

# Ancestral contributions to contemporary European complex traits

Davide Marnetto<sup>1,2\*</sup>, Vasili Pankratov<sup>1</sup>, Mayukh Mondal<sup>1</sup>,  
 Francesco Montinaro<sup>1,3</sup>, Katri Pärna<sup>1,4</sup>, Leonardo Vallini<sup>5</sup>,  
 Ludovica Molinaro<sup>1</sup>, Lehti Saag<sup>1</sup>, Liisa Loog<sup>1</sup>, Sara Montagnese<sup>6</sup>,  
 Rodolfo Costa<sup>5,7,8</sup>, Mait Metspalu<sup>1</sup>, Anders Eriksson<sup>1</sup>,  
 Luca Pagani<sup>1,5\*</sup>

August 3, 2021

<sup>1</sup> Institute of Genomics, University of Tartu, Riia 23b, 51010, Tartu, Estonia.

<sup>2</sup> present address: Department of Neurosciences ‘Rita Levi Montalcini’, University of Turin, C.so Massimo d’Azeglio 52, 10126, Torino, Italy.

<sup>3</sup> Department of Biology, University of Bari, Bari, Italy

<sup>4</sup> Department of Epidemiology, University Medical Center Groningen, Hanzeplein 1, 9713 GZ, Groningen, The Netherlands.

<sup>5</sup> Department of Biology, University of Padova, Via Ugo Bassi 58/B, 35131, Padova, Italy.

<sup>6</sup> Department of Medicine, University of Padova, Padova, Italy.

<sup>7</sup> Institute of Neurosciences, National Research Council (CNR), Padova, Italy.

<sup>8</sup> Faculty of Health and Medical Sciences, University of Surrey, Guildford, UK.

\* Corresponding Authors e-mail: [davide.marnetto@unito.it](mailto:davide.marnetto@unito.it) (DM); [lp.lucapagani@gmail.com](mailto:lp.lucapagani@gmail.com) (LP)

## Abstract

The contemporary European genetic makeup formed in the last 8000 years as the combination of three main genetic components: the local Western Hunter-Gatherers, the incoming Neolithic Farmers from Anatolia and the Bronze Age component from the Pontic Steppes. When meeting into the post-Neolithic European environment, the genetic variants accumulated during their three distinct evolutionary histories mixed and came into contact with new environmental challenges.

Here we investigate how this genetic legacy reflects on the complex trait landscape of contemporary European populations, using the Estonian Biobank as a case study.

For the first time we directly connect the phenotypic information available from biobank samples with the genetic similarity to these ancestral groups, both at a genome-wide level and focusing on genomic regions associated with each of the 27 complex traits we investigated. We also found SNPs connected to pigmentation, cholesterol, sleep, diastolic blood pressure, and body mass index (BMI) to show signals of selection following the post Neolithic admixture events. We recapitulate existing knowledge about pigmentation traits, corroborate the connection between Steppe ancestry and height and highlight novel associations. Among others, we report the contribution of Hunter Gatherer ancestry towards high BMI and low blood cholesterol levels.

Our results show that the ancient components that form the contemporary European genome were differentiated enough to contribute ancestry-specific signatures to the phenotypic variability displayed by contemporary individuals in at least 11 out of 27 of the complex traits investigated here.

# 1 Introduction

Since its origins, ancient human genetics showed that the current European genetic landscape formed only recently, in the last 8000 years, as the combination of three main genetic components: 1) the local Western Hunter-Gatherers (WHG), 2) the incoming Neolithic Farmers from the Near East (Anatolia\_N) and 3) the Bronze Age component from the Pontic Steppes, often identified with the Yamnaya culture (Yamnaya)<sup>1-3</sup>. As a result, any modern European population is a combination of at least these three components, in variable proportions depending on its particular genetic history and geographic location. Before their arrival in Europe, the ancestors of these three components evolved in different areas and environments for thousands of years, hence differentiating through neutral genetic drift but also adapting to the different climatic, nutritional and pathogenic conditions. When coming together into the post-Neolithic European environment, the genetic variants accumulated during their three distinct evolutionary histories admixed and came into contact with new environmental challenges.

Previous research efforts have indeed characterized evolutionary events specific to these populations which putatively affected their phenotype and appearance, through the tracking of few highly characterized SNPs<sup>4-6</sup> or polygenic scores<sup>7,8</sup>. Nevertheless, while the first approach is limited in the number of variants analyzed and largely blind with regards to complex polygenic traits, the second builds on population-dependent effect sizes<sup>9,10</sup> estimated in Genome Wide Association Studies (GWAS). Such summary statistics, especially in their genome-wide aggregations, may lead to directional bias and lower predictive accuracy in populations different to the one where the GWAS study was performed<sup>11-15</sup> and have sometimes led to ambiguous results about polygenic adaptation<sup>16-19</sup>.

Here we capitalize on the Estonian Biobank by measuring the relative genetic distance of contemporary individuals to a given ancestry and associating it with their phenotype, thus measuring the influence of these ancient genetic sources on the complex traits distribution of contemporary Europeans. By connecting directly the phenotypic information with the genetic

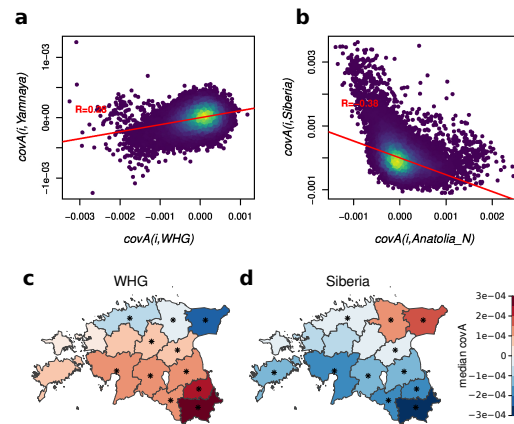
1 similarity to these ancestral groups we avoid the step of producing and interpreting association  
2 summary statistics which might compress information or produce the spurious results men-  
3 tioned above. For the same reason, our conclusions are applicable to contemporary individuals  
4 of European ancestry, where the phenotypes were collected. Conversely, using them to extrap-  
5 olate features of ancient populations, although tempting, should be done with caution due to  
6 the interaction of their genetic legacy with a radically different lifestyle and environment.

