




Article

# More Rule than Exception: Parallel Evidence of Ancient Migrations in Grammars and Genomes of Finno-Ugric Speakers

Patrícia Santos <sup>1,2</sup>, Gloria González-Fortes <sup>2</sup>, Emiliano Trucchi <sup>3</sup>, Andrea Ceolin <sup>4</sup> , Guido Cordoni <sup>5</sup>, Cristina Guardiano <sup>4</sup> , Giuseppe Longobardi <sup>6</sup> and Guido Barbujani <sup>2,\*</sup> 

<sup>1</sup> CNRS, UMR 5199—PACEA, Université de Bordeaux, Bâtiment B8, Allée Geoffroy Saint Hilaire, 33615 Pessac, France; sntprc1@unife.it

<sup>2</sup> Dipartimento di Scienze della Vita e Biotecnologie, Università di Ferrara, 44121 Ferrara, Italy; gloria.gonzalezfortes@unife.it

<sup>3</sup> Department of Life and Environmental Sciences, Marche Polytechnic University, 60131 Ancona, Italy; e.trucchi@univpm.it

<sup>4</sup> Dipartimento di Comunicazione ed Economia, Università di Modena e Reggio Emilia, 42121 Reggio Emilia, Italy; ceolin@unimore.it (A.C.); cristina.guardiano@unimore.it (C.G.)

<sup>5</sup> School of Veterinary Medicine, University of Surrey, Guildford GU2 7AL, UK; g.cordoni@surrey.ac.uk

<sup>6</sup> Department of Language and Linguistic Science, University of York, York YO10 5DD, UK; giuseppe.longobardi@york.ac.uk

\* Correspondence: g.barbujani@unife.it

Received: 30 October 2020; Accepted: 9 December 2020; Published: 11 December 2020



**Abstract:** To reconstruct aspects of human demographic history, linguistics and genetics complement each other, reciprocally suggesting testable hypotheses on population relationships and interactions. Relying on a linguistic comparative method based on syntactic data, here we focus on the non-straightforward relation of genes and languages among Finno-Ugric (FU) speakers, in comparison to their Indo-European (IE) and Altaic (AL) neighbors. Syntactic analysis, in agreement with the indications of more traditional linguistic levels, supports at least three distinct clusters, corresponding to these three Eurasian families; yet, the outliers of the FU group show linguistic convergence with their geographical neighbors. By analyzing genome-wide data in both ancient and contemporary populations, we uncovered remarkably matching patterns, with north-western FU speakers linguistically and genetically closer in parallel degrees to their IE-speaking neighbors, and eastern FU speakers to AL speakers. Therefore, our analysis indicates that plausible cross-family linguistic interference effects were accompanied, and possibly caused, by recognizable demographic processes. In particular, based on the comparison of modern and ancient genomes, our study identified the Pontic-Caspian steppes as the possible origin of the demographic processes that led to the expansion of FU languages into Europe.

**Keywords:** genomes; syntax; genetic and linguistic distances; human migrations; phylogenies

## 1. Introduction

Darwin proposed that linguistic diversity along human history tends to be correlated with the biological differentiation of populations [1]. Indeed, factors isolating populations from each other, such as geographical distance and barriers to migration, are likely to promote both biological and cultural divergence, whereas factors favoring contacts have the opposite effect at both levels [2–5]. In fact, despite elite dominance and other processes of horizontal language transmission creating local mismatches [6], parallel genetic and linguistic changes have often appeared as the rule rather than the

exception [2,4,7–10]. This implies that linguistic diversity may offer a set of testable hypotheses about the demographic processes shaping genetic diversity, and vice versa.

In this study, through a multidisciplinary approach comparing grammars and genomes, we contribute to a better understanding of population diversity, both cultural and biological, in western/Central Eurasia. We focus on Altaic- (hereafter: AL) Finno-Ugric- (FU) and Indo-European- (IE) speaking populations, with a special emphasis on FU speakers. The reason is that FU has appeared as a possible exception to the general gene-language correspondence. Indeed, its monophyletic unity was acknowledged linguistically already in the 18th century [11] and (unlike AL, a linguistically controversial unit: see [12] for a summary) remains virtually undisputed (with the possible caveats in [13]), but FU-speaking populations fail to be identified as a genetic group [14]. In particular, the westernmost FU-speaking populations in northern and especially Central Europe have been shown to display a peculiar exception to the conclusion that, in Europe, grammars are better predictors of genetic distances than geography [4]. This exception is worth some further investigation.

Until recently, comparative studies of genes and languages suffered from serious limitations, simply because of the data available. On the one hand, only seldom were whole genomes considered in these comparisons. On the other, classical etymological comparison of vocabulary items, still normally used to reconstruct phylogenetic history even in modern quantitative studies (see, e.g., [15–17]), work well within a language family, but words cannot be used for broader comparisons: for, by definition, across different language families there are no recognizable common etymologies (i.e., lexical cognates; see Ref. [18] for an important attempt to remedy some of these problems). However, the Parametric Comparison Method (PCM) [19–24]), which explores the phylogenetic information contained in the generative rules of syntax, has in principle overcome the limitations of vocabulary-based taxonomic methods: through parameters, i.e., abstract and universally definable syntactic polymorphisms, the PCM quantifies language differences/similarities across languages and even language families into a synthetic measure.

In this work, we want to find out if the apparent lack of gene-language correlation among FU-speaking populations may hide some deeper congruence determined at possibly different stages by the relationships of FU with its IE and AL neighbors and by their migration history. To do so, we will estimate the degree of similarity between languages from the three families using the PCM.

The amount of secondary contact between individual pairs of languages can in principle be measured from lexical comparison. Indeed, lexical borrowings between Indo-European and Uralic [25–27] have been used to show historical contact between the two families. In principle, by counting for each pair of languages how many loanwords have been exchanged in either direction and reducing them to a single figure, one could estimate the amount of borrowing and use it as a distance measure to be compared with genomic data. However, presently, no general database summing up this kind of information in an exhaustive and uniform fashion exists, nor an algorithm for non-arbitrarily performing such calculations.

In this respect, complementing traditional lexical insights with comparison of syntactic parameters presents some particular advantages: a dataset for many relevant languages is already available [24], its historical information has been evaluated, and a syntactic distance is readily computable for every pairwise relationship within the whole set of languages used; such a distance sums up both potential vertical and horizontal signals into a single figure, paralleling well the genetic distances among the corresponding populations, which can be influenced by both splits and successive admixture.

In the course of the analysis, some individual taxonomic insights from classical linguistic levels (Ref. [28,29] and the other references mentioned above), will also be resorted to and compared with the syntactic signal, whenever informative.

The conclusions we reach about the geographical and demographical history of FU will have some indirect consequences on the current debate about the prehistory of IE speakers. As for the latter, let us recall that, despite a long tradition of studies, it is still debated whether early IE languages came into Europe from the Pontic-Caspian steppes (and spread West in the Bronze Age [30,31]) or from Anatolia

(and spread with the dispersal of early Neolithic farmers [16,32,33]). Thus, we compared the syntax of several AL- FU- and IE-speaking populations with the available genome-wide data, both contemporary and ancient, in the area of interest. Of particular relevance was one Bronze-Age population from the Pontic-Caspian steppe, the Yamnaya, the likely source of the Bronze-Age migration leading to a Westwards diffusion of DNA of Central Asian origin and, according to some authors, of IE languages in Europe [34–37]. By contrast, a recent analysis of Asian genomes suggested that the spread of IE languages in South Asia and Anatolia may have little, if anything, to do, with migration from the Pontic-Caspian steppes [38]. An analogous uncertainty surrounds the homeland of early FU speakers, whether in the river Volga basin or further East, in Siberia [39].

Our multidisciplinary approach comparing grammars and genomes will ultimately help us better frame the evolution of this cultural and biological diversity in western/Central Eurasia, reinforcing the idea of widespread congruence between the two types of variables.

## 2. Materials and Methods

### 2.1. Genomic Dataset

The dataset analyzed in this study comprises the high-coverage sequenced genomes of 45 individuals from 17 populations from Eurasia (Supplementary Table S1). The samples were collected from Pagani et al. (2016) [40] and downloaded from the public database ENA (European Nucleotide Archive). For the sake of equal representation, a random subset of three individuals per population was chosen for populations with a larger sample size, to perform all the analyses.

Ancient and modern Genome-wide SNP array data from Ref. [41] were used to estimate Outgroup  $f_3$ -statistic and  $qpAdm$  analysis (Supplementary Tables S2 and S3, respectively).

### 2.2. Dataset Preparation

Samples from Ref. [40] were in Complete Genomics MasterVar format files (reads mapped against the human genome reference hg19/GRCh37). The MasterVar file was converted into a Variant Call Format (VCF) by the `cgatool mkvf` (version 1.8.0.1) from Complete Genomics. The VCF file created only contains SNP variants with a high confidence ( $>40$  dB). All the VCF files from the different individuals were merged using `BCFtools` (version v1.6-36) merged with the option “-m none” to output the multiallelic sites in different lines. All duplicated variants were excluded from the data. The VCF files were phased using `SHAPEIT2` (version v2.r837) using the 1000 Genomes phase 3 haplotypes as a reference panel, as recommended. Heterozygous sites not present in the 1000 Genomes data were left unphased. In the end, genotypes were obtained for 11,931,455 autosomal SNPs.

### 2.3. Principal Component Analysis

A general description of genetic variation was obtained by Principal Component Analysis (PCA). QTLtools [42] (version v1.1) was used on scaled and centered genotype data on relatively independent (50 Kb distance) and non-rare variants (minor allele frequency = 0.05).

### 2.4. Genomic Distances

Weir and Cockerham's genomic distances between populations were calculated by the 4P software [43] (version 1.0). Genomic regions that may be under selection were masked using `bedtools subtract` (version v2.26) and variants with a missing call rate exceeding 10% were excluded, resulting in a total of 9,881,752 autosomal SNPs.

### 2.5. Linguistic Dataset

For the analysis of linguistic data by PCM, we used the 94 binary parameters and their settings defining properties of nominal structures for 69 modern Eurasian languages recently employed in [24] and accessible from <https://github.com/AndreaCeolin/FormalSyntax>. Theoretical background and

more technical details about the structure of the parameter system and the settings of the individual values are found in [44–46].

The original dataset of 69 languages has been reduced to a subset of 34 IE, FU and AL languages, to improve resolution on the 17 populations for which genetic data are available and their neighbors.

