

Dutch population structure across space, time and GWAS design

Ross P Byrne^{1,*}, Wouter van Rheenen², Project MinE ALS GWAS Consortium³, Leonard H van den Berg², Jan H Veldink², Russell L McLaughlin^{1,*}

* To whom correspondence should be addressed: Ross P Byrne and Russell L McLaughlin, Complex Trait Genomics Laboratory, Smurfit Institute of Genetics, Trinity College Dublin, Dublin D02 DK07, Republic of Ireland

Email: rbyrne5@tcd.ie; mclaugr@tcd.ie

1 Smurfit Institute of Genetics, Trinity College Dublin, Dublin D02 DK07, Republic of Ireland

2 Department of Neurology and Neurosurgery, Brain Center Rudolf Magnus, University Medical Center Utrecht, Utrecht 3584 CX, The Netherlands

3 A list of Project MinE ALS GWAS Consortium authors and affiliations appears in Supplementary note 1

1 **We studied fine-grained population genetic structure and demographic change across the Netherlands using**
2 **genome-wide single nucleotide polymorphism data (1,626 individuals) with associated geography (1,422**
3 **individuals). We applied ChromoPainter/fineSTRUCTURE, identifying 40 haplotypic clusters exhibiting**
4 **strong north/south variation and fine-scale differentiation within provinces. Clustering is tied to country-wide**
5 **ancestry gradients from neighbouring lands and to locally restricted gene flow across major Dutch rivers.**
6 **Despite superexponential population growth, north-south structure is temporally stable, with west-east**
7 **differentiation more transient, potentially influenced by migrations during the middle ages. Within Dutch**
8 **and international data, GWAS incorporating fine-grained haplotypic covariates are less confounded than**
9 **standard methods.**

10 The Netherlands is a densely populated country on the northwestern edge of the European continent, bounded by
11 Germany, Belgium and the North Sea. The country is divided into twelve provinces and has a complex demographic
12 history, with occupation by several Germanic peoples since the collapse of the Roman Empire, including the
13 Frisians, the Low Saxons and the Franks. Over 17 million individuals now inhabit this relatively small region
14 (41,500km²), making it one of the most densely populated countries in Europe. Despite its small geographical size,
15 previous genetic studies of the people of the Netherlands have demonstrated coarse population structure that
16 correlates with its geography, as well as apparent heterogeneity in effective population sizes across provinces^{1,2}.
17 These observations suggest that the demographic past of the Dutch population has left residual signatures in its
18 present regional genetic structure; however this has not been fully explained in the context of neighbouring
19 populations and thus far the use of unlinked genetic markers have limited the resolution at which this structure can
20 be described. This resolution limit also confines the extent to which the confounding effects of population structure
21 can be controlled in genomic studies of health and disease such as genome-wide association studies (GWAS). As
22 these studies continue to seek ever-rarer genetic variation with ever-increasing cohort sizes, intricate understanding
23 and fine control of population structure is becoming increasingly relevant, but increasingly challenging³.

24 Recent studies have showcased the power of leveraging shared haplotypes to uncover and characterise previously
25 unrecognised fine-grained genetic structure within populations, yielding novel insights into the demographic
26 composition and history of Britain and Ireland⁴⁻⁷, Finland⁸, Japan⁹, Italy¹⁰ and Spain¹¹. Haplotype sharing has also

27 revealed genetic affinities between populations, enabling inference of historical admixture events using modern
28 populations as proxies for ancestral admixing sources¹². Furthermore, geographic information can be integrated to
29 model genetic similarity as a function of spatial distance¹³ to infer demographic mobility within or between
30 populations; one approach uses the Wishart distribution to estimate and map a surface of effective migration rates
31 based on deviations from a pure isolation by distance model¹⁴, allowing migrational cold spots to be inferred which
32 may derive from geographical boundaries such as rivers and mountains. Almost half of the area of the Netherlands
33 is reclaimed from the sea and its contemporary land surface is densely subdivided by human-made waterways and
34 naturally-occurring rivers, including the Rhine (Dutch: *Rijn*), Meuse (*Maas*), Waal and IJssel. These rivers have
35 been speculatively linked to genetic differentiation between northern and southern Dutch subpopulations in previous
36 work¹; however the explicit relationship between Dutch genetic diversity and movement of people within the
37 Netherlands has not been directly modelled.