7 We started by selecting 27 complex traits of interest, for which we have sufficient data in the  
8 Estonian Biobank, a collection of samples from a relatively homogeneous European popula-  
9 tion which is among the ones with the highest fraction of remnant Hunter Gatherer genomic  
10 component and additionally includes a Siberian (Siberia) component associated with Iron Age  
11 movements<sup>20,21</sup>. In order to associate a phenotype to the contribution of a specific ancient  
12 European ancestry we introduce *covA*, the covariance between allele frequencies in contempo-  
13 rary individuals and a given ancestral population with respect to the contemporary and ancient  
14 average frequencies (see Methods and Supplementary Notes for further details). We computed  
15 *covA* for each pair of Estonian individuals and ancestries, defined models in which each trait is  
16 predicted by *covA* of a specific ancestry and used them to elucidate ancestry/trait associations.  
17 We refined our approach by focusing on the ancestry similarity patterns in genomic regions  
18 potentially connected to each trait according to GWAS catalog<sup>22</sup>. Based only on the SNPs  
19 contained in such regions, we then measured *covA* as above and used it as a predictor to model  
20 traits, also in comparison with random genomic sets with matching size. Finally we set out to  
21 independently analyze if those regions that are associated with the genetic contribution of a  
22 specific ancestry also experienced a post-admixture selective pressure.

## 23 2 Results

### 24 2.1 *covA* measures similarity with ancestral groups

25 We computed *covA* for each pair of Estonian individuals and ancestries among WHG, Anato-  
26 lia\_N, Yamnaya and Siberia using manually curated and other ancient individuals shortlisted  
27 by genetic and chronological proximity (see Methods and Table S1). By observing *covA* joint  
28 distributions in Figure 1a,b (see Figure S1 for all combinations) we can see that as expected,  
29 *covAs* calculated on the various ancestries are strongly interdependent, mainly because they  
30 include as term the average ancestral frequency and partly because of varying grades of simi-  
31 larity among the ancestries for historical demographic reasons. In particular they tend to be  
32 negatively correlated except for *covA* for Yamnaya being associated with *covA* for WHG, re-  
33 flecting complex demographic relationships between the two, due to WHG-like Eastern Hunter  
34 Gatherer ancestry presence in Yamnaya<sup>2,3,23</sup>. Even if by European standards Estonia can be  
35 considered relatively genetically uniform, as recently shown in Pankratov *et al.* [24] the south-  
36 eastern inland counties tend to have higher haplotype sharing with Latvians, Lithuanians and  
37 Russians compared with the rest of the country, and especially the northern coast: this is  
38 reflected by median *covA* for WHG being higher in those Estonian counties, see Figure 1c.  
39 Conversely, as shown by median *covA* for Siberia in Figure 1c, the Siberian ancestry seems to



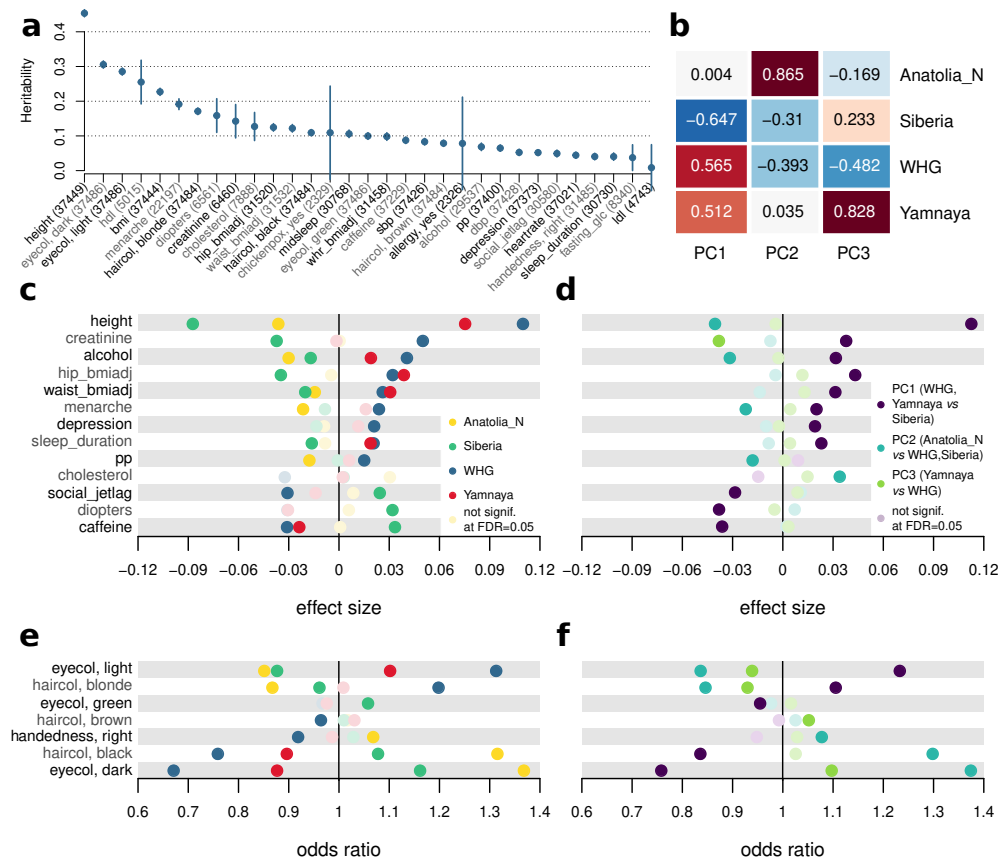
**Figure 1: *covA* distributions.** **a,b** *covA* joint distributions for two ancestry couplings. Each dot is an individual, dots in denser areas are lighter. The red line shows a linear regression, with its R coefficient. **c,d** *covAs* for WHG and Siberia across Estonian counties. Color indicates median *covA* computed in each county, with the sign reflecting excess or lack of a given ancestry, while asterisks indicate those counties for which the *covA* distribution is significantly different than the rest of Estonia (two-tailed Wilcoxon-Mann-Whitney test,  $p \leq 0.001$ )

be more abundant in the north-east, consistently with Finnish ancestry shown in Pankratov *et al.* [24]. Yamnaya and Anatolia\_N *covAs* are instead more evenly distributed (Figure S2).