Data were available for the main FU subfamilies (Ugric, Volgaic, Permic, Balto-Finnic), with the exception of Lapp (Saami). For three languages, Mari (Cheremiss), Udmurt (Votyak) and Khanty (Ostyak), we encoded two diastatic variants; the two Even (Tungusic, AL) varieties are instead diatopic. For details, see Ref. [24]. The relevant IE languages belong to three subfamilies, namely, Indo-Iranian, Germanic and Slavic (see Supplementary Figure S1).

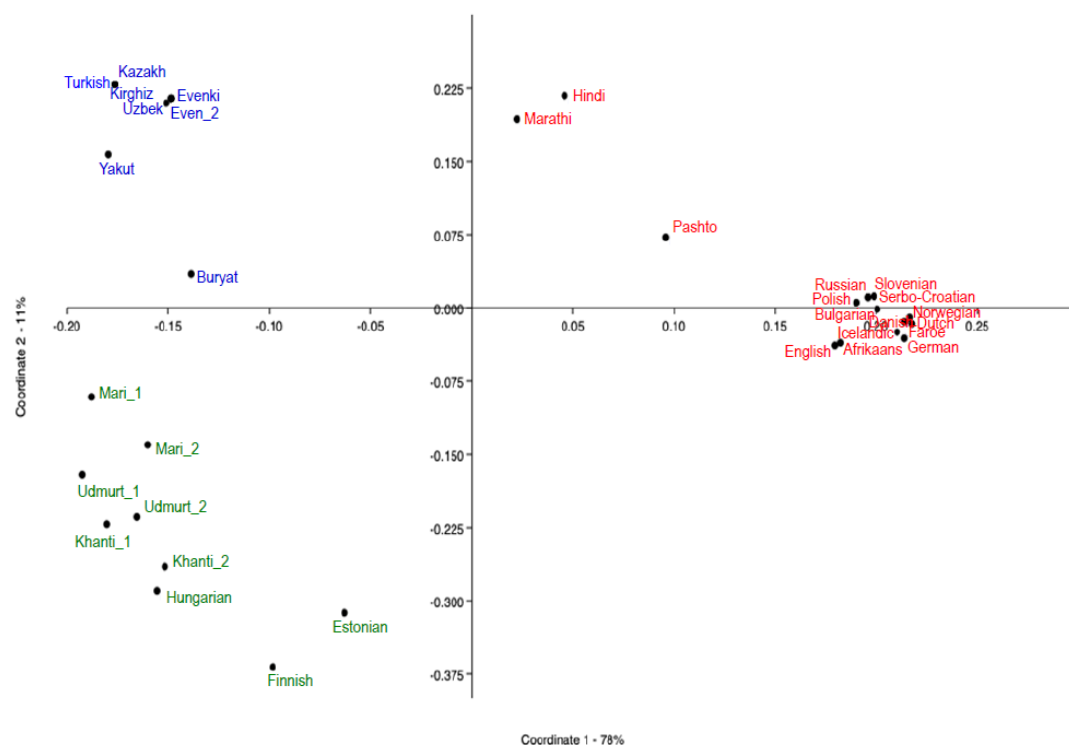
## 2.6. Linguistic Distances and Phylogenies

A matrix of pairwise syntactic distances was derived from the data matrix for the 34 relevant languages, by means of a Jaccard–Tanimoto formula (see [24]) reported below:

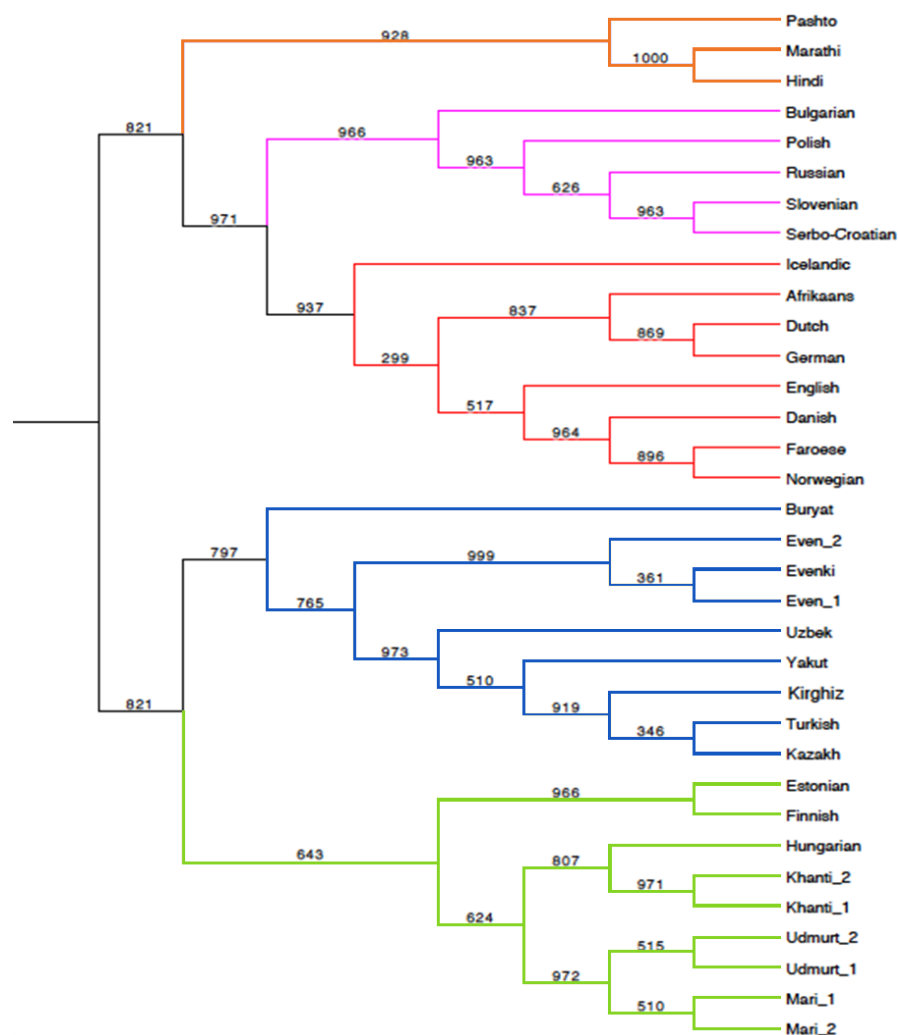
$$\Delta \text{Jaccard (A,B)} = [N_{-+} + N_{+-}] / [N_{-+} + N_{+-} + N_{++}]$$

where A and B are languages,  $N_{xy}$  indicates the number of positions where the string A has value X and B has value Y. The binary strings are interpreted as indicative of the presence (+) or absence (−) of traits (one per position in the string).

The distances inferred are summarized in Supplementary Figure S2 and visualized in a heatmap (Supplementary Figure S3). By means of a Principal Coordinate Analysis (PCoA, also called MDS, Multidimensional Scaling), calculated using the software PAST [47], we visualized the syntactic relationships between languages (Figure 1). We also represented the syntactic data in tree form through a UPGMA tree (Figure 2) using the bootstrapping procedure described in [24], in combination with the software PHYLIP [48] and Mesquite [49], and through a character-based Bayesian tree, using BEAST v1.83 [50] (Supplementary Figure S4).



**Figure 1.** Principal Coordinate Analysis (PCoA) from the syntactic distances in 34 Eurasian populations. Language groups coded as follow: Finno-Ugric (**green**), Altaic (**blue**), Indo-European (**red**).



**Figure 2.** UPGMA tree inferred from the Jaccard syntactic distances. Bootstrap values, base = 1000, at the nodes. Orange = Indo-Iranian (IE—Indo-European), pink = Slavic IE, red = Germanic IE, Blue = (AL—Altaic), Green = (FU—Finno-Ugric)

## 2.7. ChromoPainter and fineSTRUCTURE

ChromoPainter [51] (version v2) is a method to quantify distances between individual genomes. This method uses sampled chromosomes as “donors” and matches (or “paint”) other chromosomes to the donors’ DNA, thus quantifying similarities among individuals based on shared blocks of SNPs. In the heatmap, each square represents the number of DNA segments that each row (recipient) copies from each column (donor).

We used ChromoPainter output to cluster individuals into genetically homogeneous groups using fineSTRUCTURE [51] (version 2.1.3), a powerful approach for inference of fine-scale population structure from haplotype data. Each individual is presented as a matrix of non-recombining genomic chunks received from a set of multiple donors. Clusters of individuals are then inferred from the patterns of similarities among these copying matrices, by a Bayesian approach, and the tree is finally plotted using FigTree (version 1.4.2).

## 2.8. Outgroup $f_3$ -Statistics

We performed an  $f_3$  analysis using the *qp3Pop* package in ADMIXTOOLS (version 412). The outgroup  $f_3$ -statistic ( $X, Y$ ; Outgroup) is a function of shared branch length between two genomes, say  $X$  and  $Y$ , in the absence of admixture with the outgroup.  $Y$  is extracted from a set of individuals,

among whom we look for the most closely related to the individual under exam (X). Throughout the analysis we used the African Yoruba as an outgroup that we assumed to diverge from population X before all the other populations were analyzed. In this analysis, high values of  $f_3$  indicate that X and Y are genetically closer.

The modern samples from Pagani et al. (2016) [40] used in this study were merged with the Yamnaya, Anatolian, Sintashta and Nganasan individuals from Ref. [38] and used as source populations. Variants with a missing call rate exceeding 10% were excluded, resulting in 249,286 SNPs suitable for the analysis.

## 2.9. Modelling Admixture

Using the *qpAdm* package in ADMIXTOOLS (version 412), we estimated the proportions of ancestry in a *Test* population deriving from a mixture of three reference populations by leveraging shared genetic drift with a set of outgroup populations. The reference populations used were: Yamnaya, Anatolia and Nganasan (used here as a proxy for the genetically still undescribed Siberian population). As outgroup populations, we used: Han, Mbuti, Karitiana, Ulchi and Mixe. The detail: YES parameter, was set, which reports a normally distributed Z-score for the goodness of fit of the model (estimated with a Block Jackknife).

## 3. Results

### 3.1. Linguistic Analyses

#### Syntactic Comparison

The PCoA inferred from syntactic data (Figure 1) shows a first, neat division between the IE languages, with positive values of the first component (accounting for 78% of variation), and FU and AL, all found in the left area of the graph. In that area, the second PC (accounting for 11% of variation) separates FU from AL. In sum, each group appears to form a well-defined cluster. While the clouds corresponding to IE and AL are compact, although with individual outliers (Indo-Iranian and Buryat, respectively), the FU languages appear more scattered. Finnish and even more so Estonian fall particularly close to the IE languages. Such a resemblance between the Balto-Finnic group of FU and IE emerges even more neatly in the Bayesian phylogenetic analysis (Supplementary Figure S4), where the Balto-Finnic node joins the IE cluster rather than the FU one. The second important aspect that emerges from the PCoA is a split between IE and the other two groups, which might in turn hint at some closer FU–AL relatedness.