38 The Dutch have previously received special interest as a model population^{1,2} and form a major component of
39 substantial ongoing efforts to better understand human health, disease, demography and evolution. For example, at
40 the time of writing, over 10% of all studies listed in the NHGRI-EBI genome-wide association study (GWAS)
41 catalogue¹⁵ include the Netherlands in their “Country of recruitment” metadata. As well as offering insights into
42 demography and human history, refined population genetic studies are important to identify and adequately control
43 confounding effects in genomic studies of health and disease, especially if rare variants or spatially structured
44 environmental factors contribute substantially to variance in phenotype¹⁶. In this study, we harness shared
45 haplotypes to examine the fine-grained genetic structure of the Netherlands. We show that Dutch population
46 structure is much stronger than previously recognised, and is ancient and persistent over time. The strength and
47 stability of the observed structure appears to be tied to the relationship of the Netherlands to neighbouring lands and
48 to its own internal geography, and has likely been shaped over history by migration, but preserved in recent
49 generations by enduring sedentism of genetically similar individuals within regions. This has led to genetic structure
50 that demonstrably confounds GWAS; however through analysis of the Netherlands and more extensive international
51 data¹⁷, we show that using shared haplotypes as GWAS covariates significantly reduces this confounding over
52 standard single-marker methods.

53 Results

54 The genetic structure of the Dutch population

55 We mapped the haplotypic coancestry profiles of 1,626 Dutch individuals using ChromoPainter¹⁸ and clustered the
56 resulting matrix using fineSTRUCTURE¹⁸, identifying 40 genetic clusters at the highest level of the hierarchical tree
57 which segregated with geographical provenance. We explored the clustering from the finest (k=40) to the coarsest
58 level (k=2), settling on k=16 as it captured the major regional splits sufficiently with little redundancy. Clusters at
59 this level were robustly defined by total variation distance (TVD) and fixation index (F_{ST} ; Figure 1a); remarkably,
60 some F_{ST} values between Dutch clusters were comparable in magnitude to estimates between European countries
61 (calculated using data from reference 19; Supplementary table 1). Some clusters had expansive geographical ranges
62 (for example NHFG, representing individuals from North Holland, Friesland and Groningen), while others neatly
63 distinguished populations on a sub-provincial level (for example, NBE and NBW, representing east and west regions

