# ARTICLE

Check for updates

# Expansion by migration and diffusion by contact is a source to the global diversity of linguistic nominal categorization systems

Marc Allassonnière-Tang [1✉], Olof Lundgren[2], Maja Robbers[3], Sandra Cronhamn[2], Filip Larsson[2], One-Soon Her[4], Harald Hammarström[3] & Gerd Carling[2]

Languages of diverse structures and different families tend to share common patterns if they are spoken in geographic proximity. This convergence is often explained by horizontal diffusibility, which is typically ascribed to language contact. In such a scenario, speakers of two or more languages interact and influence each other's languages, and in this interaction, more grammaticalized features tend to be more resistant to diffusion compared to features of more lexical content. An alternative explanation is vertical heritability: languages in proximity often share genealogical descent. Here, we suggest that the geographic distribution of features globally can be explained by two major pathways, which are generally not distinguished within quantitative typological models: feature diffusion and language expansion. The first pathway corresponds to the contact scenario described above, while the second occurs when speakers of genetically related languages migrate. We take the worldwide distribution of nominal classification systems (grammatical gender, noun class, and classifier) as a case study to show that more grammaticalized systems, such as gender, and less grammaticalized systems, such as classifiers, are almost equally widespread, but the former spread more by language expansion historically, whereas the latter spread more by feature diffusion. Our results indicate that quantitative models measuring the areal diffusibility and stability of linguistic features are likely to be affected by language expansion that occurs by historical coincidence. We anticipate that our findings will support studies of language diversity in a more sophisticated way, with relevance to other parts of language, such as phonology.

[1] EA UMR 7206 - MNHN/ CNRS/ Université de Paris, Paris, France. [2] Lund University, Lund, Sweden. [3] Uppsala University, Uppsala, Sweden. [4] Tunghai University, Taichung, Taiwan. ✉email: marc.allassonniere-tang@mnhn.fr

The distribution of linguistic features in the more than 7000 languages of the world (Hammarström, 2016; Hammarström et al., 2019) reflects a scenario where some features may have emerged and spread by horizontal diffusion, whereas others are represented by vertical stability within their lineage. Generally, different feature types vary with respect to their inherent stability (Nichols, 1992; Dediu and Cysouw, 2013), which may reflect their functional role and cognitive preference. In the evolutionary dynamics of language, high stability implies that a feature has high gain and low death rates (attractor feature) whereas low stability implies that a feature has high gain and loss rates (unstable feature), or alternatively low gain and high death rates (recessive feature). Due to their cognitive preference, features of high stability can be both stable in lineage and diffuse by contact, but as a rule, features bound by morphology show a tendency to higher stability in the lineage (Carling and Cathcart, 2021). Both lexicon and grammar vary with respect to their inherent stability (Haspelmath and Tadmor, 2009; Dediu and Cysouw, 2013) but in general, more grammaticalized features of grammar have higher stability rates than more lexical features, and more frequent grammatical and lexical features have higher stability rates than less frequent features (Thomason and Kaufman, 1988; Wilkins, 1996; Matras, 2009). Even though lexical morphemes can be borrowed at varying degrees, grammatical morphemes are very seldom borrowed (Matras and Sakel, 2007). The most frequent lexical items of basic vocabulary have high stability rates and are usually not borrowed (Greenhill et al., 2017), but a majority of the lexicon has lower stability rates and is subject to borrowing at varying degrees (Haspelmath, 2009; Carling et al., 2019). Grammaticality can be viewed as a continuum, ranging from the most grammatical items of grammar (frequent function words of low transparency) to the least grammatical items of the lexicon (cultural and non-frequent content words of high transparency) (Matras and Sakel, 2007). Even though stability is a property that is independent of the grammar-lexicon axis, we expect to find the most grammaticalized items in the domain of high-stable and preferred features, whereas the least grammaticalized items of the lexicon are expected to be less stable in diachrony. While the distribution of every linguistic feature is likely to be shaped by both horizontal diffusion and vertical stability, few analyses based on real data have been proposed to examine how these two pathways simultaneously shape the distribution of specific language features in languages of the world. We aim at filling this gap by providing a case study on nominal classification systems.