In the UPGMA tree (Figure 2), languages from the same family, IE, FU and AL, neatly cluster together without exception; FU languages form a monophyletic cluster within which the Balto-Finnic (Finnish and Estonian) and Ugric (Hungarian and Khanty) families are well identified, with the addition of a further node comprising geographically closer Udmurt (Permian) and Mari (Volgaic) [29,52]. The latter node occurs closer to Ugric than to Balto-Finnic, in disagreement with some traditional, though not uncontroversial [52], classifications.

The outlying positions of Balto-Finnic (Finnish and Estonians), Indo-Iranian (Pashto, Marathi and Hindi) and Mongol (Buryat) within the three groups are also visualized in the Heatmap of the syntactic distances (Supplementary Figure S3).

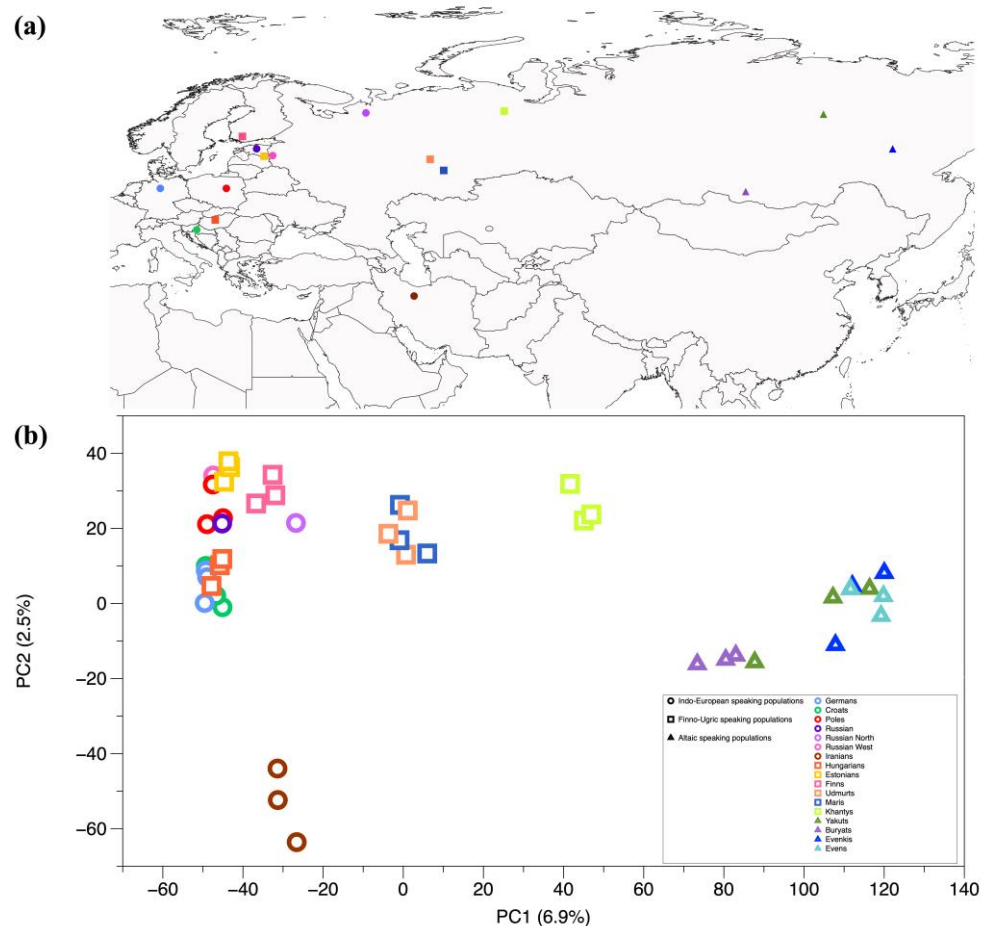
### 3.2. Genetic Comparison

#### 3.2.1. Population Structuring in Eurasia

We selected 17 populations—seven speaking IE, six FU and four AL languages—for which whole-genome data were available (Figure 3a; Supplementary Table S1). The first principal component (Figure 3b) mostly reflects geography and separates eastern from western Eurasian populations, whereas the second component separates western Eurasians along a north–south cline. The AL-speaking



populations fall into a single cluster along the first PC axis. The European IE-speaking populations form a cluster along the PC2 axis, separated from the Iranians, the latter belonging to the Asian group of IE languages.

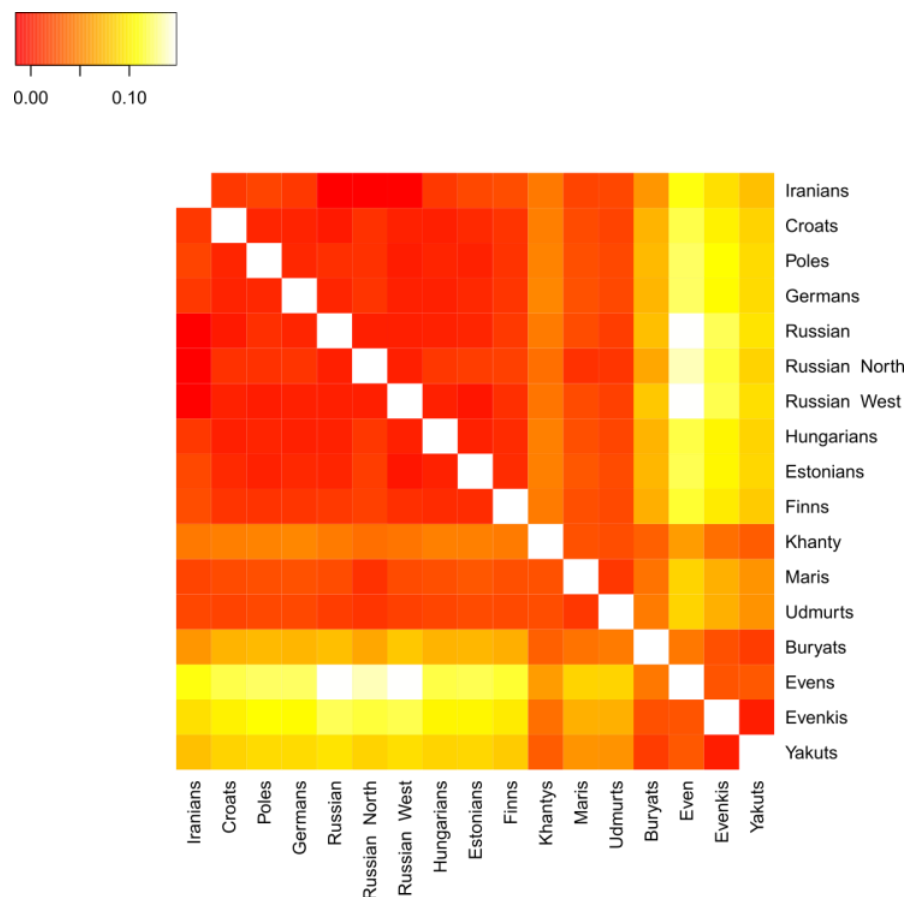


**Figure 3.** Geographical locations and Principal Component Analysis (PCA) of genomic variation. Populations speaking an IE, FU and AL language are represented by circles, squares and triangles, respectively. (a) Geographical locations of the samples in this study. (b) Projection on two dimensions of the main components (PCA) of genomic variation in IE, FU and AL speaking populations.

Conversely, the FU-speaking populations are scattered along the PC1 axis. Estonians fall within the IE diversity at the negative end of the X-axis, while Finns occupy an intermediate position between the IE speakers and the FU-speaking Udmurt and Mari people, i.e., the modern inhabitants of the Pontic steppes (Figure 3b).

### 3.2.2. Genetic Distances between Populations

Next, we calculated genetic distances ( $F_{st}$ ) between pairs of populations (Figure 4). All AL and IE speaking populations are genetically closer to other populations of their language family than to populations belonging to a different family. Instead, that is not the case for the FU speakers; all of Estonians, Finns and Hungarians are genetically closer to their respective European neighbors speaking IE. In addition, among the eastern populations, the Mari and Udmurt seem genetically more similar to the other Europeans than to the AL speakers. Exceptions are the easternmost and Trans-Uralic Khanty (Ostyaks), which seem equally close to Mari, Udmurt and most of the AL speakers. This observation can be reconciled with historical data, which place the origins of the Khanty people in the Russian steppes followed by a northward migration into western Siberia in about 500 AD (500 BCE) [53].



**Figure 4.** Pairwise genetic distances between Eurasian populations. Darker colors indicate that populations are genetically closer, whereas lighter colors indicate that populations are genetically distant.