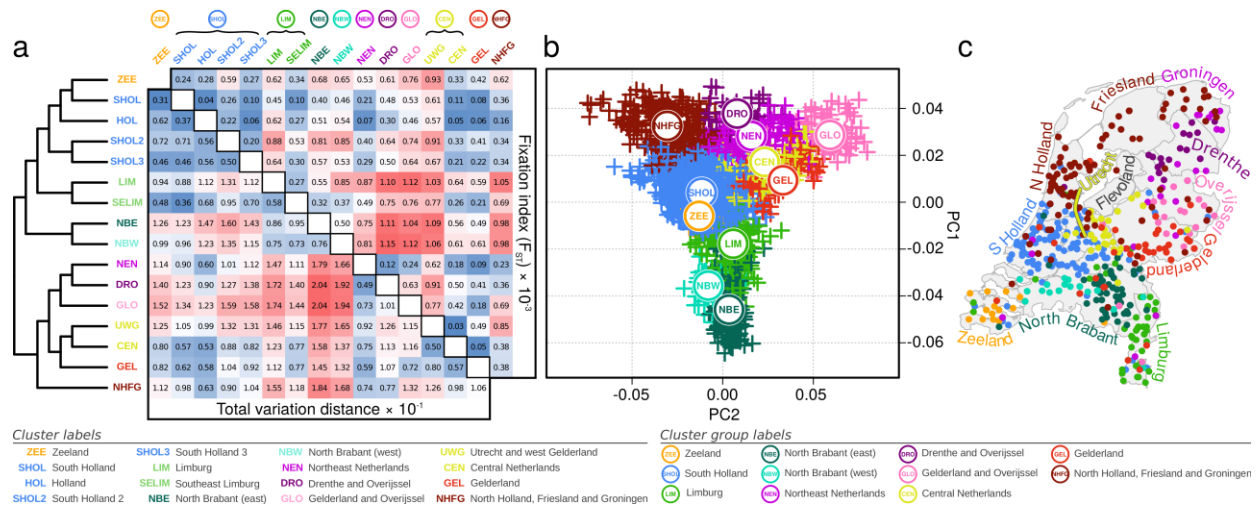


Figure 1 The genetic structure of the people of the Netherlands. (a) fineSTRUCTURE dendrogram of ChromoPainter coancestry matrix showing clustering of 1,626 Dutch individuals based on haplotypic similarity. Associated total variation distance (TVD) and fixation index statistics between clusters are shown in the matrix. Permutation testing of TVD yields $p < 0.001$ for all cluster pairs, indicating that clustering is non-random. Cluster labels derive from Dutch provinces and are arranged into cluster groups for genetically and geographically similar clusters (circled labels). (b) The first two principal components (PCs) of ChromoPainter coancestry matrix for all individuals analysed. Points represent individuals and are coloured and labelled by cluster group. (c) Geographical distribution of 1,422 sampled individuals, coloured by cluster groups defined in (a). Labels represent provinces of the Netherlands. Map boundary data from the Database of Global Administrative Areas (GADM; <https://gadm.org>).

64 of North Brabant). For visualisation we projected the ChromoPainter coancestry matrix in lower dimensional space
 65 using principal component analysis (PCA; Figure 1b) and assigned cluster labels based on majority sampling
 66 location (available for 1,422 individuals), arranging neighbouring and genetically similar clusters into cluster
 67 groups, as with previous work⁶. The first principal component (PC) of coancestry followed a strong north-south
 68 trend (latitude vs mean PC1 per town $r^2 = 0.52$; $p = 6.8 \times 10^{-72}$) with PC2 generally explained by a west-east gradient
 69 (longitude vs mean PC2 per town $r^2 = 0.29$; $p = 3.4 \times 10^{-33}$).

70 As previously observed in different populations⁶, the distribution of individuals in this genetic projection generally
 71 resembled their geographic distribution (Figure 1c), with some exceptions. For example, North Brabant is
 72 geographically further north than Limburg, but is further separated by PC1 from northern clusters. We explored the
 73 possibility that this could instead be explained by relative ancestral affinities to neighbouring lands by modelling the
 74 genome of each Dutch individual as a linear mixture of European sources (obtained from reference 19) using
 75 ChromoPainter, retaining source groups that best matched Dutch individuals for at least 5% of the genome⁴ (Figure
 76 2). The resulting profiles of German, Belgian and Danish ancestries were significantly autocorrelated ($p_{DE}, p_{BE} <$
 77 0.0001 ; $p_{DK} < 0.001$; Moran's I and Mantel's test) and spatially arranged along geographical directions S66°W,
 78 N73°E and S73°E respectively, approximately corresponding to declining ancestry gradients directed away from the
 79 German and Belgian borders and the North Sea boundary (Figure 2; $r_{DE}^2 = 0.31$; $r_{BE}^2 = 0.35$; $r_{DK}^2 = 0.12$; $p_{DE} =$
 80 9.4×10^{-119} ; $p_{BE} = 2.7 \times 10^{-133}$; $p_{DK} = 1.1 \times 10^{-39}$). The spatial distribution of French ancestry was
 81 comparatively uniform, with only a modest correlation due east ($r_{FR}^2 = 0.014$; $p_{FR} = 9.5 \times 10^{-6}$). The general trend
 82 across the Netherlands was thus of complementary Belgian and German ancestral affinities, decaying with distance
 83 from the respective borders. North Brabant, however, showed a greater Belgian profile than Limburg, despite
 84 similar, substantial Belgian frontiers in both Dutch provinces. Conversely, the German ancestry profile of Limburg