## 2.2 Connecting complex trait variation to genome-wide ancestry similarity

We examined 27 complex traits (31 if considering pigmentation variants) which were corrected and adjusted for covariates (including sex, age, genotyping platform and others, see Table S2), and expecting varying degrees of influence from genome-wide ancestry depending on their heritability as captured by our dataset(Figure 2a).

As shown above, *covA* exhibits a high correlation across ancestries. Thus we avoided implementing a model with largely multicollinear predictors including *covA* for all ancestries and instead adopted separate models for each ancestry, complementing them with a regression on *covA* PCs (Figure 2b). While *covAs* (Figure 2c,e) highlight the overall excess or lack of certain ancestries in relation with a given phenotype but are largely intertwined, PCs (Figure 2d,f) can be interpreted as independent axes defined by 2 or 3 *covAs* (Figure 2b). Being independent variables in a comprehensive predictive model, they provide a clue to disentangle the potentially collinear *covA* signal and can be reliably used to evaluate significance. When applying this approach, at least one *covA*-based PC had a significant coefficient (coefficient  $p$  value significant at Benjamini-Hochberg FDR=0.05) in the 16 traits shown in Figure 2c-f out of 27 tested. Furthermore, it is also visible how WHG and Yamnaya tend to be linked with the phenotypic



**Figure 2: Genome-wide ancestry-trait associations.** **a** All traits analyzed and their estimated heritability after covariate adjustment. Numbers in parentheses indicate the number of unrelated samples for which phenotypic information was available for each trait. **b** Loading matrix for genome wide *covAs* and their PCs. PCs can be interpreted as axes defined by 2 or 3 *covAs*. **c-f** Genome-wide *covA* or *covA*-based PCs estimated coefficients for traits which have at least one significant PC coefficient.  $\beta$  for continuous **c,d** and Odds Ratios (OR) for categorical **e,f** traits. Pastel dots are deemed not significant at Benjamini-Hochberg FDR = 0.05 (double-sided coefficient *p* value) **c,e**  $\beta$ /ORs of *covA* for a specific ancestry in a model including it together with socioeconomic covariates. Independent models are run for different *covAs*; colors label the probed ancestry. **d,f**  $\beta$ /ORs of first three *covA* PCs in a model including them together with socioeconomic covariates. The legend also describes an interpretation of the PCs.

1 ranges in a similar fashion. As an example, the genomes of taller individuals tend to be more  
2 similar to WHG and Yamnaya, while the opposite is true for Anatolia.N and Siberia. PCs  
3 are largely consistent with this result, with the PC discriminating Yamnaya and WHG lacking  
4 significance.

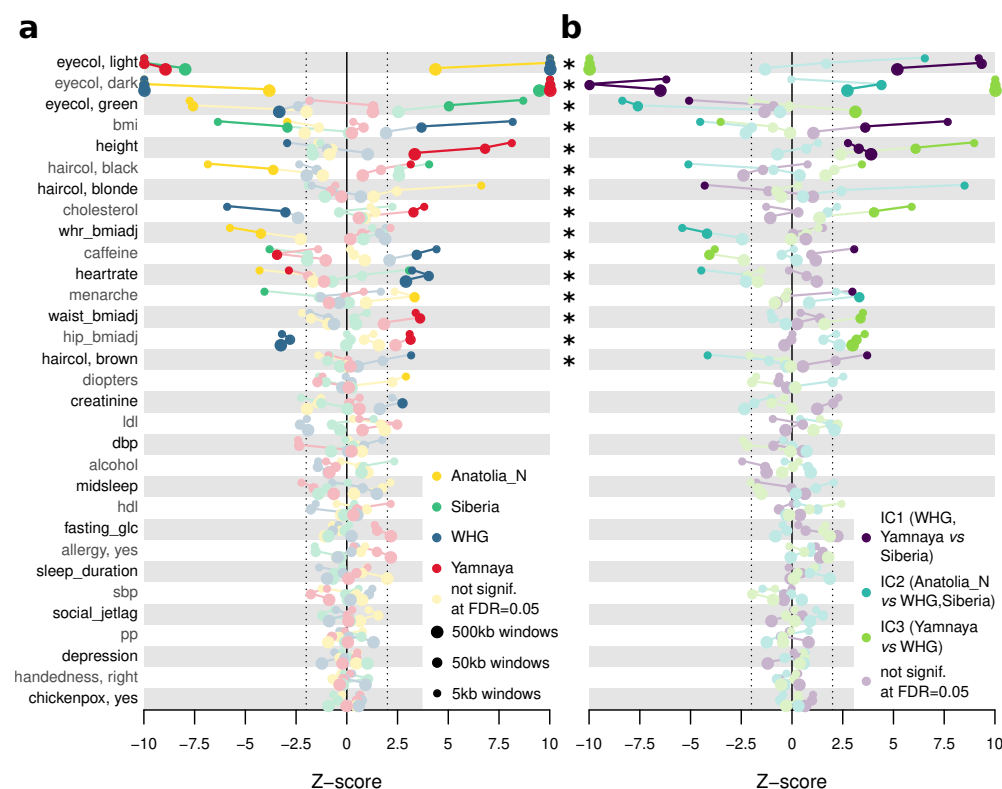
5 As we are not controlling for genotype-based PCs in order not to hide potential genome-wide  
6 signals, we run the risk of obtaining spurious ancestry/trait associations not caused by genet-  
7 ics. This is due to uneven ancestry similarity across Estonia concurrent with geographically  
8 associated socio-economic differences that can influence a trait. Even if such risk is reduced  
9 by the extensive sampling of a relatively uniform population, small differences tied to histor-  
10 ical reasons<sup>24</sup> are still visible in *covA* (see Figure 1c,d, Figure S2 ). Therefore, we include a  
11 city/countryside residency covariate in the models, defined as 1 for people living in Tallinn’s  
12 county (the wealthiest and most populous) and 0 otherwise, and a covariate for educational  
13 attainment, which is a good proxy for family socioeconomic status<sup>25,26</sup>. This control allows us  
14 to suggest a significant influence of genomic ancestry on the 16 traits in Figure 2c-f, even when  
15 geographical and social stratification is present.