Globally, there are three major types of linguistic systems that mirror the cognitive process of categorizing objects within our environment (Lakoff and Johnson, 2013: 162–163; Kemmerer, 2014, 2017a, 2017b). The first type is grammatical gender (Corbett, 1991, 2013), such as the masculine/feminine distinction in French or the masculine/feminine/neuter distinction in German. The second type is noun classes (Corbett, 1991; Grinevald and Seifart, 2004), such as the semantic-based distinction of more than 15 classes in Swahili. The third type is classifiers (Aikhenvald, 2000; Grinevald, 2015), such as the shape-based distinctions in Mandarin (see Supplementary material 1.1 for further details on the definitions). On a grammaticality continuum, gender and noun class markers are thus typical examples of 'grammatical items', while classifiers are relatively closer to 'lexical items', or 'content words'. At a system level, gender is the most grammaticalized system with a few classes of low transparency, triggering agreement in all or a subset of lexemes of a language. Noun class is a less grammaticalized system, involving more classes of higher transparency, but also triggering agreement and targeting all or a subset of lexemes of a language. Classifiers represent the least grammaticalized system, involving a higher number of transparent

markers, targeting selected lexemes of a language and not triggering agreement (see Supplementary material 1.2 for further details). Based on this premise, the existing literature suggests that, on the one hand, classifiers are more easily diffused (horizontally) across language families than gender and noun class (Nichols, 1992: 32, 2003; Wichmann and Holman, 2009: 54–55; Seifart, 2010; Greenhill et al., 2017). On the other hand, in terms of vertical inheritability within languages of the same family, grammatical gender and noun class systems are much more stable than classifiers (Nichols, 2003; Greenhill et al., 2017; Allassonnière-Tang and Dunn, 2020). Studies indicate that grammatical gender hardly ever arises in the course of language contact (Stolz and Levkovych, 2021). However, little quantitative data have been provided to investigate the results of the dynamics of these factors on the distribution of nominal categorization systems worldwide (Seifart, 2010: 730). As an example, classifiers may diffuse faster horizontally. However, their low stability (Kilarski and Allassonnière-Tang, 2021) might counterbalance this fast diffusion, while a slow diffusion of grammatical gender and noun classes might be counterbalanced by their stable inheritability.

We constructed a database of 3077 languages annotated with the presence/absence of gender, noun class, and classifier systems. This database is the first contribution of this paper, as it exceeds by far the existing databases on classifiers and/or gender/noun class in terms of size. As an example, data from the *World Atlas of Language Structures* (Dryer and Haspelmath, 2013) have a sample of 400 languages for classifiers (Gil, 2013) and 257 languages for gender/noun class (Corbett, 2013). The data were compiled by automatic data extraction and checked manually according to precisely defined linguistic criteria for identifying the presence/absence of different nominal classification systems. Data were first extracted from language grammars and grammar sketches using a lightweight keyword-extraction technique (Supplementary material 1.3). Thereupon, manual checking was performed for each individual language and feature, using the Gramfinder tool as an aid for navigating through grammars more effectively (Supplementary material 1.4). Our data (Fig. 1) show that classifiers are more frequent than gender and noun class. Within the data, 26.5% (814/3077) of the languages have classifiers, while 20.1% (634/3077) have gender and 10.3% (317/3077) have noun classes. We can also see that 46.6% (1434/3077) of the languages do not have any of the three systems. In terms of geographic areas, our findings match the existing literature. Classifiers are mainly found in Asia (Gil, 2013) and gender in Europe (Corbett, 2013). Africa has a mixed picture of noun class and gender. Languages that have both gender and classifier systems are mainly found in South America and Papua New Guinea (Seifart, 2010).