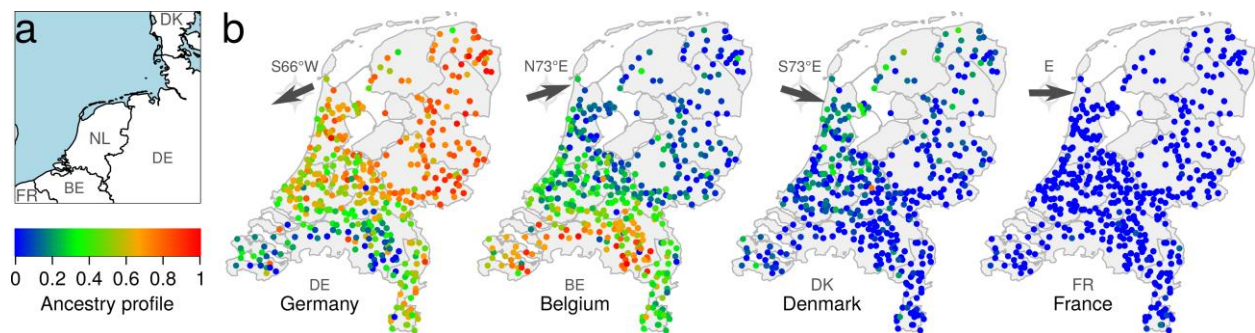


Figure 2 The ancestry profile of the Netherlands. (a) The Netherlands and its geographical relationship to neighbouring lands. (b) German, Belgian, Danish and French haplotypic ancestry profiles for 1,422 Dutch individuals. Arrows indicate the predominant directions along which the ancestry gradients are arranged across the Netherlands. Map boundary data from the Database of Global Administrative Areas (GADM; <https://gadm.org>) and Natural Earth (<https://naturalearthdata.com>).

85 greatly exceeded that of North Brabant, reflecting its 200-kilometre border with Germany and centuries of
86 consequent demographic contact and likely genetic admixture.

87 Genome flux and stasis in the Netherlands

88 To explore temporal trends in Dutch population structure we called genomic segments of pairwise identity-by-
89 descent (IBD) using RefinedIBD²⁰. An IBD haplotype sharing matrix is conceptually similar to a ChromoPainter
90 coancestry matrix²¹, but trades some sensitivity to be more explicitly interpretable. As IBD segment length is
91 inversely related to age^{22,23}, different length intervals can inform on structure at different time depths. Total pairwise
92 IBD between Dutch individuals mirrored the structure observed with ChromoPainter (Figure 3a), with 8 distinct
93 clusters identified in the IBD sharing matrix that broadly segregated with geography and recapitulated some of the
94 important splits obtained from fineSTRUCTURE, most strikingly the west-east split in North Brabant.
95 Decomposing total IBD by centiMorgan (cM) length into short (1-3 cM), medium (3-5 cM) and long (5-7 cM) bins,
96 we observed stability over time of north-south structure and the emergence of west-east structure embedded in 3-5
97 cM segments (Figure 3b), corresponding to an expected time depth around 1,120 years ago²³. As this date and the
98 structure observed is dependent on the (arbitrary) thresholds set for IBD segment length bins, we have also provided
99 an interactive environment in which Dutch population structure can be explored across a range of IBD segment bins
100 (bioinf.gen.tcd.ie/ctg/nlibd).