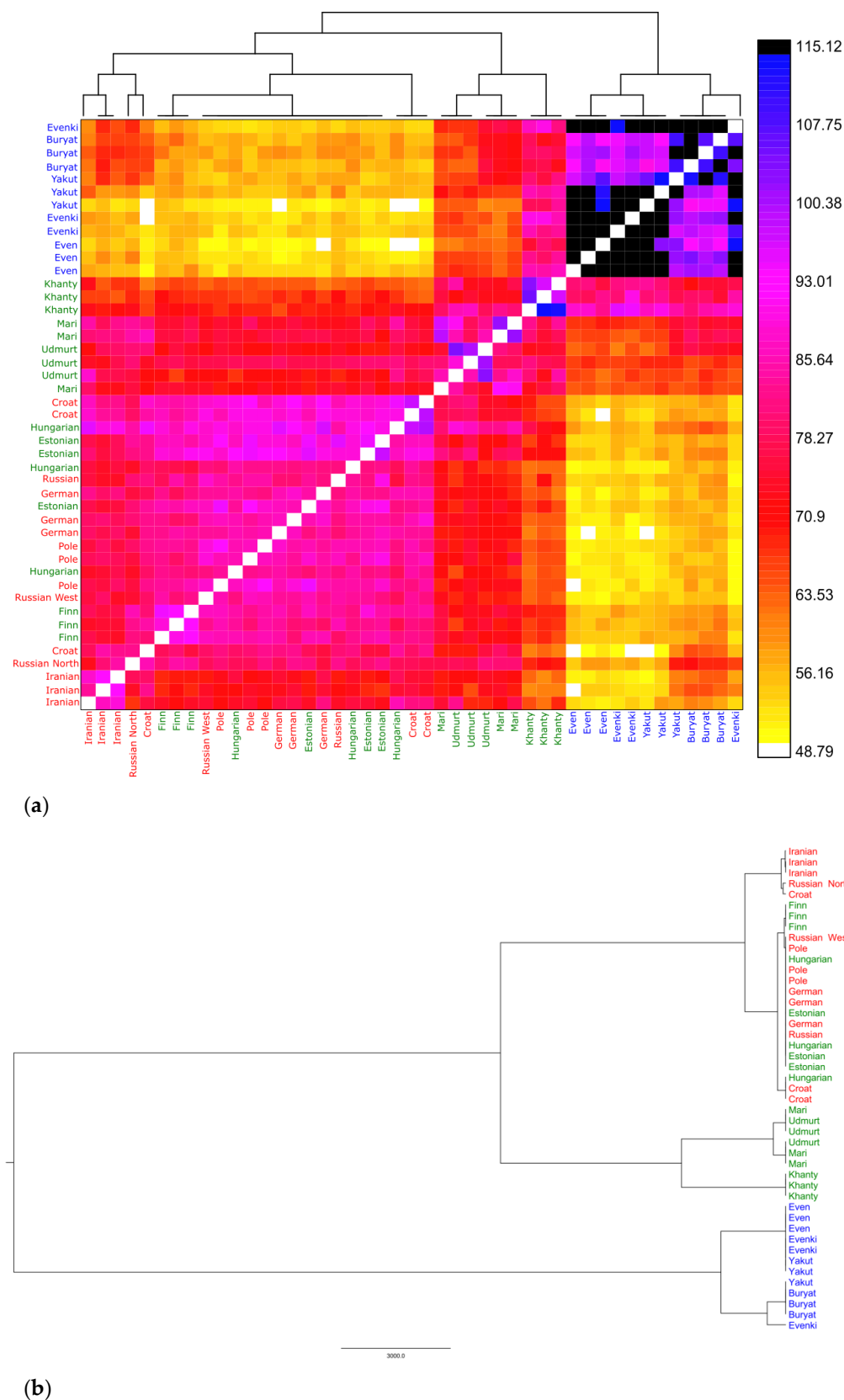
### 3.2.3. Shared Haplotypes

In the analysis of genetic distances, each single-nucleotide polymorphism is independently considered, regardless of its association with other polymorphisms. To analyze the patterns of population resemblance in finer detail, we thus moved to the haplotype level, using ChromoPainter and fineSTRUCTURE (Figure 5). This approach does not depend on prior information on sample groupings and operates instead with data-driven natural groups defined by patterns of haplotype sharing.

This approach also led us to identify three main genetic groups, broadly corresponding to the three main language families. However, as already observed in the *Fst* analysis, there were exceptions. The western FU-speaking populations (Estonians, Finns and Hungarians) seem to mainly share co-ancestry with the other Europeans, regardless of the language spoken. Conversely, among the eastern FU speakers, Udmurt, Mari and Khanty, there is a high level of haplotype sharing. In addition, this analysis revealed for the first time some co-ancestry of Finns (and partly Estonians' and Hungarians') with AL speakers of Siberian origin.

The evolutionary tree inferred from these data (fineSTRUCTURE cluster analysis; Figure 5b) shows two deep splits, the first isolating all AL speakers, and the second separating eastern FU speakers from a group composed by western FU and IE speakers. All this could even point to different ancestries for the FU-speaking populations, with phenomena of horizontal language diffusion leading them to a shared linguistic identity. However, lexical analyses and, in a more modulated fashion, even the syntactic ones support an original FU linguistic unity, later fragmented by northward and westward migrations and contacts. To better understand these results, we resorted to ancient DNA.





**Figure 5.** Estimates of shared ancestry between Eurasian individuals. (a) Co-ancestry heatmap. Each of the 51 individuals is represented as a row, where each pixel represents the level of co-ancestry (higher for darker colors) shared with each of the other individuals. (b) fineSTRUCTURE cluster analysis obtained from the co-ancestry matrix. Red = IE; Green = FU; Blue = AL.

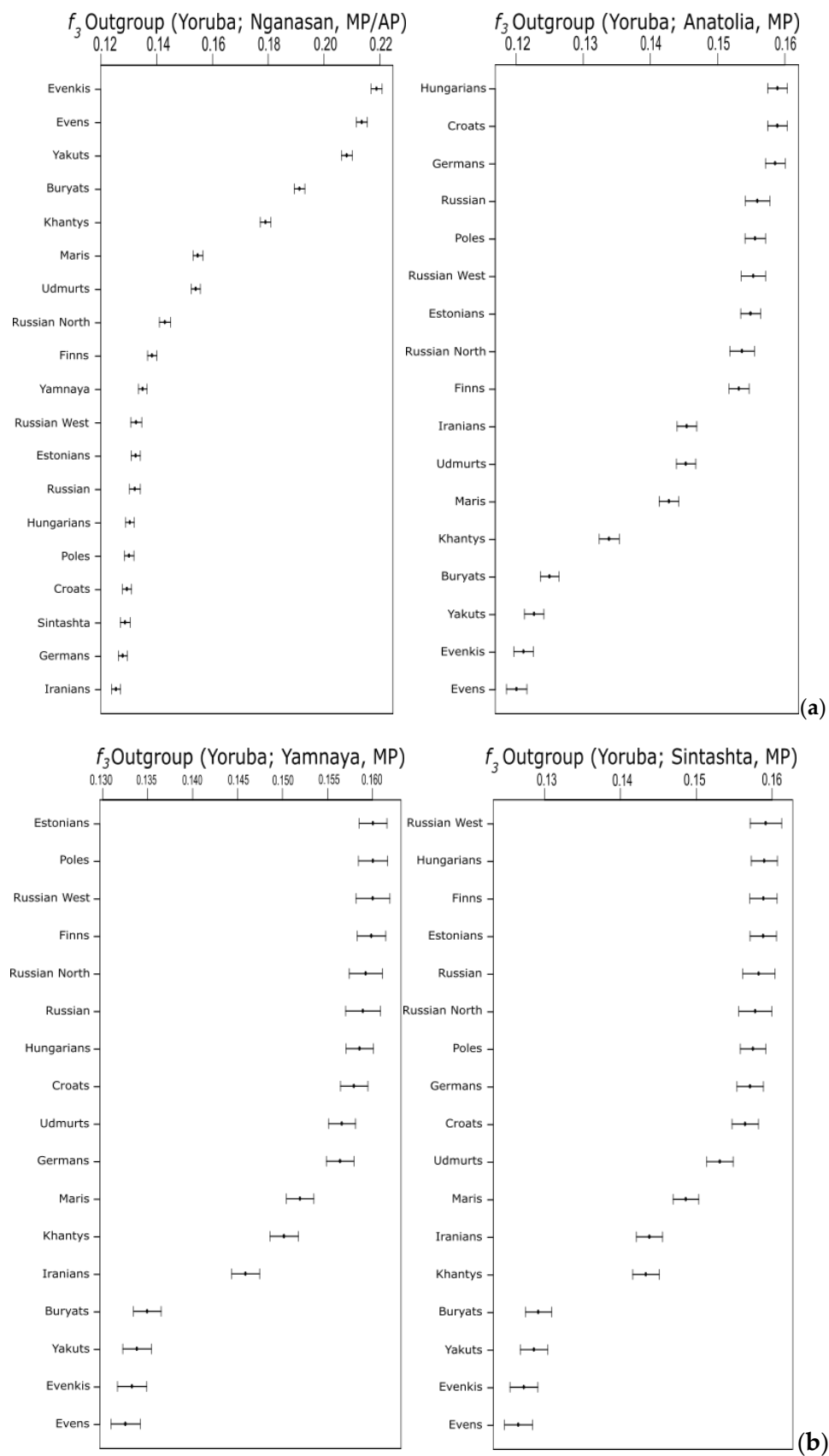
### 3.2.4. Affinities between Modern and Ancient Populations

Our genetic analysis showed the Udmurt and Mari to be closer than the Khanty to European populations (Figures 4 and 5a). We hypothesize that this observation may be related to shared ancestry with Yamnaya, an ancient pastoralist population that lived in the current Udmurt and Mari territories, around the Pontic-Caspian steppes, and that expanded into Central and western Europe in the third millennium BCE, contributing a Caucasian genomic component that nowadays is widespread in Europeans [35,37]. We tested for genetic continuity from the ancient Steppe populations, Yamnaya (~4700 yBP) and the more recent Sintashta (~3900 yBP) on the one hand [35,36], to current Udmurt and Mari on the other. An ancient Anatolian sample [54] was also included in our tests, potentially accounting for the genetic legacy of early farmers from the Near East.

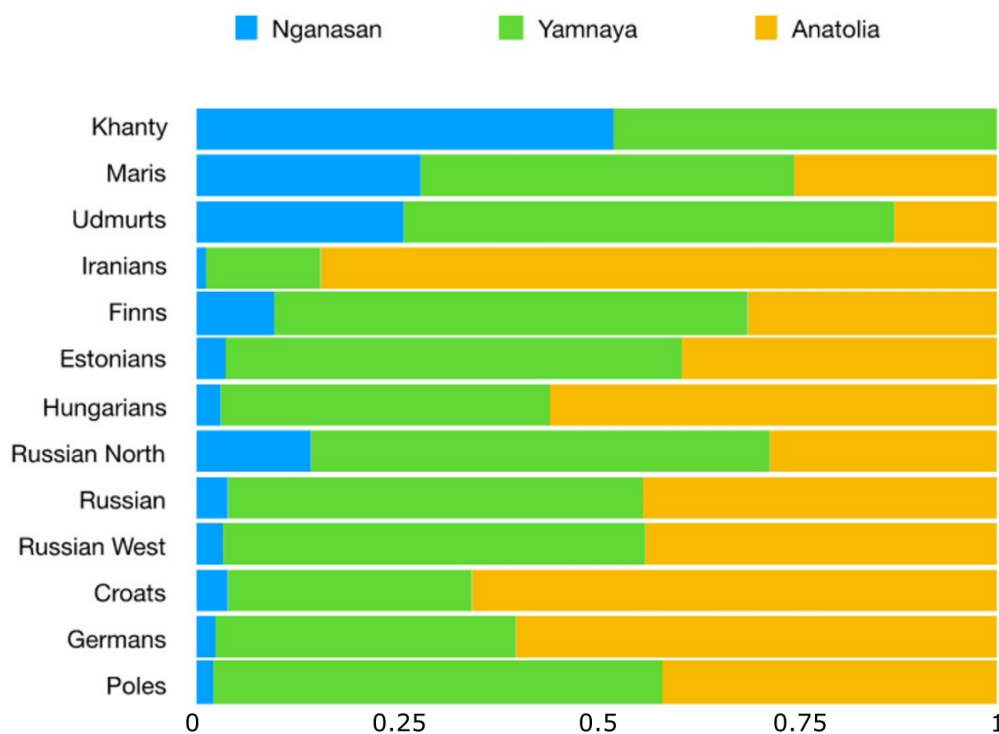
We formulated outgroup  $f_3$ -statistics of the form  $f_3(AP, MP; Yoruba)$ , where  $AP$  was represented in turn by each of the three ancient populations, and  $MP$  was each of the modern samples in our dataset (Figure 6 and Supplementary Figure S5). In general, we found all ancient samples to share more genetic drift with modern Europeans and Russians than with non-European populations. Among the eastern populations, the Udmurt and Mari are the ones sharing the most genetic drift with Yamnaya and Sintashta; on the other hand, the Iranians (IE) are the Asian sample closest to the Anatolian farmers, in agreement with recent findings [37]. In addition, within the European populations, the  $f_3$  values show opposite trends for the Anatolian and the Yamnaya/Sintashta, the former sharing more genetic drift with southern and Central Europeans (Croats and Germans) and the latter being closer to Northeast Europeans, including the FU-speaking Estonians and Finns, once again in general agreement with previous findings (e.g., Ref. [35]). It is interesting to notice the peculiar behavior of the Hungarians. They appear much closer to the ancient Anatolians than to the Yamnaya, which is common among southern European populations; however, they are the modern Europeans sharing most genetic drift with the Sintashta. This may be indicative of a relatively more recent genetic contact between them and the Steppe populations, i.e., after the process leading to the spread of the Yamnaya component into Europe.