## 16 **2.3 Phenotype-associated genomic regions show specific similarity pat-** 17 **terns**

18 To narrow down the signal emerging from the genome-wide analyses we defined three sets  
19 of candidate regions by considering windows of 5kb, 50kb or 500kb centered around GWAS  
20 catalog<sup>22</sup> hits for appropriate categories (see Methods and Table S3). As shown in Figure S3,  
21 these genomic regions harbor a higher heritability intensity ( $h^2/\text{Mb}$ ) than the whole genome,  
22 supporting their appropriateness as candidate regions for the traits of interest.

23 We then asked whether candidate regions for a given trait showed significantly different coef-  
24 ficients when compared to 50 size-matching random genomic sets, and found it true in 11 out  
25 of 27 traits(double-sided Z-test, Benjamini-Hochberg FDR = 0.05), see Z-scores in Figure 3.  
26 This analysis has the advantage of naturally controlling for all potential confounders that apply  
27 to the genome in its entirety, e.g. social, economic and cultural statuses as introduced in the  
28 previous section, thus allowing us to not include any covariates. In addition, this analysis pin-  
29 points genetic signals that are likely to be functionally connected to the trait. Among others,  
30 blood cholesterol levels are shown to be positively correlated with similarity to Yamnaya in  
31 cholesterol-associated regions with respect to the rest of the genome, while the opposite is true  
32 for WHG.

33 Again, to better interpret the signal and avoid multicollinearity, we transformed *covAs* with  
34 the loadings yielded by the PC analysis on whole genome *covAs* (Figure 2b). This, though  
35 not returning actual PCs in each candidate region, drastically reduces the collinearity (highest  
36 Variance Inflation Factor=1.62 in hair color 50kb candidate regions), while allowing simpler  
37 interpretation and, crucially, cross-region comparisons required for Z-scores computation. In-  
38 deed this analysis confirms the significance of the association between cholesterol levels and the  
39 Yamnaya-WHG axis previously mentioned. In contrast to our genome-wide results, candidate



**Figure 3: Ancestry-trait association on candidate regions.** **a** Z-scores of *covA* coefficients, the color refers to the ancestry probed. **b** Z-scores of coefficients associated with *covA* independent components (IC) computed with whole genome-based *covA* PC loadings. Each color is associated with one of the three ICs. For each trait we show the Z-score of the standardized coefficient associated with candidate regions against a distribution of 50 random genomic regions of matching size. Candidate regions are determined around GWAS hits for appropriate traits as windows with three different widths: 5 (small dot), 50 (medium dot) and 500 (large dot) kilobases. Pastel dots are deemed not significant at Benjamini-Hochberg FDR = 0.05,  $p$  value from double-sided Z-test; asterisks mark traits to be considered significant according to **b**; dotted lines correspond to absolute Z-scores = 2.



regions no longer yield concordance between WHG and Yamnaya trends across the traits spectrum, both when considering *covA* and their independent components (IC), suggesting a higher specificity of this refined approach.

## 2.4 Selection signatures at candidate regions with ancestry/trait association

So far we only explored associations between a given trait and a local excess of a given ancestry. The observed local admixture unbalance points to a role of that ancient contribution in explaining a given phenotype. However, these results alone do not show whether after the admixture event the incoming genetic material also underwent a selective sweep within the recipient population, altering population-wide allele frequencies as investigated in Mathieson *et al.* [5]. In other words, the local admixture imbalances we detected so far are not necessarily transferred to the whole population.

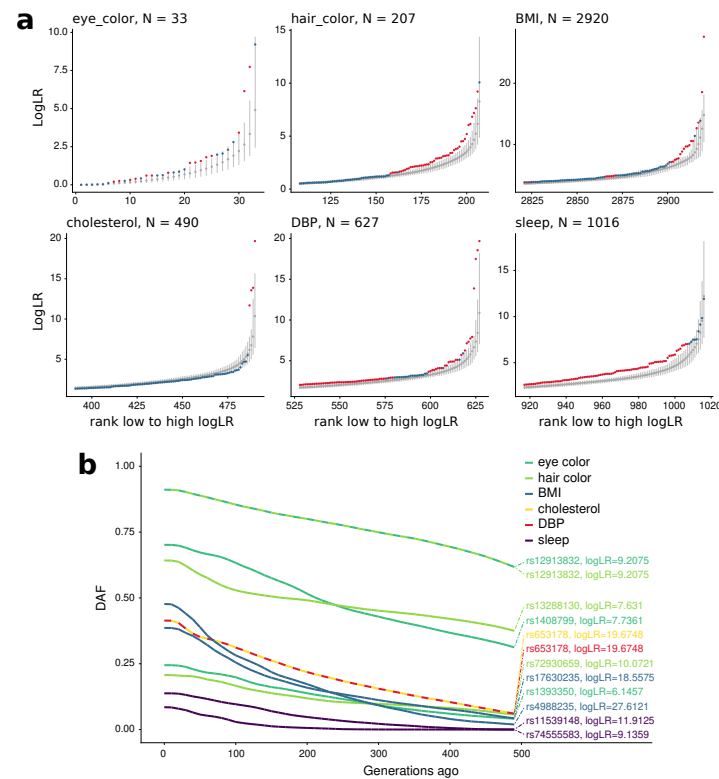
We therefore independently asked whether the phenotypes that showed differential contribution from different ancestries exhibit signs of recent natural selection. We applied CLUES<sup>27</sup> to the list of GWAS hits used above as index for our candidate regions to obtain per-SNP evidence of recent (up to 500 generations ago) natural selection, and to see which phenotypes show enrichment in SNPs with strong selection signals compared to a random set of GWAS hits. Out of the genomic regions responsible for ancestry/trait association shown in Figure 3, pigmentation-related SNPs (eye and hair color) showed extremely high CLUES logLR values (Figures 4a, S4) in accordance with previous results<sup>5,8,28</sup>, as well as SNPs related to BMI and cholesterol, pointing to ongoing or recent selection at these loci. Diastolic blood pressure (DBP) and sleep-related SNPs also showed the same extreme signature, but the candidate regions encompassing them did not reach significance in ancestry/trait association.