101 Although these observations could potentially be biased by power to detect population structure in longer and
102 shorter bins, the temporally volatile west-east structure contrasts with the stability and persistence of old north-south
103 structure and possibly represents a genomic signature of historical demographic flux in the region and its
104 surrounding lands. With this in mind, we investigated possible admixture from outside demographic groups using
105 GLOBETROTTER¹² with 4,514 European individuals¹⁹ representing modern proxies for admixing sources. Across
106 the Dutch sample, significant migration dating to 1088 CE (95% c.i. 1004-1111 CE) was inferred with the major
107 contributing source best modelled by modern Germans and the minor source best modelled by southern European
108 groups (France, Spain) (Table 1). This is supported by single-marker ADMIXTURE component estimates showing
109 that the Netherlands has the closest profile to Germanic groups (Supplementary figure 1) and is consistent with the
110 ancestry profile gradients detailed in Figure 2. The timing of the inferred 11th century event was stable across Dutch
111 fineSTRUCTURE clusters (to varying degrees of confidence), suggesting that the signal represents an important

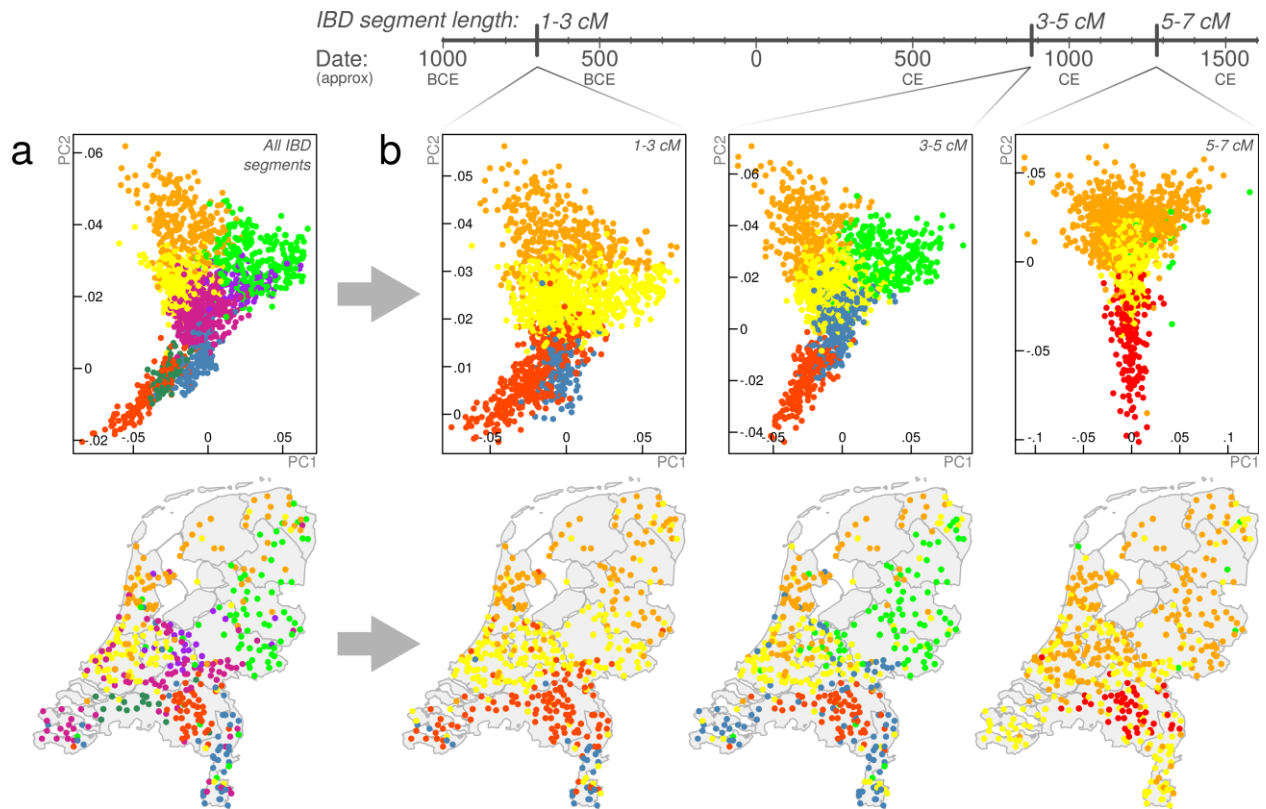












Figure 3 The changing genomic structure of the Dutch population over time. (a) Principal component (PC) analysis of pairwise total identity-by-descent (IBD) for 1,626 Dutch individuals (top) and their geographical provenance (bottom). Points represent individuals and are coloured by cluster assignment (mclust on pairwise IBD matrix). (b) PCs (top) and geographical provenance (bottom) for pairwise sharing of 1-3, 3-5 and 5-7 centiMorgan (cM) IBD segments, corresponding to point estimates of expected time depths at approximately 2,700, 1,120 and 720 years ago, respectively. Time depths for IBD segment bins have wide distributions²³; expected values presented here should be interpreted as a guide only and the changing west-east structure over time does not necessarily reflect (for instance) a precisely-timed admixture event. Map boundary data from the Database of Global Administrative Areas (GADM; <https://gadm.org>).

112 period in the establishment of the modern Dutch genome (Table 1); however, given the state of demographic flux in
 113 Europe at the time, its exact historical correlate is open to interpretation. Notably, a significant admixture event with