However, the ratios we found deviate from the existing literature in several respects. First, the largest available databases report 140 (35.0%, 140/400) classifier languages (Gil, 2013) and 112 (43.5%, 112/257) gender/noun class languages worldwide (Corbett, 2013). Our data show that the actual ratio of classifier and gender/noun class is much lower. Moreover, in our data, the number of classifier languages and that of gender languages are fairly close, and both are much higher than noun class languages. This is intriguing since noun class and gender systems are often considered to be the same type of system (Corbett, 2013) and under such a premise, one would expect that the distribution should be similar. To further investigate the diffusibility and heritability of each of the three systems, we used Delaunay neighbors and phylogenetic neighbors to assess the areal and phylogenetic cohesion of these systems (Supplementary material 2.1). This method mirrors the structure of existing methods for the assessment of diffusibility and heritability of linguistic features (Parkvall, 2008; Wichmann and Holman, 2009; Dediu and Cysouw, 2013; Cathcart et al., 2018; Murawaki and Yamauchi,
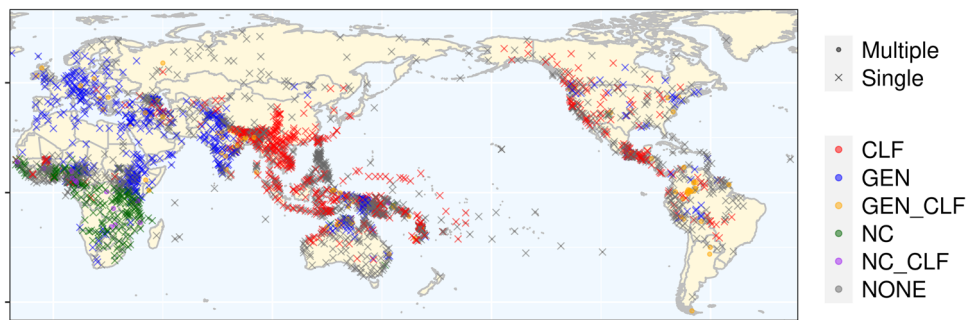
**Fig. 1 Nominal classification systems in languages of the world.** The abbreviations are interpreted as follows. CLF classifier, GEN gender, NC Noun Class. Note that Gender and Noun Class are mutually exclusive according to our coding policy (see Supplementary material 1).

2018; Nikolaev, 2019). If languages found within the same geographic area tend to share the same feature, the feature has a strong geographic cohesion and it implies that the feature diffuses geographically. If languages found within the same phylogenetic branch share the same feature, it has strong phylogenetic cohesion, which indicates that the feature is robustly inherited within the language family. The results (Fig. 2a) from Wilcoxon rank-sum tests with continuity correction show that the geographic cohesion of classifiers ($\mu = 0.5$, $m = 0.5$) is significantly smaller than the geographic cohesion of gender ($\mu = 0.6$, $m = 0.7$, $w = 207207$, $p < 0.001$) and noun classes ($\mu = 0.7$, $m = 0.8$, $w = 79545$, $p < 0.001$), while the geographic cohesion of gender is significantly smaller than the geographic cohesion of noun classes ($w = 82972$, $p < 0.001$). Furthermore, the phylogenetic cohesion of classifiers ($\mu = 0.6$, $m = 0.7$) is significantly smaller than the phylogenetic cohesion of gender ($\mu = 0.8$, $m = 1$, $w = 195,193$, $p < 0.001$) and noun classes ($\mu = 0.9$, $m = 1$, $w = 79,472$, $p < 0.001$), while the phylogenetic cohesion of gender is significantly smaller than the phylogenetic cohesion of noun classes ($w = 84,686$, $p < 0.001$).

The results support our theory in showing that gender and noun class have a stronger phylogenetic cohesion and thus a stronger heritability than classifiers. However, our results also demonstrate that the geographic cohesion of gender and noun class is higher than that of classifiers. Furthermore, the geographic and phylogenetic cohesion of noun class is higher than that of gender. This suggests that factors other than horizontal diffusibility and vertical stability might have either inflated the distribution of gender languages or diminished the distribution of noun class languages. Such factors can be family-specific since the measure of geographic cohesion does not control language families. Universal preferences could potentially have an effect too, as universal preferences converge more easily since they are preferred by the human processor. Research indicates that more grammaticalized systems are cognitively preferred and more learnable, because of the higher cognitive load of distinguishing more categories (Bentz and Winter, 2013). Due to limitations of data and method, our study does not account for these possible explanations.