The recent and putatively ongoing nature of the inferred selective pressure on the six traits shown in Figure 4a is further exemplified by the steep increase in derived allele frequencies over time inferred for the top 3 SNPs of each trait and shown in Figure 4b. These include some loci previously shown to be selected in West Eurasians (rs4988235 at MCM6/LCT<sup>29</sup>, pigmentation-related SNPs at HERC2/OCA2, TYRP1, TYR, TPCN2<sup>8,28,30</sup>, rs653178 at ATXN2<sup>31</sup>) and some other yet to be explored: rs17630235, rs11539148, rs74555583.

## 3 Discussion

Putting together the genome-wide, region-specific and selection results, the emerging picture points to a different role of each ancestry in having contributed to the phenotype landscape of contemporary Europeans. As a whole, the most affected traits include pigmentation and anthropometric traits together with blood cholesterol levels, caffeine consumption, heart rate and age at menarche.





**Figure 4: Selection signatures.** **a** CLUES log likelihood ratios (logLR) values distribution for GWAS hits for six selected phenotypes. For each phenotype at most 100 top SNPs with highest logLR values and the corresponding ranks from the random GWAS hits distribution are shown. Grey dots show mean values for each rank in the background distribution while the whiskers show the 5-95 percentile range. The logLR values for tested SNPs are shown in red or blue depending on whether the value lies above the 95th percentile of the values from the background distribution with a given rank. Number of tested SNPs for each phenotype are shown in panel titles. **b** Maximum likelihood estimates of derived allele frequency trajectories for top 3 SNPs with highest logLR values for each phenotype. When more than one SNPs come from the same locus, only the top-scoring SNP is shown.

1 In particular WHG ancestry is linked to lower cholesterol levels, higher BMI and putatively con-  
2 tributed light (but not green) eye color to the contemporary Estonian population. Importantly,  
3 these associations stand when carefully considering *covA* ICs and, in addition, loci associated  
4 with these features also appear to have undergone selection in Estonians. Secondly, although  
5 WHG seems to have an association with hip circumference, caffeine consumption and brown  
6 hair pigmentation, these evidences are ambiguous.

7 An enriched Yamnaya ancestry in the pigmentation candidate regions, in contrast with the  
8 genome wide analysis, is linked to dark eye and hair colors, consistently with what inferred  
9 from aDNA data from the Baltic region<sup>6</sup>. This ancestry is also linked to a strong build, with  
10 high stature (in agreement with previous literature<sup>5,7</sup>) and large hip and waist circumferences,  
11 both at genome-wide and region-specific levels, but also high cholesterol concentrations when  
12 focusing on candidate regions. The associations of Yamnaya and WHG ancestries to respec-  
13 tively higher and lower cholesterol levels, together with the clues of selection at loci connected  
14 to cholesterol and BMI, add a critical element to the knowledge of post-neolithic dietary adap-  
15 tation<sup>6,32,33</sup> and might have important health-related implications.

16 Caffeine consumption, although having significant associations, is difficult to connect to a spe-  
17 cific ancestry: Yamnaya ancestry seems to be linked with lower consumption, whereas the  
18 direction of Siberia and WHG associations depends on the genomic regions included in the  
19 analysis.

20 An enriched Anatolia\_N ancestry in the pigmentation candidate regions has implications op-  
21 posite to Yamnaya, again in contrast with the genome-wide signal. This recurring localized  
22 peculiarity of pigmentation loci possibly reflects selection specific to strong GWAS hits as al-  
23 ready seen for skin pigmentation<sup>8</sup>. Notably, Anatolia\_N enrichment in trait-related genomic  
24 regions is connected with a reduced BMI-corrected waist/hip ratio and heart rate. After con-  
25 sidering *covA* ICs, this connection between Anatolia\_N and heart rate seems to be the one  
26 driving the apparent associations of all other ancestries.

27 Lastly, the Siberia ancestry is connected with dark eye and hair pigmentation, but also green  
28 eye color and lower age at menarche. Again, even if this last trait has ambiguous associations  
29 with Anatolia\_N and WHG ancestries, *covA* ICs provide a clue to disentangle their interactions  
30 in favour of a more robust connection with the Siberia ancestry.

31 Some ancestry/trait associations that were not considered significant at a genome-wide level,  
32 are instead discovered when comparing candidate regions to the rest of the genome, possibly due  
33 to the higher sensitivity of this approach. On the other hand, the opposite happens for alcohol  
34 consumption, depression, sleep duration, social jetlag, dieters, pulse pressure, creatinine levels.  
35 This might be due to a misleading or incomplete tagging of the actual functional regions by  
36 the GWAS catalog hits, or to an incomplete correction of socioeconomic and other non-genetic  
37 factors. In case of sleep-connected traits and DBP, the reported signal of recent or ongoing  
38 selection for loci associated to these phenotypes suggests a yet more complex picture.