 114 a major Danish source was inferred between 759 and 1290 CE in the NHFG cluster group (representing Dutch
 115 northern seaboard provinces); this period spans a historical period of recorded Danish Viking contact and rule in
 116 northern Dutch territories.

117 In addition to influence from outside populations, the population structure detailed in Figure 1 and Figure 3 has
 118 likely been shaped by independent demographic histories within the Netherlands. In support of this, we noted that
 119 short (1-2 cM) IBD segments shared between northern clusters and provinces outnumbered those shared between
 120 southern clusters and provinces (Supplementary figure 2), and, as observed previously², northern provinces shared
 121 more short segments with southern provinces than southern provinces shared amongst themselves. Together, these
 122 results suggest that the north had a smaller ancestral effective population size (N_e) than the south and is probably
 123 derived from an ancient or historical founder event forming the northern population from a subset of southerners.
 124 We formally characterised ancestral trajectories in N_e between the north and the south of the Netherlands using the
 125 nonparametric method IBDNe²⁴ for the entire Dutch sample and two subsamples representing the principal
 126 fineSTRUCTURE north/south split (Figure 4a), retaining a random sample of 641 individuals from each group. We

Table 1 GLOBETROTTER date and source estimates for admixture into the Netherlands.

| Cluster group | Conclusion | Minor | Major | Prop | Date CE | 95% c.i. CE | p |
|---|-------------------|------------|--------|------|---------|-------------|---|
|  | one-date multiway | SPA-FRA(2) | GER(5) | 0.25 | 1169 | 1086-1244 | 0 |
|  | one-date-multiway | FRA(8) | GER(5) | 0.4 | 1172 | 771-1773 | 0 |
|  | one-date-multiway | FRA(8) | GER(5) | 0.4 | 1085 | 939-1262 | 0 |
|  | one-date-multiway | GER(5) | BEL(5) | 0.34 | 1013 | 668-1383 | 0 |
|  | one-date | SPA-FRA(2) | GER(5) | 0.19 | 1172 | 925-1364 | 0 |
|  | one-date-multiway | FRA(8) | GER(5) | 0.16 | 1390 | 1116-1932 | 0 |
|  | one-date | SPA-FRA(2) | GER(5) | 0.14 | 1128 | 893-1306 | 0 |
|  | one-date | SPA-FRA(2) | GER(5) | 0.18 | 1049 | 854-1244 | 0 |
|  | one-date | SPA-FRA(2) | GER(5) | 0.17 | 1189 | 1046-1391 | 0 |
|  | one-date | GER(9) | DEN(5) | 0.36 | 1060 | 759-1290 | 0 |
| ALL | one-date | SPA-FRA(2) | GER(5) | 0.25 | 1088 | 1004-1111 | 0 |