Contrary to what could be expected, the modern FU inhabitants of the Russian steppes, Mari and Udmurt, appear more distant from Yamnaya than Estonians and Finns. One possible explanation would be the presence, in their genomes, of a Siberian-related component, known to be widespread in contemporary Central and North Asian populations [55–58]. We tested for its presence in our samples by modelling Nganasan, a population of residual speakers of a moribund Samoyed language (i.e., distantly related to FU) from the Taymyr Peninsula, as a proxy of the carriers of this Siberian component (as also in Refs. [33,38]). We did find support for the presence of such a Siberian component among Mari and Udmurt; the outgroup  $f_3$  statistics of the form  $(Nganasan, MP/AP; Yoruba)$  showed that Udmurt and Mari are indeed closer to Nganasan than Yamnaya, which shared similar  $f_3$  values with other European population with regards to Nganasan. Figure 6b shows a clear trend; the Nganasans share more genetic drift with all AL speakers, followed by Udmurt and Mari, and then by European populations, no matter if FU- or IE speakers.

To further test whether the peculiar genetic position of the Udmurt and Mari is really associated with the higher presence of a Siberian genetic component in their genome, we ran a  $qpAdm$  analysis (Figure 7 and Supplementary Table S4). All the FU-speaking populations were successfully modelled as a mixture of Yamnaya, Anatolian and Nganasan-related ancestry, with the exception of the Khanty, who seem to have no Anatolian ancestry. In particular, the Mari and Udmurt genomes appear to contain a large component that can be related with a Siberian genetic ancestry, confirming our expectations. Furthermore, this Siberian ancestry is present, at low though non-negligible percentages, in the western FU-speaking Finns (but less saliently in Estonians).



**Figure 6.** Outgroup  $f_3$ -statistic analysis. Shared genetic drift between ancient and modern (MP) populations. (a) Shared genetic drift between Anatolian, Yamnaya, Sintashta, (b) Nganasan and modern/ancient populations.



**Figure 7.** Admixture proportions from three sources estimated using *qpAdm*. Sources used were Nganasan, Yamnaya and Anatolia (percentages and chi-square values are shown in the Supplementary Table S4).

## 4. Discussion

### 4.1. Syntactic Diversity

Syntax distinguishes IE, FU and AL languages quite well, although IE and AL have single outliers (Indo-Iranian and Buryat, respectively). Conversely, the FU family turns out to be less compact, in spite of the greater geographic spanning and population size of IE, and of the weaker lexical evidence purportedly supporting AL (see Ref. [12] for the state of the debate). The whole family appears scattered and in some structural contiguity with their eastern and western neighbors.

A previous study comparing Uralic lexical data, including the FU speakers, had suggested that some secondary contact played a role in the divergence of these languages [59]. The scattered pattern is now more clearly observed and measurable through the cross-family application of our syntactic analysis. The outlying position of Estonian and Finnish among FU languages is evident (especially see Figure 1 and Supplementary Figures S3 and S4). As for the other western FU language, Hungarian, qualitative analysis shows that the language shares some parameter values with IE, as opposed to the rest of FU. Yet, such similarities do not emerge in the trees, possibly reflecting the much later arrival in Europe of the Hungarian language [60].

The very distribution of similarities and differences in the syntactic parameters suggests that the scattering of FU languages is likely to be secondary, i.e., due to cultural contacts. Indeed, there is no evidence of potential convergence of Khanty, Udmurt, Mari with IE. In addition, the main syntactic changes detaching Finnish and especially Estonian on the one hand, and Hungarian on the other, from the other FU languages are that they are: A. different from each other; B. unidirectional, i.e., of a kind that is often acquired but hardly reversed; C. shared with neighboring IE languages at the time of the contact with the respective FU languages [24]. This tends to exclude that these properties might be ancestral (proto-Uralic) and lost in the more eastern varieties because of recent convergence with Asian languages. The similarities of eastern varieties with AL languages are, instead, more ambiguous as to whether they may be shared inheritance or a secondary effect.

In sum, our syntactic phylogenetic analysis supports the original wisdom that FU has been a monophyletic cluster, and is well compatible with the traditional view that the western FU languages have reached Europe from the East at some ancient point; but syntax also detects and measures the pattern of secondary similarities with neighboring languages.

#### 4.2. Genome Diversity

The three main population groups identified by the linguistic analysis are also biologically differentiated; however, while IE and AL samples form distinct genetic clusters, both in the PCA and ChromoPainter analyses, a peculiar pattern emerges within the FU language family. While the Khanty show affinities with a well-differentiated cluster comprising all AL speakers, the other FU speakers appear to be part of a broad group, including all IE-speaking individuals. In particular, the western FU speakers, namely, Finns, Estonians and Hungarians, are genetically closer to IE populations in Europe than to the Asian UR-speaking populations. Estonians and Finns also share more ancestry with each other than with the Hungarians. This genetic similarity can reflect: (i) a different source of steppe ancestry in the Hungarians (more closely related with the Sintashta) than in Finns and Estonians (genetically closer to the Yamnaya) (Figure 6a); and/or (ii) a lower contribution of Siberian ancestors to the Hungarian genomes than to the Estonians and especially the Finns (Figure 6b).

#### 4.3. Comparison of Genetic and Linguistic Results

Judging whether or not linguistic and genetic data mirror each other may be partly a matter of taste. However, there is little doubt that the syntactic and genomic findings of this study match and corroborate each other. In five out of six cases, linguistic and genetic evidence were consistent (Table 1), the only exception being the third one.

**Table 1.** Synopsis of the main results of this study.

	Syntax	Modern Genomes	Ancient Genomes
1	AL languages form a cluster	AL speakers form a cluster	Higher Siberian component in AL speakers than in all the other populations
2	Indo-Iranian languages distinct from European IE languages	Indo-Iranian speakers distinct from other IE speakers	Higher Anatolian component in Indo-Iranian speakers than in other IE speakers
3	FU languages separated from IE and AL	In the tree, FU speakers and IE speakers fall in the same cluster	Yamnaya and Anatolian components similar in western FU speakers and their European IE-speaking neighbors
4	Estonian closer to IE and more distant than Finnish from other FU languages	Estonians closer to IE speakers than Finns	Siberian component lower in Estonians than in Finns
5	Mari, Khanty and Udmurt closer to AL than to IE languages	Mari, Khanty and Udmurt speakers more distant from IE speakers than Finns, Estonians and Hungarians	Higher Siberian component in Mari, Khanty and Udmurt speakers than in any other FU population
6	Easternmost FU Khanty least distant from easternmost Yakut of all AL languages	Khanty speakers halfway between the Mari/Udmurt speakers and eastern AL populations	Khanty speakers have the Siberian and Yamnaya component, but no Anatolian one

Note that ancient Siberian ancestry is (here and elsewhere: Refs. [55,58]) approximated by a modern population, Nganasans.

In this field, however, exceptions are as interesting as the rules, as they call our attention to phenomena that need be further investigated. By looking into the syntactic features of western FU languages, and into their speakers' genomes, we could recognize peculiar processes affecting the demographic history of people speaking Estonian, Finnish and Hungarian.

Indeed, the Bayesian syntactic tree matches the strong similarity between IE and Balto-Finnic revealed by the genomic tree and PCA, but, on the whole, syntax supports the FU unity to a stronger extent than genetics, and neatly recognizes the Ugric group (Hungarian and Khanty). On the contrary, at the genetic level, the FU-speaking populations cluster according to geography, with Hungarian speakers close to Central Europeans, Khanty speakers close to their eastern AL-speaking neighbors, and the steppe-dwelling Mari and Udmurt speakers in an intermediate position. This result suggests that syntax can also capture secondary demographic events (e.g., population admixture), which genetics can identify only if they have entailed substantial demographic change.

In particular, syntax shows more limited secondary effects on Hungarian from its IE geographic neighbors and preserves well its historical similarity with Khanty. As we shall discuss later in this paper, historical data suggest that the establishment of Hungarian in Central Europe was the product of an episode of elite dominance, i.e., a deep change of language with limited demographic impact [4]. If so, the genomes of modern speakers of Hungarian were affected only marginally by the phenomena that radically replaced their language.

#### 4.4. Demographic Scenarios: Linguistic and Genetic Evidence

The conclusions above, especially as summed up in the Synopsis of Table 1, reinforce our trust in the usefulness of gene/language comparison as a heuristic tool. We consider now what such approach tells us about the prehistory of a portion of western Eurasia.

While little is known about the ancient demographic history of AL populations, data are now available from pre-historic peoples of western Eurasia and the Near East. Analysis of genomes from pre-historic inhabitants of the Russian Steppes, the Yamnaya [36,37], identified a westward migration of people who contributed an ancestral component today widespread in Europe. Although no linguistic evidence was presented in that influential study, the authors linked the westward Yamnaya migration with the expansion of the IE languages. This hypothesis is based on two implicit but reasonable premises: one, that IE languages have been uninterruptedly attested in all the core of Europe for over two millennia now, and, two, that a high congruence (and consequent reciprocal predictability) of genetic and linguistic diversity is characteristic of Indo-European-speaking Europe (see [4], for example).

At the same time, other studies have related the presence of the Estonian and Finnish languages in Europe, instead, to a northward migration of people of ultimately Siberian ancestry [55,58].