We hypothesize that our observations can be explained by the difference between two feature dispersal mechanisms, which we label *feature diffusion* and *language expansion*. Feature diffusion implies that features spread by areal diffusion, which does not require relocation or diffusion of languages or speakers. Language expansion implies a spread of features, which depends on either the long-distance movement of language communities or the small-scale movement of foreign speakers into a local community, which may lead to a subsequent language shift or death of local languages (Neureiter et al., 2021). While all linguistic features are spread by both mechanisms, the weight of language expansion
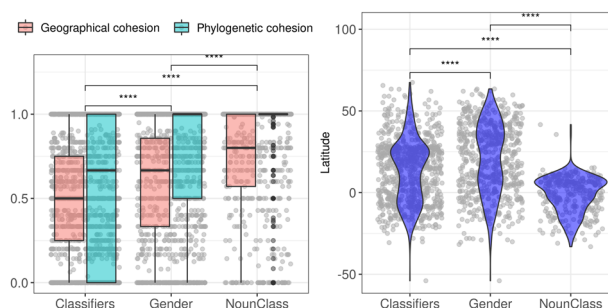


**Fig. 2 The distribution of different measures of classifier, gender, and noun class systems. a** Shows the distribution and average of geographic and phylogenetic cohesion of classifiers, gender, and noun class. Gender and noun class are stronger than classifiers in terms of geographic cohesion and in particular in phylogenetic cohesion. **b** Shows the distribution of classifiers, gender, and noun class with respect to latitude. Classifiers and gender are more evenly distributed than noun class, which is concentrated in the southern hemisphere.

and feature diffusion may vary. Besides other influencing factors such as high birth rates and low death rates, a high areal cohesion could depend on both feature diffusion and language expansion. The same may be true for a high phylogenetic cohesion (features may potentially diffuse within families). In our study, we hypothesize that the more grammaticalized gender and noun class systems spread historically mainly through language expansion, giving gender and noun class a larger magnitude of language expansion. Contrary to this, we hypothesize that the more lexical classifier systems relied more on feature diffusion.

As an additional exploration toward this hypothesis, we consider the 'Continental Axis Theory' (Diamond, 1997), which suggests that humans tend to migrate east and west rather than north and south to stay within similar climatic conditions for the sake of farming (Greenhill, 2014; Güldemann and Hammarström, 2020). This suggests that language expansion is likely to be stronger along latitudes. When visualizing the distribution of latitude across languages with different features in our data (Supplementary material 2.2, Fig. 2b), we see that noun classes are concentrated in a specific range of latitude, which likely represents Africa, while both classifier and gender languages have a larger range of latitude. This shows that the spreading dynamics of gender are in fact much more similar to classifiers than noun class. Again, this large range of latitude could have resulted from feature diffusion and/or language expansion for both features. To further verify which of the two scenarios is more likely, we consider the diversity of language families across languages with classifier, gender, and noun class systems.

Language contact is largely determined by complex and sparsely documented social factors (Hickey, 2010; Bowern et al., 2011). Nevertheless, if the distribution of a feature is mostly influenced by feature diffusion, the feature is more likely to be found across languages from different language families located in geographic proximity. This can be explained by the fact that the diffusion would happen by language contact and diffuse from a language to its geographic neighbors, with few restrictions of family affiliation (Coupé et al., 2013). In our case study, classifiers are expected to diffuse more than gender and noun class. We thus expect that classifiers are more likely to be found across different language families in the same area. As for gender and noun class, if they expand more by language expansion, we expect that a gender or noun class language is less likely to have geographic neighbors from different language families since languages from different families are more likely to have been pushed away and/ or replaced by the family with gender. To investigate this hypothesis, we divide the world map into 3267 grids (Derungs et al., 2018). For each grid and for each feature, we count the number of language families represented by languages within the grid with the feature in question (Supplementary material 2.3). Our data show that the family density of classifier languages is indeed higher than the family density of gender languages ($w = 294{,}410$, $p < 0.001$) and noun class languages ($w = 171{,}006$, $p < 0.001$). The data also show that the family density of gender languages is significantly higher than the family density of noun class languages ($w = 117{,}264$, $p < 0.001$).