39 A general caveat about significance levels observed in this study is that as we refrain from  
40 reducing traits by arbitrary choices, even testing multiple alternatives of the same trait, we  
41 expose ourselves to inflated false negatives. Complex traits are often interdependent for biolog-

ical reasons; therefore, when correcting for multiple testing, this risk is intrinsic to this type of analysis. We deemed it best to acknowledge and control it by avoiding overly stringent multiple testing corrections as Bonferroni. In addition, as highly significant traits tend to have higher heritability, it is likely that our analysis might not have enough statistical power for poorly heritable traits.

Taken together, our results show that the ancient components that form the contemporary European landscape were differentiated enough at a functional level to contribute ancestry-specific signatures on the phenotypic variability displayed by contemporary individuals irrespectively to which target population one may examine. In particular, when looking at Estonians, for 11 out of 27 traits surveyed here we could confirm a significant relationship between presence of a given ancestry in genetic regions associated with a given phenotype and how this is expressed by contemporary individuals. While showing that both autochthonous (WHG) and incoming groups contributed genetic material that shapes the phenotype landscape observed today, we also demonstrated that a subset of these loci further underwent positive selection in the last 500 generations. Although not determining whether the selected alleles (and phenotypes) were predominantly contributed by the autochthonous or incoming groups, by connecting genotypic ancestry and complex traits measured in a large dataset, our results reveal both neutral and adaptive consequences of the post-neolithic admixture events on the European phenotype landscape.

## 4 Methods

### 4.1 Sample selection and ancient European grouping

We used 50,353 sequenced or genotyped individuals from the Estonian Biobank<sup>34</sup> as contemporary Estonian sampleset. After removing second-degree relatives ( $\pi\text{-hat} > 0.25$ ) we obtained a subset of 37,952 individuals and used it as a scaffold to perform a PC Analysis (PCA) with Eigensoft-6.1.4. Other individuals were projected on the same PCA space. Outliers identified in this process (with parameters `numoutlieriter: 5` `numoutlierevec: 10` `outliersigmathreshold: 6`) were discarded. Samples that on the first round of genome-wide *covAs* were more distant than 8 Interquartile Ranges (IQR) from the upper or lower quartile against any of the ancestries were also discarded, resulting in 49811 individuals included in our sample set. For each trait of interest we first removed individuals with missing data for traits and covariates and subsequently discarded second-degree relatives.

To define ancestral European groups we started from the Allen Ancient DNA Resource (AADR) V44.3 merged with present-day individuals typed on the Human Origins array (see Data Availability section). From this set we defined a manually curated core set for each ancestral group, then performed a PCA on a space defined by modern Eurasian and North African individuals west of Iran (included), where the ancient samples were projected. We expanded these core sets to other individuals from AADR dataset using multi-dimensional ellipses with diameters

1 equal to 3 core set SDs. We used 4 dimensions: the annotated dating and the first 3 PCs  
2 generated above. With this process we selected 90 WHG, 92 Anatolia\_N, 74 Yamnaya S1.  
3 In addition, from the ones available from the same dataset, we took 7 samples as representa-  
4 tive of the broader Siberian ancestry, assuming any Siberian individual would be equidistant  
5 to the other ancestral European groups: S\_Even-3.DG, S\_Even-1.DG, S\_Even-2.DG, Bur1.SG,  
6 Bur2.SG, Kor1.SG and Kor2.SG.

## 7 4.2 Phenotypes treatment and heritability

8 Continuous traits were treated as specified in Table S2 and regressed against the covariates  
9 according to the same table. Individuals with traits or covariates more distant than 4 IQRs  
10 from the upper or lower quartile were considered as outliers and discarded. The heritability  
11 was computed using LDAK 5.0<sup>35</sup>. First we computed a kinship matrix with the LDAK-Thin  
12 Model: we thinned down SNPs on the non-related sample set defined above with parameters  
13 `--window-prune .98 --window-kb 100`, then used `--calc-kins-direct` with the resulting  
14 weights and `--power .25`. Finally we estimated heritability using REML solver.

## 15 4.3 *covA* definition

*covA* is the covariance in allele frequency ( $p$ ) within a contemporary individual  $i$  (i.e. its allele dosage) with the ancestral group of interest  $j$ , computed respectively against the allele frequency  $p_C$  of the contemporary population  $C$  and the average frequency  $p_A$  in all the  $A$  ancient groups:

$$covA(i, j) = (p_i - p_C)(p_j - p_A) \quad (1)$$

16 *covA* is expected to be high when the allele frequencies of the individual  $i$  and the ances-  
17 try  $j$  are similar in comparison with the differences within the contemporary population and  
18 across the ancestries that contributed to its genetic makeup. *covA* can be computed across the  
19 genome or for specific regions of interest, averaging over the contribution of multiple SNPs. See  
20 Supplementary Notes and Figures S5, S6 for further discussion of *covA* properties.

## 21 4.4 Predicting traits with *covA* and *covA*-based PCs

22 We fitted each standardized trait with a model including one standardized *covA* and, in case  
23 of the genome-wide analysis, socioeconomic covariates as described in the result section. This  
24 analysis was restricted to samples for which socioeconomic covariates were defined, i.e. 38,996  
25 samples (including relatives): the actual sample size for this analysis is therefore less than  
26 reported in Figure 2a and Table S2. The standardized coefficient ( $\beta$  or effect size), or the Odds  
27 Ratio (OR) were used to assess ancestry/trait association for continuous and categorical traits  
28 respectively. In particular, categorical traits were transformed to  $\{0, 1\}$  where 1 stands for  
29 the specified category and 0 for all the others. In addition, each trait was regressed against

1 three *covA*-based PCs, which explained all *covA* variability. PCs were standardized and included  
2 together as predictors, socioeconomic variables were again added as covariates. In the candidate  
3 regions analysis, we adopted exactly the same steps, performing individual regressions for all  
4 the *covAs* and coupling this with a model including all PC-transformed *covAs*. Notably, we  
5 transformed all *covAs* using the loadings of the whole genome *covA* PCs, obtaining components  
6 that were largely independent, yet not strictly principal. Furthermore, to evaluate association  
7 we used coefficient Z-scores computed against the same statistics extracted from 50 random  
8 genomic sets with matching size.