Our multidisciplinary analysis seems to necessarily point to a more intricate scenario. Consider two points, again: first, a highest prevalence of the ancestral Yamnaya component is now attested precisely among Estonians and Finns (cf. Figure 7); second, the genetic/linguistic comparisons summarized in Table 1 have uncovered that, at a finer-grained analysis, a substantial gene/language congruence exists within the FU family.

Linguistically, it is clear that speakers of AL, IE and FU languages (along with their Uralic cousins, the Samoyeds) have long formed three separate groups, which had contacts leading to various degrees of linguistic and biological exchange. The data cannot exclude that some even more ancient unity—or close contact—may have involved at least two of the groups, proto-Uralic and proto-Altaic speakers [24,39]. Specifically, our syntactic database confirms the existence of a neat Ugric subfamily of FU, and detects remarkable, though measurably different, degrees of Indo-Europeanization for Finnish and Estonian.

Now, the classical and solid linguistic relatedness of Mari and Udmurt with the Balto-Finnic languages (partly obscured in the syntactic phylogenies because of the well-justified parametric interference of the latter with IE), as well as the genetic relationships of these modern populations (Mari and Udmurt in the steppes and Estonians and Finns in north-eastern Europe) with the ancient Yamnaya (Figures 6a and 7), suggests a demographic and linguistic northward expansion of people with steppe-related ancestry into the Baltic area. This is consistent with the fact that Balto-Finnic preserves ancient Indo-Iranian (hence southern) IE loanwords more than other Uralic varieties and with the general hypothesis of an expansion of the FU languages from the Volga river basin [26,28,61,62].



A likely link of the expansion of FU languages into northeastern Europe with that of the Yamnaya (and not direct Siberian) genomic ancestry is of a chronological order. Some possible dates of the FU diversification, estimated between 5000 and 4000 yBP in [26,60,62], coincide in time with the first appearance of the Yamnaya genomic component in ancient European populations [35,36,41]. Particularly in the Baltic region, analysis of ancient DNA has dated the first contacts between Yamnaya migrants and the local communities around the early Bronze Age (5000–4500 yBP), involving Baltic hunter–gatherers with no Neolithic ancestry [63]. Such chronological evidence overlaps with the early dates of lexical borrowing between IEs and Balto-Finns mentioned above [25–27].

A reconstruction of the lexical history of the FU languages suggests a later northward expansion from the Baltic area to southern Finland, possibly around 2000–1600 yBP [60]. This expansion was accompanied by the separation of the northern (Finnish) from the southern (Estonian) group, which may have involved two successive linguistic steps [29]. It is this last diversification that overlaps in time with the first appearance of the Siberian genomic component in ancient remains around the Baltic area [58]. Therefore, ancient DNA and archaeological studies agree in suggesting that people related to the Yamnaya culture moved from the coastal Baltic areas into southern Finland around 2000–1600 yBP (i.e., much later than their linguistic presence around Baltic is first attested by their lexical exchanges with Slavic and Germanic tribes), where they first came into contact with the ancestors of modern Saami (Lapps) [57,58,60,64,65]. Traces of this contact, and of the limited admixture that must have followed, are still detectable in the genomes of Finns [66] (also in Figure 7) more than in the Estonian ones. The different degrees of syntactic similarity we were able to quantify in this study between Estonian and Finnish with respect to IE, on the one hand, and the other FU languages, on the other (Figure 1 and Supplementary Figure S2), seem in agreement with this secondary migration model (although by themselves they could not be informative about the direction of such exchanges, whether South-to-North or North-to-South).

Finally, our genomic analysis is readily compatible with a second event of expansion of the FU languages (possibly through the Russian steppes) into Europe without involving Siberian mediation. That is the case of the FU speakers in Hungary. There is historical evidence that at the beginning of the Medieval era, the language spoken in nowadays Hungary was still Late Latin (at least as an official language), later subject to the effects of Slavic, Germanic and Avar invasions [67]. The main linguistic shift can be approximately dated around 895–905 AD, when people coming from the East conquered Hungary, imposing their own Ugric language [67,68]. Ancient DNA studies of the invaders have shown that they were genetically close to the Sintashta of the steppes, and apparently unrelated with Siberian ancestors [69], in fine agreement with our genetic analysis. Therefore, the presence of a FU language in Europe is not correlated with the presence of a Siberian component in the DNA of its speakers.

#### 4.5. Speculations on the Diffusion of IE into Europe

Linguists and archaeologists have long discussed the timing and modes of spread of IE languages in Europe. Gimbutas [30] associated it with the westward spread of the Kurgan culture, from the Pontic steppes during the Bronze age, whereas Renfrew [33] saw it as a consequence of the Neolithic farmers' demic diffusion from Anatolia (see also Refs. [16] and [70]). These alternatives (the Steppe and the Anatolian hypothesis, respectively) are paralleled at the genetic level, by studies supporting demic dispersal, respectively, of Yamnaya-related populations in the Early Bronze Age [35,36], and of Anatolia-related populations during the Neolithic transition [2,3,7,38,71,72].

The genomic similarity between the Yamnaya and the first FU speakers of Europe may be difficult to reconcile with the view that the Yamnaya were also the first who introduced IE languages in Europe, as suggested by studies of genomic data not supported by linguistic analyses [35,36]. One possibility, supported by a study of Iberia [73], is that the arrival in Europe of the Yamnaya genomic component did not necessarily entail the same linguistic changes in all areas. In the absence of adequate data to formally test this hypothesis, we still may speculate that the small, but non-negligible,

ancestry component associated with the Anatolian Neolithic [38] among the Yamnaya may reflect the previous northward gene flow from the Near East into the Pontic steppes. If so, it would be possible to reconcile genetic evidence for the Neolithic demic diffusion from the Near East, linguistic evidence on a Near East center of IE diffusion [16,33,38,70,74], and data suggesting a role of Yamnaya people in spreading both IE [35,36,75] and FU (this study) languages by imagining the existence of some linguistic diversity within the Yamnaya-like populations and concluding that IE languages have entered Europe in two moments and by two routes. The first one would correspond to the main Neolithic expansion, Northwest into southern and then Central Europe, but also North, towards the Pontic Steppes. The linguistic impact of this migration would have not been the same for all people in the Pontic steppes; some would retain their original FU languages, and some would acquire an IE language. The former would then mostly move towards the Baltic Sea area, whereas the latter would correspond to the IE-speaking populations dispersing in Central Europe in the Bronze Age [36,76], giving rise to the Bell Beaker and Corded Ware cultures.

## 5. Conclusions

This study exemplifies how appropriate quantitative and qualitative tools allow one to measure cross-family language variation, offering a novel insight into human prehistory and generating testable hypothesis for large-scale genomic analyses. Of course, we must warn about the risk of over-interpreting correlations between languages and genes, especially in the absence of accurate dates of linguistic diversification and expansion, which are not yet well established. Furthermore, full inference of complex processes requires the study of broader datasets than available for the present study.

Nonetheless, on the whole, our analysis, taking advantage of various linguistic features, including the syntactic ones, which can be compared across families and are sufficiently stable in time, suggests that Darwin's prediction of a general correspondence between biological evolution and language transmission is still generally valid, and that exceptions to this rule are both limited (more than it may appear from simply relying on traditional and non-quantitative language taxonomies, such as, e.g., in [68]) and extremely useful for a detailed reconstruction of human past.

**Supplementary Materials:** The following are available online at <http://www.mdpi.com/2073-4425/11/12/1491/s1>, Figure S1. Approximate geographical location of the 34 languages considered. Figure S2. Syntactic Distances. Figure S3. Heatmap from the syntactic distances. Dark red represents maximum distance, dark blue minimum distance. Figure S4. Bayesian phylogeny (BEAST) from the syntactic dataset. Figure S5. Outgroup  $f_3$ -statistics analysis. Table S1. Whole-genome samples collected for the populations under study. Table S2. Ancient DNA samples used in this study. Table S3. Human Origins data on present-day humans used in this study. Table S4. Statistics of the  $qpAdm$  models.

**Author Contributions:** Conceptualization, G.L., C.G. and G.B.; methodology, P.S., G.G.-F., G.L., A.C., C.G. and G.B.; software, P.S. and G.C.; formal analysis, P.S., G.G.-F., E.T., A.C. and G.C.; investigation, P.S., G.G.-F., E.T., A.C. and G.C.; genetic data curation, P.S.; linguistic data A.C., C.G. and G.L.; writing—original draft preparation, P.S., G.L. and G.B.; writing—review and editing, G.G.-F., C.G., G.L. and G.B.; visualization, P.S.; supervision, G.L. and G.B.; funding acquisition, P.S., G.L., C.G. and G.B. All authors have read and agreed to the published version of the manuscript. All listed authors meet the ICMJE criteria and all who meet the four criteria are identified as authors. We attest that all authors contributed significantly to the creation of this manuscript, each having fulfilled criteria as established by the ICMJE. We confirm that the manuscript has been read and approved by all named authors. We confirm that the order of authors listed in the manuscript has been approved by all named authors.

**Funding:** This research was funded by the ERC Adv.Gr. 295733 *Darwin's Last Challenge (LanGeLin)* 2012–2018 (PI Giuseppe Longobardi, Co-I Guido Barbujani); by the MIUR PRIN 2017K3NHHY *Models of language variation and change: new evidence from language contact* (C. Guardiano); by the French National Research Agency under the IDEX Bordeaux NETAWA Emergence project no. ANR-10-IDEX-03-02 'Out of the Core: Exploring social NETWORKS at the dawn of Agriculture in Western Asia 10 000 years ago'; and by the CNRS momentum project "Symboling and Neighboring at the Dawn of Agriculture in Europe 8000 years ago".

**Acknowledgments:** We thank Andrea Benazzo for his bioinformatics support. Part of the analyses was carried out while P.S. was Hugo Reyes-Centeno's guest at the DFG (German Research Foundation) Centre for Advanced Studies "Words, Bones, Genes, Tools: Tracking Linguistic, Cultural and Biological Trajectories of the Human Past", at the University of Tübingen. We are also indebted to Ándras Bárány, Judit Gervain, Anders Holmberg, István Kenesei, Paul Kiparsky, Katalin Kiss, and Márton Sósokuthy for help with Finno-Ugric evidence and analyses, and to Monica-Alexandrina Irimia and Nina Radkevich for assistance in collecting more language data.

**Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

## References

1. Darwin, C. *On the Origin of Species by Means of Natural Selection, or the Preservation of Favoured Races in the Struggle for Life*; John Murray: London, UK, 1859.
2. Sokal, R.R. Genetic, geographic, and linguistic distances in Europe. *Proc. Natl. Acad. Sci. USA* **1988**, *85*, 1722–1726. [[CrossRef](#)]
3. Barbujani, G.; Pilastro, A. Genetic evidence on origin and dispersal of human populations speaking languages of the Nostratic macrofamily. *Proc. Natl. Acad. Sci. USA* **1993**, *90*, 4670–4673. [[CrossRef](#)] [[PubMed](#)]
4. Longobardi, G.; Ghirotto, S.; Guardiano, C.; Tassi, F.; Benazzo, A.; Ceolin, A.; Barbujani, G. Across language families: Genome diversity mirrors linguistic variation within Europe. *Am. J. Phys. Anthr.* **2015**, *157*, 630–640. [[CrossRef](#)] [[PubMed](#)]
5. Creanza, N.; Ruhlen, M.; Pemberton, T.J.; Rosenberg, N.A.; Feldman, M.W.; Ramachandran, S. A comparison of worldwide phonemic and genetic variation in human populations. *Proc. Natl. Acad. Sci. USA* **2015**, *112*, 1265–1272. [[CrossRef](#)] [[PubMed](#)]
6. Renfrew, C. Archaeology, Genetics and Linguistic Diversity. *R. Anthropol. Inst. Gt. Br. Irel.* **1992**, *27*, 445. [[CrossRef](#)]
7. Cavalli-Sforza, L.L.; Piazza, A.; Menozzi, P.; Mountain, J. Reconstruction of human evolution: Bringing together genetic, archaeological, and linguistic data. *Proc. Natl. Acad. Sci. USA* **1988**, *85*, 6002–6006. [[CrossRef](#)] [[PubMed](#)]
8. Poloni, E.S.; Passarino, G.; Santachiara-Benerecetti, A.S.; Langaney, A.; Excoffier, L.; Poloni, E. Human Genetic Affinities for Y-Chromosome P49a,f/TaqI Haplotypes Show Strong Correspondence with Linguistics. *Am. J. Hum. Genet.* **1997**, *61*, 1015–1035. [[CrossRef](#)] [[PubMed](#)]
9. Belle, E.M.S.; Barbujani, G. Worldwide analysis of multiple microsatellites: Language diversity has a detectable influence on DNA diversity. *Am. J. Phys. Anthr.* **2007**, *133*, 1137–1146. [[CrossRef](#)]
10. Henn, B.M.; Botigué, L.R.; Gravel, S.; Wang, W.; Brisbin, A.; Byrnes, J.K.; Fadhlouli-Zid, K.; Zalloua, P.A.; Moreno-Estrada, A.; Bertranpetit, J.; et al. Genomic ancestry of North Africans supports back-to-Africa migrations. *PLoS Genet.* **2012**, *8*, e1002397. [[CrossRef](#)]
11. Gyarmathi, S. *Affinitas Linguae Hungaricae Cum Linguis Fennicae Originis Grammaticae Demonstrata*; J.C. Dieterich: Göttingen, Germany, 1799.
12. Ceolin, A. Significance testing of the Altaic family. *Diachronica* **2019**, *36*, 299–336. [[CrossRef](#)]
13. Marcantonio, A. *The Uralic Language Family: Facts, Myths and Statistics*; Blackwell: Oxford, UK, 2002.
14. Nettle, D.; Harriss, L. Genetic and Linguistic Affinities between Human Populations in Eurasia and West Africa. *Hum. Biol.* **2003**, *75*, 331–344. [[CrossRef](#)] [[PubMed](#)]
15. Ringe, D.; Warnow, T.; Taylor, A. Indo-European and computational cladistics. *Trans. Philol. Soc.* **2002**, *100*, 59–129. [[CrossRef](#)]
16. Gray, R.D.; Atkinson, Q.D. Language-tree divergence times support the Anatolian theory of Indo-European origin. *Nature* **2003**, *426*, 435–439. [[CrossRef](#)] [[PubMed](#)]
17. Gray, R.D.; Drummond, A.J.; Greenhill, S.J. Language phylogenies reveal expansion pulses and pauses in Pacific settlement. *Science* **2009**, *323*, 479–483. [[CrossRef](#)] [[PubMed](#)]
18. Jäger, G. Support for linguistic macrofamilies from weighted sequence alignment. *Proc. Natl. Acad. Sci. USA* **2015**, *112*, 12752–12757. [[CrossRef](#)]
19. Longobardi, G. Methods in parametric linguistics and cognitive history. *Linguistic Var. Yearb.* **2003**, *3*, 101–138. [[CrossRef](#)]
20. Guardiano, C.; Longobardi, G. Parametric comparison and language taxonomy. In *Grammaticalization and Parametric Variation*; Battlori, M., Picallo, C., Roca, F., Eds.; Oxford University Press: Oxford, UK, 2005.
21. Longobardi, G.; Guardiano, C. Evidence for syntax as a signal of historical relatedness. *Lingua* **2009**, *119*, 1679–1706. [[CrossRef](#)]
22. Longobardi, G.; Guardiano, C.; Silvestri, G.; Boattini, A.; Ceolin, A. Toward a syntactic phylogeny of modern Indo-European languages. *J. Hist. Linguist.* **2013**, *3*, 122–152. [[CrossRef](#)]

23. Longobardi, G.; Guardiano, C. Phylogenetic reconstruction in syntax: The Parametric Comparison Method. In *The Cambridge Handbook of Historical Syntax*; Ledgeway, A., Roberts, I., Eds.; Cambridge University Press: Cambridge, UK, 2017; pp. 241–272. ISBN 9781107049604.
24. Ceolin, A.; Guardiano, C.; Irimia, M.-A.; Longobardi, G. Formal syntax and deep history. *Front. Psychol.* **2020**, *11*, 2384.
25. Kylstra, A.D.; Hahmo, S.-L.; Hofstra, T.; Nikkilä, O. Lexikon der 2110 Älteren Germanischen Lehnwörter in den Ostseefinnischen Sprachen. In *Band I: A-J*; Rodopi: Amsterdam, The Netherlands, 1991.
26. Koivulehto, J. The earliest contacts between Indo-European and Uralic speakers in the light of lexical loans. In *Early Contacts between Uralic and Indo-European: Linguistic and Archaeological Considerations*; Carpelan, C., Parpola, A., Petteri, K., Eds.; Suomalais-Ugrilainen Seura: Helsinki, Finland, 2001; pp. 235–263.
27. Sammalahhti, P. The Indo-European loanwords in Saami. In *Early Contacts between Uralic and Indo-European: Linguistic and Archaeological Considerations*; Carpelan, C., Parpola, A., Petteri, K., Eds.; Suomalais-Ugrilainen Seura: Helsinki, Finland, 2001; pp. 397–415.
28. Kallio, P. Phonetic Uralisms in Indo-European. In *Early Contacts between Uralic and Indo-European: Linguistic and Archaeological Considerations*; Carpelan, C., Parpola, A., Petteri, K., Eds.; Suomalais-Ugrilainen Seura: Helsinki, Finland, 2001; pp. 221–234.
29. Kallio, P. The Diversification of Proto-Finnic. In *Fibula, Fabula, Fact: The Viking Age in Finland*; Ahola, J., Frog, C.T., Eds.; Studia Fennica: Helsinki, Finland, 2014; pp. 155–168.
30. Gimbutas, M. The Three Waves of Kurgan People into Old Europe, 4500–2500 BC. *Archives Suisses D'anthropologie Générale* **1979**, *43*, 113–137.
31. Anthony, D. *The Horse, the Wheel and Language: How Bronze-Age Riders from the Eurasian Steppes Shaped the Modern World*; Princeton University Press: Princeton, NJ, USA, 2007.
32. Ammerman, A.J.; Cavalli-Sforza, L.L. *The Neolithic Transition and the Genetics of Populations in Europe*; Princeton University Press: Princeton, NJ, USA, 1984.
33. Renfrew, C. *Archaeology and Language: The Puzzle of Indo-European Origins*; Jonathan Cape: London, UK, 1987.
34. Heggarty, P. Indo-European and the Ancient DNA Revolution. In Proceedings of the Workshop on Indo-European Origins Held at the Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany, 2–3 December 2013.
35. Haak, W.; Lazaridis, I.; Patterson, N.; Rohland, N.; Mallick, S.; Llamas, B.; Brandt, G.; Nordenfelt, S.; Harney, E.; Stewardson, K.; et al. Massive migration from the steppe was a source for Indo-European languages in Europe. *Nature* **2015**, *522*, 207–211. [[CrossRef](#)]
36. Allentoft, M.E.; Sikora, M.; Sjögren, K.-G.; Rasmussen, S.; Rasmussen, M.; Stenderup, J.; Damgaard, P.B.; Schroeder, H.; Ahlström, T.; Vinner, L.; et al. Population genomics of Bronze Age Eurasia. *Nature* **2015**, *522*, 167–172. [[CrossRef](#)]
37. Narasimhan, V.M.; Patterson, N.; Moorjani, P.; Rohland, N.; Bernardos, R.; Mallick, S.; Lazaridis, I.; Nakatsuka, N.; Olalde, I.; Lipson, M.; et al. The formation of human populations in South and Central Asia. *Science* **2019**, *365*, eaat7487. [[CrossRef](#)]
38. De Barros Damgaard, P.; Martiniano, R.; Kamm, J.; Moreno-Mayar, J.V.; Kroonen, G.; Peyrot, M.; Barjamovic, G.; Rasmussen, S.; Zacho, C.; Baimukhanov, N.; et al. The first horse herders and the impact of early Bronze Age steppe expansions into Asia. *Science* **2018**, *360*, eaar7711. [[CrossRef](#)]
39. Janhunen, J. Indo-Uralic and Ural-Altaic: On the diachronic implications of areal typology. In *Early Contacts between Uralic and Indo-European: Linguistic and Archaeological Considerations*; Carpelan, C., Parpola, A., Petteri, K., Eds.; Suomalais-Ugrilainen Seura: Helsinki, Finland, 2001; pp. 207–220.
40. Pagani, L.; Lawson, J.; Jagoda, E.; Mörseburg, A.; Clemente, F.; Hudjashov, G.; DeGiorgio, M.; Eriksson, A.; Saag, L.; Wall, J.; et al. Genomic analyses inform on migration events during the peopling of Eurasia. *Nature* **2016**. [[CrossRef](#)]
41. Mathieson, I.; Lazaridis, I.; Rohland, N.; Mallick, S.; Patterson, N.; Roodenberg, S.A.; Harney, E.; Stewardson, K.; Fernandes, D.; Novak, M.; et al. Genome-wide patterns of selection in 230 ancient Eurasians. *Nature* **2015**, *528*, 499–503. [[CrossRef](#)]
42. Delaneau, O.; Ongen, H.; Brown, A.A.; Fort, A.; Panousis, N.I.; Dermitzakis, E.T. A complete tool set for molecular QTL discovery and analysis. *Nat. Commun.* **2017**, *8*, 15452. [[CrossRef](#)]
43. Benazzo, A.; Panziera, A.; Bertorelle, G. 4P: Fast computing of population genetics statistics from large DNA polymorphism panels. *Ecol. Evol.* **2015**, *5*, 172–175. [[CrossRef](#)]