Additional evidence consistent with our hypothesis relates to the geographic coverage of language families (Supplementary material 2.4). If we count the pairwise distance between all languages of each family and compare the mean distance across all possible pairs (normalized on a scale from 0 to 1), we find that the Indo-European language family has the largest geographic coverage, almost twice as large as the second-ranked Austronesian family (the ratio is 1–0.62). Accordingly, the top 10 families with a large geographic coverage are Indo-European (1), Austronesian (0.62), Eskimo-Aleut (0.51), Afro-Asiatic (0.48), Turkic (0.46), Atlantic-Congo (0.41), Tungusic (0.40), Mongolic-Khitan (0.34), Athabaskan-Eyak-Tlingit (0.34), and Pama-Nyungan (0.33). Several of these families did experience large-scale migration of their speakers and are also well-known for their gender and/or noun class systems (e.g., Indo-European, Afro-Asiatic, Atlantic-Congo). This is also consistent with our hypothesis, as the expansion of these families could have indirectly contributed to the expansion of their gender and/or noun class systems. As an example, the expansion of the Indo-European family (Mallory and Adams, 2006; Anthony, 2007; Carling, 2019) is likely to be one of the main factors that contributed to the expansion of grammatical gender in Eurasia, while the expansion of the Niger-Congo languages (and in particular Bantu) is likely to have played a major role in the expansion of noun classes in Africa (Hepburn-Gray, 2020).

Finally, language speakers tend to stay in similar environments when they migrate (Nichols, 1992; Gray and Jordan, 2000; Ramat, 2012; Hock and Jospeh, 2019); therefore, if gender and noun class languages spread more by language expansion, we expect to find less variance within the natural environment surrounding their location. If classifiers spread more by feature diffusion, the spread is expected to be more independent of environmental factors. We thus assume that features spreading by language expansion should have a smaller variance of the environmental factors that facilitate migration and farming (Antunes et al., 2020). We investigate environmental factors that are less likely to vary across geographical areas (Moore et al., 2002; Pacheco Coelho et al., 2019; Antunes et al., 2020). As an example, the mean temperature varies drastically across geographical areas, which is likely to

affect its variance. We thus select these three environmental factors: low variation of elevation, distance to water bodies, and rainfall. Low variation of elevation is generally more suitable for farming, as topographically complex areas largely correspond to versatile ecosystems and may pose restrictions on settlement options (Hassan, 1975; Kavanagh et al., 2018), while accessibility to an adequate source of water (either by river or rainfall) is also one of the basic conditions considered when expanding and finding new settlements.

River access seems to influence distributions of hunter-gatherer languages (Derungs et al., 2018), which can partly be explained by high supplies of protein-rich foods provided in these regions (Hassan, 1975). Rivers may act as physical barriers that keep linguistic groups apart, as well as resource providers, and most importantly as means of transportation. River density is described as giving rise to high language diversity with a good fit for Africa (Axelsen and Manrubia, 2014). Riverine transportation options might also increase contact among groups which, depending on the navigability of the riverine network, can decrease language diversity for example in North America (Pacheco Coelho et al., 2019). Feature diffusion along river networks can be assessed areally by route inference (Ranacher et al., 2017). Precipitation variables can correlate significantly with linguistic diversity (Moore et al., 2002; Collard and Foley, 2002; Hua et al., 2019). Continuous precipitation provides reliable bases for production and thereby decreases ecological risk (Nettle, 1996) and favors small-scale linguistic groups (Amano et al., 2014). More precisely, precipitation seasonality seems more influential than total annual or monthly precipitation or production factors such as mean growing season (Hua et al., 2019).

The exact time of the spreads is unknown for most of our families, but given their reconstructed age range of 4000–8000 years (Nichols, 2008; Greenhill et al., 2017), the spreads most likely occurred during the mid-Holocene period. Therefore, we use mid-Holocene projections for the selected environmental factors (Derungs et al., 2018), which are shown in Fig. 3. The results are consistent with our hypothesis, as for each of the three environmental factors, classifier languages have the largest variance, followed by gender languages, and then by noun class languages (as shown by Quantile dispersion, Levene tests, and Conover tests. Further details are included in Supplementary material 2.5).