Minor and **Major** represent inferred proxy admixing sources. **Prop** represents estimated minor admixture proportion. Admixing sources are derived from ChromoPainter/fineSTRUCTURE clustering of 4,514 European reference individuals (Methods); labels represent principal country of origin (SPAin, FRAnce, GERmany, BELgium, DENmark) with cluster numbers arbitrarily assigned within countries.

127 also characterised historical N_e within individual Dutch provinces for which genotypes for more than 40 individuals
 128 were available. Countrywide, N_e has grown superexponentially over the past 50 generations in the Netherlands
 129 (Figure 4a) and has been consistently lower in the north than the south. Despite this, the pattern of growth in
 130 northern and southern groups was identical, with a steady exponential growth up to around 1650 CE, when a major
 131 uptick in growth rate was observed. This corresponds to a period of substantial economic development in the
 132 Netherlands over the 17th century known to historians as the Dutch Golden Age. Preceding this period, historical N_e
 133 estimates for the entire country and for northern/southern groups showed only a modest response to the Black Death
 134 (*Yersinia pestis* plague pandemic) of the 14th century which claimed up to 60% of Europe's population²⁵.
 135 Conversely, N_e estimation within individual Dutch provinces revealed a much more detectable impact (Figure 4b).

136 Genomic signatures of Dutch mobility

137 We noted that long (>7 cM) IBD segments, which capture recent shared ancestry, were almost always shared within
 138 genetic clusters (and provinces), and rarely between (Supplementary figure 2). This indicates a propensity for
 139 genetically similar individuals (relatives) to remain mutually geographically proximal, suggesting a degree of
 140 sedentism that has likely influenced Dutch population structure over time. It has also previously been argued that
 141 genetic structure in the Netherlands may be partially rooted in geographic obstacles imposed by the country's major
 142 waterways¹ so we explicitly modelled genetic similarity as a function of geographic distance using EEMS¹⁴ to infer
 143 migrational hot and cold spots (Figure 5). The resulting effective migration surface showed several apparent barriers