## 9 4.5 Candidate genomic regions

10 We downloaded GWAS hits from GWAS catalog<sup>22</sup> (date of download: 20/11/2020) and then ex-  
11 tracted for each trait a set of hits connected to it filtering on the reported trait ("TRAIT/DISEASE"  
12 field) or selecting the appropriate trait in the Experimental Factor Ontology (EFO) field, as  
13 specified in Table S3. Then we took windows of 5, 50 and 500 Kbs centered on the selected  
14 hits and merged them where overlapping, obtaining three sets of candidate regions for each  
15 trait. To perform the Z-score analysis, for each of them we obtained 50 matching window sets  
16 randomly placed across the genome.

## 17 4.6 Testing for signals of positive selection

18 In order to test individual SNPs for signatures of positive selection we utilized the Relate/CLUES  
19 pipeline<sup>27,36</sup>. This was applied on a curated subset of 1800 unrelated samples; further details  
20 on its application are described in Relate/CLUES Supplementary Methods. CLUES was run  
21 once for each of the 14,712 unique GWAS hits for traits analyzed here with a derived allele  
22 frequency (DAF) above 1% and passing the 1000 Genomes strict mask. To obtain an expected  
23 distribution we randomly sampled 10,000 GWAS hits from the GWAS catalog meeting the same  
24 conditions and ran CLUES for positions not present among the 14,712 SNPs. Next, for each  
25 phenotype we compared its distribution of the logLR values to that of random GWAS hits. We  
26 took 1000 random subsets (with replacement) from the 10,000 logLR values each of the same  
27 length as the number of GWAS hits for a given phenotype and ranked the logLR values from  
28 lowest to highest within each subset. In this way we obtained 1000 values for each logLR rank  
29 from 1 to  $N$  where  $N$  is the number of SNPs analyzed for a given phenotype. For each rank we  
30 calculated the mean and the 5<sup>th</sup> and 95<sup>th</sup> percentiles. Finally, we rank SNPs within each trait  
31 and compare each logLR value to the mean and 5<sup>th</sup> – 95<sup>th</sup> percentiles range for the correspond-  
32 ing rank of the background distribution. As we are interested in deviations in the higher ranks  
33 we focus on the top 100 ranks for each phenotype. Such an approach is conservative as we are  
34 testing not against presumably neutral SNPs but against random GWAS hits that are shown  
35 to be enriched in signals on natural selection compared to random SNPs in the genome<sup>36</sup>.

## 1 Data availability

2 The datasets analyzed during the current study are publicly available and can be accessed  
3 from the following repositories: data from Estonian Biobank at [https://genomics.ut.ee/en/](https://genomics.ut.ee/en/access-biobank)  
4 [access-biobank](https://genomics.ut.ee/en/access-biobank) (accessed with Approval Number 285/T-13 obtained on 17/09/2018 by the  
5 University of Tartu Ethics Committee); AADR plus Human Origins dataset at [https://reich.](https://reich.hms.harvard.edu/allen-ancient-dna-resource-aadr-downloadable-genotypes-present-day-and-ancient-dna-data)  
6 [hms.harvard.edu/allen-ancient-dna-resource-aadr-downloadable-genotypes-present-](https://reich.hms.harvard.edu/allen-ancient-dna-resource-aadr-downloadable-genotypes-present-day-and-ancient-dna-data)  
7 [day-and-ancient-dna-data](https://reich.hms.harvard.edu/allen-ancient-dna-resource-aadr-downloadable-genotypes-present-day-and-ancient-dna-data); GWAS catalog at <https://www.ebi.ac.uk/gwas/>.

## 8 Code Availability

9 Code for analyses performed in this paper will be accessible upon publication.

## 10 References

- 11 1. Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for present-  
12 day Europeans. *Nature* **513**, 409–13 (2014).
- 13 2. Haak, W. *et al.* Massive migration from the steppe was a source for Indo-European lan-  
14 guages in Europe. *Nature* **522**, 207–211 (2015).
- 15 3. Allentoft, M. E. *et al.* Population genomics of Bronze Age Eurasia. *Nature* **522**, 167–172  
16 (2015).
- 17 4. Olalde, I. *et al.* Derived immune and ancestral pigmentation alleles in a 7,000-year-old  
18 Mesolithic European. *Nature* **507**, 225–8 (2014).
- 19 5. Mathieson, I. *et al.* Genome-wide patterns of selection in 230 ancient Eurasians. *Nature*  
20 **528**, 499–503 (2015).
- 21 6. Saag, L. *et al.* Genetic ancestry changes in Stone to Bronze Age transition in the East  
22 European plain. *Science Advances* **7**, eabd6535 (2021).
- 23 7. Cox, S. L., Ruff, C. B., Maier, R. M. & Mathieson, I. Genetic contributions to variation  
24 in human stature in prehistoric Europe. *Proceedings of the National Academy of Sciences*  
25 *of the United States of America* **116**, 21484–21492 (2019).
- 26 8. Ju, D. & Mathieson, I. The evolution of skin pigmentation-associated variation in West  
27 Eurasia. *Proceedings of the National Academy of Sciences of the United States of America*  
28 **118**, e2009227118 (2021).
- 29 9. Visscher, P. M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation.  
30 *The American Journal of Human Genetics* **101**, 5–22 (2017).
- 31 10. Martin, A. R. *et al.* Human Demographic History Impacts Genetic Risk Prediction across  
32 Diverse Populations. *American journal of human genetics* **100**, 635–649 (2017).