44. Guardiano, C.; Longobardi, G. Parameter theory and parametric comparison. In *The Oxford Handbook of Universal Grammar*; Roberts, I., Ed.; Oxford University Press: Oxford, UK, 2017; pp. 377–398.
45. Guardiano, C.; Longobardi, G.; Cordon, G.; Crisma, P. Formal syntax as a phylogenetic method. In *The Handbook of Historical Linguistics II*; Janda, R.D., Joseph, B.D., Vance, B.S., Eds.; John Wiley & Sons, Inc.: New York, NY, USA, 2020; pp. 145–182.
46. Crisma, P.; Guardiano, C.; Longobardi, G. Syntactic parameters and language learnability. *Stud. Saggi Linguist.* **2020**, *58*, 99–130.
47. Hammer, O.; Harper, D.; Ryan, P. PAST: Paleontological Statistics Software Package for Education and Data Analysis. *Palaeontol. Electron.* **2001**, *4*, 1–9.
48. Felsenstein, J.; Felsenstein, J. *Inferring Phylogenies*; Sinauer Associates: Sunderland, MA, USA, 2004.
49. Maddison, W.P.; Maddison, D.R. Mesquite: A Modular System for Evolutionary Analysis; Version 1. 2004. Available online: <http://www.mesquiteproject.org> (accessed on 11 December 2020).
50. Bouckaert, R.; Vaughan, T.G.; Barido-Sottani, J.; Duchêne, S.; Fourment, M.; Gavryushkina, A.; Heled, J.; Jones, G.; Kühnert, D.; De Maio, N.; et al. BEAST 2.5: An advanced software platform for Bayesian evolutionary analysis. *PLoS Comput. Biol.* **2019**, *15*, e1006650. [[CrossRef](#)]
51. Lawson, D.J.; Hellenthal, G.; Myers, S.; Falush, D. Inference of population structure using dense haplotype data. *PLoS Genet.* **2012**, *8*, 11–17. [[CrossRef](#)]
52. Salminen, T. Problems in the taxonomy of the Uralic languages in the light of modern comparative studies. In *Лингвистический беспредел: Сборник Статей к 70-летию АИ Кузнецовой*; Издательство Московского университета: Moscow, Russia, 2002; pp. 44–55.
53. Pimenoff, V.N.; Comas, D.; Palo, J.U.; Vershubsky, G.; Kozlov, A.; Sajantila, A. Northwest Siberian Khanty and Mansi in the junction of West and East Eurasian gene pools as revealed by uniparental markers. *Eur. J. Hum. Genet.* **2008**, *16*, 1254–1264. [[CrossRef](#)]
54. Lazaridis, I.; Nadel, D.; Rollefson, G.; Merrett, D.C.; Rohland, N.; Mallick, S.; Fernandes, D.; Novak, M.; Gamarra, B.; Sirak, K.; et al. Genomic insights into the origin of farming in the ancient Near East. *Nature* **2016**, *536*, 419–424. [[CrossRef](#)]
55. Tambets, K.; Yunusbayev, B.; Hudjashov, G.; Ilumäe, A.M.; Rootsi, S.; Honkola, T.; Vesakoski, O.; Atkinson, Q.; Skoglund, P.; Kushniarevich, A.; et al. Genes reveal traces of common recent demographic history for most of the Uralic-speaking populations. *Genome Biol.* **2018**, *19*, 139. [[CrossRef](#)]
56. Jeong, C.; Balanovsky, O.; Lukianova, E.; Kahbatkyzy, N.; Flegontov, P.; Zaporozhchenko, V.; Immel, A.; Wang, C.C.; Ixan, O.; Khussainova, E.; et al. The genetic history of admixture across inner Eurasia. *Nat. Ecol. Evol.* **2019**, *3*, 966–976. [[CrossRef](#)]
57. Saag, L.; Laneman, M.; Varul, L.; Malve, M.; Valk, H.; Razzak, M.A.; Shirobokov, I.G.; Khartanovich, V.I.; Mikhaylova, E.R.; Kushniarevich, A.; et al. The Arrival of Siberian Ancestry Connecting the Eastern Baltic to Uralic Speakers further East. *Curr. Biol.* **2019**, *29*, 1701–1711.e16. [[CrossRef](#)]
58. Lamnidis, T.C.; Majander, K.; Jeong, C.; Salmela, E.; Wessman, A.; Moiseyev, V.; Khartanovich, V.; Balanovsky, O.; Ongyerth, M.; Weihmann, A.; et al. Ancient Fennoscandian genomes reveal origin and spread of Siberian ancestry in Europe. *Nat. Commun.* **2018**, *9*, 1–12. [[CrossRef](#)]
59. Lehtinen, J.; Honkola, T.; Korhonen, K.; Syrjänen, K.; Wahlberg, N.; Vesakoski, O. Behind Family Trees: Secondary Connections in Uralic Language Networks. *Lang. Dyn. Chang.* **2014**, *4*, 189–221. [[CrossRef](#)]
60. Honkola, T.; Vesakoski, O.; Korhonen, K.; Lehtinen, J.; Syrjänen, K.; Wahlberg, N. Cultural and climatic changes shape the evolutionary history of the Uralic languages. *J. Evol. Biol.* **2013**, *26*, 1244–1253. [[CrossRef](#)]
61. Janhunen, J. Proto-Uralic: What, where, and when? In *The Quasiquicentennial of the Finno-Ugrian Society*; Suomalais-Ugrilainen Seura: Helsinki, Finland, 2009.
62. Kallio, P. Suomen kantakielen absoluuttista kronologiaa. *Virittäjä* **2006**, *110*, 2–25.
63. Jones, E.R.; Zarina, G.; Moiseyev, V.; Lightfoot, E.; Nigst, P.R.; Manica, A.; Pinhasi, R.; Bradley, D.G. The Neolithic Transition in the Baltic Was Not Driven by Admixture with Early European Farmers. *Curr. Biol.* **2017**, *27*, 576–582. [[CrossRef](#)]
64. Kivikoski, E. *Suomen Historia I: Suomen Esihistoria*; Werner-Söderström Oy: Porvoo, Finland, 1961.
65. Miettinen, T. Suomenlahden ulkosaarten esihistoriaa. In *Suomenlahden ulkosaaret: Lavansaari, Seiskari, Suursaari, Tytärsaari*; Hamari, R., Korhonen, M., Timo, M., Talve, I., Eds.; Suomalaisen Kirjallisuuden Seura: Helsinki, Finland, 1996.

66. Palo, J.U.; Ulmanen, I.; Lukka, M.; Ellonen, P.; Sajantila, A. Genetic markers and population history: Finland revisited. *Eur. J. Hum. Genet.* **2009**, *17*, 1336–1346. [[CrossRef](#)]
67. Csányi, B.; Bogácsi-Szabó, E.; Tömöry, G.; Czibula, Á.; Priskin, K.; Csösz, A.; Mende, B.; Langó, P.; Csete, K.; Zsolnai, A.; et al. Y-Chromosome Analysis of Ancient Hungarian and Two Modern Hungarian-Speaking Populations from the Carpathian Basin. *Ann. Hum. Genet.* **2008**, *72*, 519–534. [[CrossRef](#)]
68. Cavalli-Sforza, L.L. *Genes, Peoples, and Languages*; University of California Press: Berkeley, CA, USA; Los Angeles, CA, USA, 1997.
69. Neparáczi, E.; Kocsy, K.; Tóth, G.E.; Maróti, Z.; Kalmár, T.; Bihari, P.; Nagy, I.; Pálfi, G.; Molnár, E.; Raskó, I.; et al. Revising mtDNA haplotypes of the ancient Hungarian conquerors with next generation sequencing. *PLoS ONE* **2017**, *12*, e0174886. [[CrossRef](#)]
70. Bouckaert, R.; Lemey, P.; Dunn, M.; Greenhill, S.J.; Alekseyenko, A.V.; Drummond, A.J.; Gray, R.D.; Suchard, M.A.; Atkinson, Q.D. Mapping the origins and expansion of the Indo-European language family. *Science* **2012**, *337*, 957–960. [[CrossRef](#)]
71. Menozzi, P.; Piazza, A.; Cavalli-Sforza, L. Synthetic Maps of Human Gene Frequencies in Europeans. *Science* **1978**, *201*, 786–792. [[CrossRef](#)]
72. Sokal, R.R.; Oden, N.L.; Legendre, P.; Fortin, M.-J.; Kim, J.; Thomson, B.A.; Vaudor, A.; Harding, R.M.; Barbujani, G. Genetics and Language in European Populations. *Am. Nat.* **1990**, *135*, 157–175. [[CrossRef](#)]
73. Olalde, I.; Mallick, S.; Patterson, N.; Rohland, N.; Villalba-Mouco, V.; Silva, M.; Dulas, K.; Edwards, C.J.; Gandini, F.; Pala, M.; et al. The genomic history of the Iberian Peninsula over the past 8000 years. *Science* **2019**, *363*, 1230–1234. [[CrossRef](#)]
74. Chikhi, L.; Nichols, R.A.; Barbujani, G.; Beaumont, M.A. Y genetic data support the Neolithic demic diffusion model. *Proc. Natl. Acad. Sci. USA* **2002**, *99*, 10008–10013. [[CrossRef](#)]
75. Ning, C.; Wang, C.-C.; Gao, S.; Yang, Y.; Zhang, X.; Wu, X.; Zhang, F.; Nie, Z.; Tang, Y.; Robbeets, M.; et al. Ancient Genomes Reveal Yamnaya-Related Ancestry and a Potential Source of Indo-European Speakers in Iron Age Tianshan. *Curr. Biol.* **2019**, *29*, 2526–2532.e4. [[CrossRef](#)]
76. Tassi, F.; Vai, S.; Ghirotto, S.; Lari, M.; Modi, A.; Pilli, E.; Brunelli, A.; Susca, R.R.; Budnik, A.; Labuda, D.; et al. Genome diversity in the Neolithic Globular Amphorae culture and the spread of Indo-European languages. *Proc. R. Soc. B Biol. Sci.* **2017**, *284*, 20171540. [[CrossRef](#)]

**Publisher’s Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



© 2020 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).