As a summary, we show that the geographical distribution of nominal classification systems is likely to have been influenced by the mechanisms of language expansion. Evidence from language



**Fig. 3 The values of distance to river, precipitation in the wettest quarter, and the standard deviation of elevation for languages of each nominal classification system, in which classifier languages have the largest variance.** The values of the three environmental factors have been normalized to a scale from 0 to 1.

family density, the geographical coverage of language families, and variance of environmental factors highlight the importance of distinguishing between the two mechanisms of language expansion and feature diffusion. These two mechanisms are generally not distinguished in quantitative assessments of the horizontal stability of linguistic features; however, they can lead to similar results, while in fact telling a drastically different story about the importance of grammaticalization in the diffusibility and the stability of the analyzed linguistic features. Our study also points out the importance of testing such an assumption for other linguistic features and other factors. We demonstrate how the effects of language expansion could be investigated in a case study of nominal classification systems. We encourage future studies to replicate the analysis on other features related to phonology, syntax, semantics, among others, to compare their dynamics in spreading. We also encourage the building of evolutionary models to take into account the impact of non-linguistic factors, such as language expansion, along with linguistic factors, such as grammaticalization, so that the spreading dynamics of linguistic features are modeled in a more accurate way.

## Data availability

Supplementary Information is available for this paper at the journal website and at the repository https://github.com/marctang/Diversity_NominalCategorization.

## References

Aikhenvald AY (2000) Classifiers: a typology of noun categorization devices. Oxford University Press, Oxford

Allassonnière-Tang M, Dunn M (2020) The evolutionary trends of grammatical gender in Indo-Aryan languages. Lang Dyn Change 11(2):211–240

Amano T, Sandel B, Eager H, Bulteau E, Svenning J-C, Dalsgaard B, Rahbek C, Davies RG, Sutherland WJ (2014) Global distribution and drivers of language extinction risk. Proc R Soc B: Biol Sci 281:20141574

Anthony DW (2007) The horse, the wheel, and language: how Bronze-Age riders from the Eurasian steppes shaped the modern world. Princeton University Press, Princeton

Antunes N, Schiefenhövel W, d'Errico F, Banks WE, Vanhaeren M (2020) Quantitative methods demonstrate that environment alone is an insufficient predictor of present-day language distributions in New Guinea. PLoS ONE 15:e0239359

Axelsen JB, Manrubia S (2014) River density and landscape roughness are universal determinants of linguistic diversity. Proc R Soc Lond B: Biol Sci 281:1–9

Bentz C, Winter B (2013) Languages with more second language learners tend to lose nominal case. Lang Dyn Change 3:1–27

Bowern C et al. (2011) Does lateral transmission obscure inheritance in hunter-gatherer languages? PLoS ONE 6:e25195

Carling G (eds.) (2019) The Mouton Atlas of languages and cultures: vol. 1. Mouton De Gruyter, Berlin

Carling G, Cathcart C (2021) Reconstructing the evolution of Indo-European grammar. Language 97:561–598.

Carling G, van de Weijer J, Cronhamn S, Johansson N, Farren R (2019) The Cultural Lexicon of Indo-European in Europe: quantifying stability and change. In: Kroonen G, Mallory JP, Comrie B (eds.) Talking Neolithic: Proceedings of the workshop on Indo-European origins. Max Planck Institute for Evolutionary Anthropology, Leipzig, pp. 39–68

Cathcart C, Carling G, Larsson F, Johansson N, Round E (2018) Areal pressure in grammatical evolution: an Indo-European case study. DIA 35:1–34

Collard IF, Foley RA (2002) Latitudinal patterns and environmental determinants of recent human cultural diversity: do humans follow biogeographical rules? Evol Ecol Res 4:371–383

Corbett GG (1991) Gender. Cambridge University Press, Cambridge

Corbett GG (2013) Number of genders. In: Dryer MS, Haspelmath M (eds.) The World Atlas of language structures online. Max Planck Institute for Evolutionary Anthropology, Leipzig

