, et al. "Mathematical elegance with biochemical realism: the
2013     covarion model of molecular evolution." Journal of Molecular Evolution 53 (2001):
2014     711-723.

2015     [5] Zhang, Menghan, et al. "Phylogenetic evidence for Sino-Tibetan origin in
2016     northern China in the Late Neolithic." Nature 569.7754 (2019): 112-115.

2017

2018 *Q9: Eighth, and final point, the paper is not nice to read, the authors should check*
2019 *their wordings, which are often hard to follow, at times with flaws in grammar, and*
2020 *it would really profit from a complete overhaul and a thorough checking by a proof*
2021 *reader.*

2022 **Replies to Q9:**

2023     We really appreciate the reviewer for pointing this out. In the revised main text,
2024 we have corrected all the typos and grammar flaws. And, we have simplified the long
2025 and wording sentences into the concise and shorten ones. Moreover, we have engaged
2026 the AJE language editing service to thoroughly polish the language of the revised
2027 manuscript (ID: Q2K9ZRSF). We hope that our revised manuscript can be more
2028 readable to native English speakers.

2029

**Replies to Q10:**

We appreciate these comments and are very grateful for the reviewer's encouragement. According to the reviewer's suggestions, we have carefully rewritten the contents about the validations of the approach and the comparison with other approaches. Moreover, we have added a more detailed description of the rationale of our approach. As supplementary, we have also redrawn the schematic diagram to more visually demonstrate the rationale and procedure of our approach. Most importantly, we have restructured the logical flow of our paper, with a focus on sharing a useful and rigorously validated approach with the science community.

Reviewers' Comments:

Reviewer #1:
Remarks to the Author:
The authors answered all my concerns and I do not have further major comments.
Minor changes that need to be addressed:
- new figure 1: typo in panel b, "langauge"
- Panel c: cordiniate, algrithm
- Panel f: unclear sentence + writings in orange are too small and cannot be read


Reviewer #2:
Remarks to the Author:
I think this is perhaps the fourth time I have reviewed this article. As I stated before, I am neither statistician nor linguist, but I detect that the authors have replied to all previous comments by the referees to the maximum extent possible. So I am happy to see the article go to press.

I am impressed by the authors' claims for their efficacy of their "language velocity field" method (at least for the 4 examples they consider), based on PCA rather than phylogenetic "family tree" distances between language subgroups, even if my understanding of all the algebraic formulae that they present is rather limited. The main point for me is that the conclusions of the authors with respect to the homelands of 4 language families that they consider are virtually identical to those I offer in my two recent books The Five-Million-Year Odyssey (Princeton 2022) and First Farmers (second edition, Wiley Blackwell 2023).

So I wish the authors the best of luck with publication and scholarly reception of their views.


Reviewer #3:
Remarks to the Author:
Dear Authors. I have now read all your comments and also had a look at the revised paper and I decided that I should no longer stand in the way, preventing your study to be published. What I would like to ask you, however, is one final thing: For transparency and for replicability, please make sure to make a RELEASE of your code on GitHub and please download this release and submit it to an open independent repository that guarantees long-term archival, such as, for example, Zenodo or Open Science Framework. Here, you will receive a DOI and you should add this DOI to your paper, so we can check the very same code you used to produce the final results that you share with us. Since GitHub itself is owned by Microsoft and Microsoft could shut it down any time they please (think of what happened to Twitter), we need to have the data and the code in public hands. This should not be too hard to do for you, so I hope you'll account for it quickly, and I will recommend the publication of your study, once these changes have been made.

As I will ask for my reports to be published along with my name, I emphasize that the fact that I agree with the publication of this study does not mean that I explicitly express full confidence in its results. It rather means that I feel that it is the best if this study is at this point shared with a larger public that can then discuss then findings in due course and may well find that they have some flaws which were overseen during the review process. I myself am not able to find these flaws by now, nor am I able to assess the quality of the study in full, due to the specifics of my own background. But I am confident that this study provides an interesting contribution to the field and therefore deserves to be published and discussed by more qualified colleagues than myself.

# Response Letter to Reviewers

**Replies to Reviewer 1:**

*Q1: The authors answered all my concerns and I do not have further major comments. Minor changes that need to be addressed: new figure 1: typo in panel b, "langauge"*

**Replies to Q1:**

We sincerely appreciate your careful examination. We have corrected this typo in the revised manuscript.

*Q2: Panel c: cordiniate, algorithm*

**Replies to Q2**

These typos have been corrected in the revised manuscript.

*Q3: Panel f: unclear sentence + writings in orange are too small and cannot be read*

**Replies to Q3**

We sincerely appreciate your comments. We have enlarged the texts in orange to ensure that they can be read clearly by the readers.

**Replies to Reviewer 2:**

*Q1: I think this is perhaps the fourth time I have reviewed this article. As I stated before, I am neither statistician nor linguist, but I detect that the authors have replied to all previous comments by the referees to the maximum extent possible. So I am happy to see the article go to press. I am impressed by the authors' claims for their efficacy of their "language velocity field" method (at least for the 4 examples they consider), based on PCA rather than phylogenetic "family tree" distances between language subgroups, even if my understanding of all the algebraic formulae that they*

*present is rather limited. The main point for me is that the conclusions of the authors with respect to the homelands of 4 language families that they consider are virtually identical to those I offer in my two recent books The Five-Million-Year Odyssey (Princeton 2022) and First Farmers (second edition, Wiley Blackwell 2023). So I wish the authors the best of luck with publication and scholarly reception of their views.*

**Replies to Q1**

We sincerely appreciate your support and affirmation all the time. Moreover, we are also very grateful for your recommendation of your two excellent books to us. We believe that the evidence mentioned within these books can greatly enhance the credibility of our conclusions.

**Replies to Reviewer 3:**

*Q1: Dear Authors. I have now read all your comments and also had a look at the revised paper and I decided that I should no longer stand in the way, preventing your study to be published. What I would like to ask you, however, is one final thing: For transparency and for replicability, please make sure to make a RELEASE of your code on GitHub and please download this release and submit it to an open independent repository that guarantees long-term archival, such as, for example, Zenodo or Open Science Framework. Here, you will receive a DOI and you should add this DOI to your paper, so we can check the very same code you used to produce the final results that you share with us. Since GitHub itself is owned by Microsoft and Microsoft could shut it down any time they please (think of what happened to Twitter), we need to have the data and the code in public hands. This should not be too hard to do for you, so I hope you'll account for it quickly, and I will recommend the publication of your study, once these changes have been made.*

**Replies to Q1**

We are deeply grateful for your support and encouragement. Your valuable suggestions and comments have greatly improved the quality of our manuscript and the transparency and replicability of our approach. Following your suggestions, we have also uploaded our R package and codes to the Zendo (https://doi.org/10.5281/zenodo.10223872).

*Q2: As I will ask for my reports to be published along with my name, I emphasize that the fact that I agree with the publication of this study does not mean that I explicitly express full confidence in its results. It rather means that I feel that it is the best if this study is at this point shared with a larger public that can then discuss then findings in due course and may well find that they have some flaws which were overseen during the review process. I myself am not able to find these flaws by now, nor am I able to assess the quality of the study in full, due to the specifics of my own background. But I am confident that this study provides an interesting contribution to the field and therefore deserves to be published and discussed by more qualified colleagues than myself.*

**Replies to Q2**

We sincerely appreciate your support and encouragement. Moreover, we are very grateful that you are willing to publish your reports with your name. We believe that your reports can provoke new thoughts among the readers.