- 1 11. Manrai, A. K. *et al.* Genetic Misdiagnoses and the Potential for Health Disparities. *New*  
2 *England Journal of Medicine* **375**, 655–665 (2016).
- 3 12. Kim, M. S., Patel, K. P., Teng, A. K., Berens, A. J. & Lachance, J. Genetic disease risks  
4 can be misestimated across global populations. *Genome Biology* **19**, 179 (2018).
- 5 13. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health  
6 disparities. *Nature Genetics* **51**, 584–591 (2019).
- 7 14. Kerminen, S. *et al.* Geographic Variation and Bias in the Polygenic Scores of Complex  
8 Diseases and Traits in Finland. *American journal of human genetics* **104**, 1169–1181  
9 (2019).
- 10 15. Marnetto, D. *et al.* Ancestry deconvolution and partial polygenic score can improve sus-  
11 ceptibility predictions in recently admixed individuals. *Nature Communications* **11**, 1628  
12 (2020).
- 13 16. Racimo, F., Berg, J. J. & Pickrell, J. K. Detecting Polygenic Adaptation in Admixture  
14 Graphs. *Genetics* **208**, 1565–1584 (2018).
- 15 17. Novembre, J. & Barton, N. H. Tread lightly interpreting polygenic tests of selection.  
16 *Genetics* **208**, 1351–1355 (2018).
- 17 18. Sohail, M. *et al.* Polygenic adaptation on height is overestimated due to uncorrected strat-  
18 ification in genome-wide association studies. *eLife* **8**, 1–17 (2019).
- 19 19. Berg, J. J. *et al.* Reduced signal for polygenic adaptation of height in UK biobank. *eLife*  
20 **8**, 1–47 (2019).
- 21 20. Tambets, K. *et al.* Genes reveal traces of common recent demographic history for most of  
22 the Uralic-speaking populations. *Genome Biology* **19**, 139 (2018).
- 23 21. Saag, L. *et al.* The Arrival of Siberian Ancestry Connecting the Eastern Baltic to Uralic  
24 Speakers further East. *Current Biology* **29**, 1701–1711.e16 (2019).
- 25 22. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association  
26 studies, targeted arrays and summary statistics 2019. *Nucleic Acids Research* **47**, D1005–  
27 D1012 (2019).
- 28 23. Damgaard, P. d. B. *et al.* 137 ancient human genomes from across the Eurasian steppes.  
29 *Nature* **557**, 369–374 (2018).
- 30 24. Pankratov, V. *et al.* Differences in local population history at the finest level: the case of  
31 the Estonian population. *European Journal of Human Genetics* **28**, 1580–1591 (2020).
- 32 25. Liu, H. Genetic architecture of socioeconomic outcomes: Educational attainment, occupa-  
33 tional status, and wealth. *Social Science Research* **82**, 137–147 (2019).
- 34 26. Morris, T. T., Davies, N. M., Hemani, G. & Smith, G. D. Population phenomena inflate  
35 genetic associations of complex social traits. *Science Advances* **6**, eaay0328 (2020).
- 36 27. Stern, A. J., Wilton, P. R. & Nielsen, R. An approximate full-likelihood method for infer-  
37 ring selection and allele frequency trajectories from DNA sequence data. *PLOS Genetics*  
38 **15**, e1008384 (2019).
- 39 28. Key, F. M., Fu, Q., Romagne, F., Lachmann, M. & Andres, A. M. Human adaptation  
40 and population differentiation in the light of ancient genomes. *Nature Communications* **7**,  
41 1–11 (2016).



- 1 29. Bersaglieri, T. *et al.* Genetic signatures of strong recent positive selection at the lactase  
2 gene. *American journal of human genetics* **74**, 1111–20 (2004).
- 3 30. Pickrell, J. K. *et al.* Signals of recent positive selection in a worldwide sample of human  
4 populations. *Genome Research* **19**, 826–837 (2009).
- 5 31. Ding, K. & Kullo, I. J. Geographic differences in allele frequencies of susceptibility SNPs  
6 for cardiovascular disease. *BMC medical genetics* **12**, 55 (2011).
- 7 32. Buckley, M. T. *et al.* Selection in Europeans on Fatty Acid Desaturases Associated with  
8 Dietary Changes. *Molecular Biology and Evolution* **34**, 1307–1318 (2017).
- 9 33. Mathieson, S. & Mathieson, I. FADS1 and the Timing of Human Adaptation to Agricul-  
10 ture. *Molecular Biology and Evolution* **35**, 2957–2970 (2018).
- 11 34. Leitsalu, L. *et al.* *International Journal of Epidemiology* **44**, 1137–1147 (2015).
- 12 35. Speed, D., Holmes, J. & Balding, D. J. Evaluating and improving heritability models using  
13 summary statistics. *Nature Genetics* **52**, 458–462 (2020).
- 14 36. Speidel, L., Forest, M., Shi, S. & Myers, S. R. A method for genome-wide genealogy  
15 estimation for thousands of samples. *Nature Genetics* **51**, 1321–1329 (2019).

## 16 Acknowledgements

17 This work is supported by the European Union through the European Regional Development  
18 Fund, project No. 2014-2020.4.01.16-0024, MOBTT53 (DM, KP, LM, LP); MOBEC008 (VP,  
19 MMo, MMe, AE); 2014-2020.4.01.16-0030 (FM, MMe); 2014-2020.4.01.15-0012 (MMe); through  
20 the Horizon 2020 research and innovation programme grant no. 810645 (VP, MMo, MMe, AE)  
21 and through the Horizon 2020 MSCA Initial Training Network, grant no. 765937 (RC). LS,  
22 MMe are supported by the Estonian Research Council through PUT PRG243. SM is supported  
23 by the STARS@UNIPD 2019 Consolidator Grant for the project CircadianCare.

## 24 Author Contributions

25 DM, LP conceived and designed the study; AE contributed in the statistical design; DM, VP  
26 performed data analyses; MMo, FM, KP, LV, LM, LP contributed to data analyses; SM, RC  
27 provided analyses and expertise about sleep traits; FM, LS, LL, MMe contributed with ancient  
28 genetics expertise; DM, LP drafted the manuscript; all authors reviewed and approved the  
29 submitted paper.

## <sup>1</sup> **Competing Interests**

- <sup>2</sup> The authors declare no competing interests.