Coupé C, Hombert JM, Marsico E, Pellegrino F (2013) Investigations into determinants of the diversity of the world's languages. In: Peng G, Shi F (eds.) Eastward flows the great river: Festschrift in honor of Professor William S-Y. Wang on his 80th birthday. City University of Hong Kong Press, Hong Kong, pp. 75–108

Dediu D, Cysouw M (2013) Some structural aspects of language are more stable than others: a comparison of seven methods. PLoS ONE 8:1–10

Derungs C, Köhl M, Weibel R, Bickel B (2018) Environmental factors drive language density more in food-producing than in hunter–gatherer populations. Proc R Soc B: Biol Sci. https://doi.org/10.1098/rspb.2017.2851

Diamond J (1997) Guns, germs, and steel. W W Norton, New York

Dryer MS, Haspelmath M (eds.) (2013) The world atlas of language structures online. Max Planck Institute for Evolutionary Anthropology, Leipzig

Gil D (2013) Numeral classifiers. In: Dryer MS, Haspelmath M (eds.) The World Atlas of language structures online. Max Planck Institute for Evolutionary Anthropology, Leipzig

Hammarström, H, Forkel, R, Haspelmath, M. (2019) Glottolog 4.1. Jena: Max Planck Institute for the Science of Human History. https://glottolog.org/

Gray RD, Jordan FM (2000) Language trees support the express-train sequence of Austronesian expansion. Nature 405:1052–1055

Greenhill SJ (2014) Demographic correlates of language diversity. In: Evans B, Bowern C (eds.) Routledge handbook of historical linguistics. Routledge, New York, pp. 557–578

Greenhill SJ, Wu C-H, Hua X, Dunn M, Levinson SC, Gray R (2017) Evolutionary dynamics of language systems. Proc Natl Acad Sci USA 114:E8822–E8829

Grinevald C (2015) Linguistics of classifiers. In: Wright JD (ed.) International encyclopedia of the social & behavioral sciences. Elsevier, Amsterdam, pp. 811–818

Grinevald C, Seifart F (2004) Noun classes in African and Amazonian languages: towards a comparison. Linguistic Typol 8:243–285

Güldemann T, Hammarström H (2020) Geographical axis effects in large-scale linguistic distributions. In: Crevels M, Muysken P (eds.) Language dispersal, diversification, and contact. Oxford University Press, Oxford, pp. 58–77

Hammarström H (2016) Linguistic diversity and language evolution. J Lang Evol 1:19–29

Haspelmath M (2009) Lexical borrowing: concepts and issues. In: Haspelmath M, Tadmore U (eds.) Loanwords in the world's languages: a comparative handbook. Mouton De Gruyter, Berlin, pp. 35–54

Haspelmath M, Tadmor U (2009) Loanwords in the world's languages: a comparative handbook. Mouton De Gruyter, Berlin

Hassan FA (1975) Determination of the size, density, and growth rate of hunting-gathering populations. Mouton De Gruyter, Berlin

Hepburn-Gray R (2020) Niger-Congo noun classes: reconstruction, historical implications, and morphosyntactic theory. Dissertation, State University of New York at Buffalo.

Hickey R (2010) Language contact: reconsideration and reassessment. In: Hickey R (ed.) The handbook of language contact. Wiley-Blackwell, Oxford, pp. 1–28

Hock HH, Joseph BD (2019) Language history, language change, and language relationship: an introduction to historical and comparative linguistics. Mouton De Gruyter, Berlin

Hua X, Greenhill SJ, Cardillo M, Schneemann H, Bromham L (2019) The ecological drivers of variation in global language diversity. Nat Commun 10:2047

Kavanagh PH, Vilela B, Haynie HJ, Tuff T, Lima-Ribeiro M, Gray RD, Botero CA, Gavin MC (2018) Hindcasting global population densities reveals forces enabling the origin of agriculture. Nat Hum Behav 2:478–484

Kilarski M, Allassonnière-Tang M (2021) Classifiers in morphology. In: Aronoff M (ed.) Oxford research encyclopedia of linguistics. Oxford University Press, Oxford, pp. 1–28

Kemmerer D (2014) Word classes in the brain: implications of linguistic typology for cognitive neuroscience. Cortex 58:27–51

Kemmerer D (2017a) Some issues involving the relevance of nominal classification systems to cognitive neuroscience: response to commentators. Lang Cogn Neurosci 32:447–456

Kemmerer D (2017b) Categories of object concepts across languages and brains: the relevance of nominal classification systems to cognitive neuroscience. Lang Cogn Neurosci 32:401–424

Lakoff G, Johnson M (2013) Metaphors we live by. University of Chicago Press, Berkeley

Mallory JP, Adams DQ (2006) The Oxford introduction to Proto-Indo-European and the Proto-Indo-European world. Oxford University Press, Oxford

Matras Y (2009) Language contact. Cambridge University Press, Cambridge

Matras Y, Sakel J (eds.) (2007) Grammatical borrowing in cross-linguistic perspective. Mouton De Gruyter, Berlin

Moore JL, Manne L, Brooks T, Neil DB, Davies R, Rahbek C, Williams P, Balmford A (2002) The distribution of cultural and biological diversity in Africa. Proc R Soc Lond Ser B: Biol Sci 269:1645–1653

Murawaki Y, Yamauchi K (2018) A statistical model for the joint inference of vertical stability and horizontal diffusibility of typological features. J Lang Evol 3:13–25

Nettle D (1996) Language diversity in West Africa: an ecological approach. J Anthropol Archaeol 15:403–438

Neureiter N, Ranacher P, van Gijn R, Bickel B, Weibel R (2021) Can Bayesian phylogeography reconstruct migrations and expansions in linguistic evolution? R Soc Open Sci 8:201079

Nichols J (1992) Linguistic diversity in space and time. University of Chicago Press, Chicago

Nichols J (2003) Diversity and stability in language. In: Joseph BD, Janda RD (eds.) The handbook of historical linguistics. Blackwell, Oxford, pp. 283–311

Nichols J (2008) Language spread rates as indicators of glacial-age peopling of the Americas. Curr Anthropol 49:1109–1117

Nikolaev D (2019) Areal dependency of consonant inventories. Lang Dyn Change 9:104–126

Parkvall M (2008) Which parts of language are the most stable? STUF—Lang Typol Univers 61:234–250

Pacheco Coelho MT et al. (2019) Drivers of geographical patterns of North American language diversity. Proc R Soc B 286:20190242

Ramat P (2012) The impact of migrations on the linguistic landscape of Europe. In: Kortmann B, van der Auwera J (eds.) The languages and linguistics of Europe. Mouton De Gruyter, Berlin, pp. 683–695

Ranacher P, Van Gijn R, Derungs C (2017) Identifying probable pathways of language diffusion in South America. https://doi.org/10.5167/UZH-137656

Seifart F (2010) Nominal classification. Language Linguist Compass 4:719–736

Stolz T, Levkovych N (2021) On the (almost im)possible emergence of grammatical gender in language-contact situations. In: Leykovych N (ed.) Susceptibility vs resistance in language-contact situations. Mouton De Gruyter, Berlin, p 3–40

Thomason SG, Kaufman T (1988) Language contact, creolization and genetic linguistics. University of California Press, Berkeley

Wichmann S, Holman EW (2009) Temporal stability of linguistic typological features. Lincom, München

Wilkins DP (1996) Natural tendencies of semantic change and the search for cognates. In: Durie M, Ross M (eds.) The comparative method reviewed: regularity and irregularity in language change. Oxford University Press, Oxford, pp. 264–304

## Author contributions

M.A.-T., H.H. and G.C. conceived and designed the study. M.A.-T., O.L., F.L., O.-S.H., H.H. and G.C. contributed to establish the data. M.A.-T. conducted the quantitative analyses. All authors were involved in discussion and interpretation of the results. All authors contributed to the writing with M.A.-T., H.H. and G.C. having leading roles; All authors contributed to the Supplementary Information.

## Competing interests

The authors declare no competing interests.

## Ethical approval

Not applicable as this study did not involve human participants.

## Informed consent

Not applicable as this study did not involve human participants.

## Additional information

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1057/s41599-021-01003-5.

**Correspondence** and requests for materials should be addressed to Marc Allassonnière-Tang.

**Reprints and permission information** is available at http://www.nature.com/reprints

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.