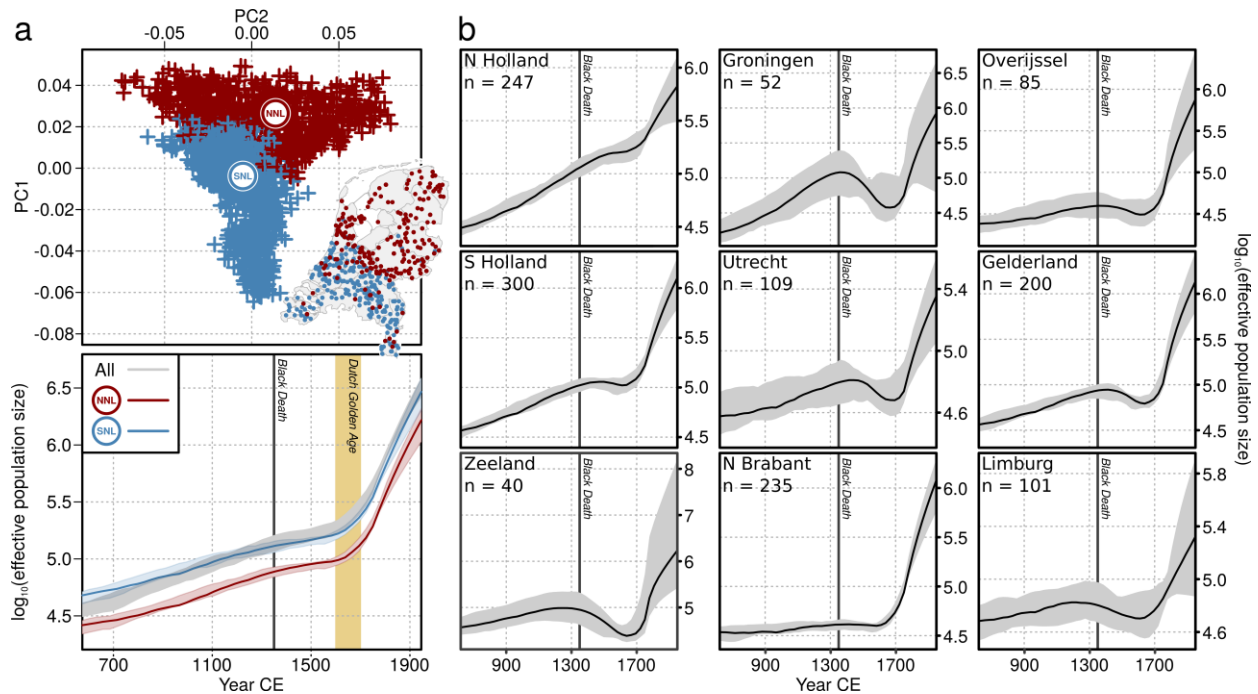


Figure 4 Dutch effective population size over time. (a) Historical change in effective population size (N_e) over the past 50 generations for all Dutch individuals and subsets of northerners and southerners. The top plot shows the principal components of ChromoPainter coancestry coloured by the first ($k=2$) fineSTRUCTURE split, which separates the Dutch population into northern (NNL) and southern (SNL) genetic clusters; inset shows geographical distribution of these individuals. The bottom plot shows growth in effective population size countrywide or per fineSTRUCTURE cluster over the past 50 generations. (b) Historical N_e trajectories for individual Dutch provinces with more than 40 individuals sampled. N_e plots show estimates \pm 95% c.i. and assume 28 years per generation and mean year of birth at 1946 CE. Map boundary data from the Database of Global Administrative Areas (GADM; <https://gadm.org>).

144 to gene flow, the strongest and most contiguous of which runs in an east-west direction across the Netherlands
 145 overlapping the courses of the Rhine, Meuse and Waal rivers. This inferred migrational boundary also
 146 approximately corresponds to the geographical division determining the principal fineSTRUCTURE split between
 147 northern and southern Dutch populations (Figure 4a) as well as the geographical boundaries between clusters
 148 inferred from ancient IBD segments (Figure 3b), suggesting that these rivers have been a historically persistent
 149 determinant of Dutch population structure.

150 GWAS confounding by fine-grained structure

151 As population structure confounds GWAS (for example due to stratification of cases and controls between
 152 subpopulations), we investigated the extent to which haplotype sharing captures confounding structure in a Dutch
 153 sample of 1,971 cases of amyotrophic lateral sclerosis (ALS) and 2,782 controls from a recent multi-population
 154 GWAS for ALS¹⁷. PCs of the haplotypic ChromoPainter coancestry matrix for these 4,753 individuals explained
 155 substantially more variance in ALS phenotype than PCs calculated from SNP genotypes alone, indicating latent
 156 structure captured by ChromoPainter that is stratified between cases and controls (Figure 6a). To estimate the extent
 157 to which this stratified structure confounds GWAS we calculated case-control association statistics using a logistic
 158 model covarying for either 20 ChromoPainter PCs or 20 SNP PCs and estimated the linkage disequilibrium (LD)
 159 score regression intercepts for both sets of resulting summary statistics. An intercept higher than 1 indicates
 160 confounding in the GWAS; Figure 6a shows that GWAS statistics calculated with ChromoPainter PCs as covariates

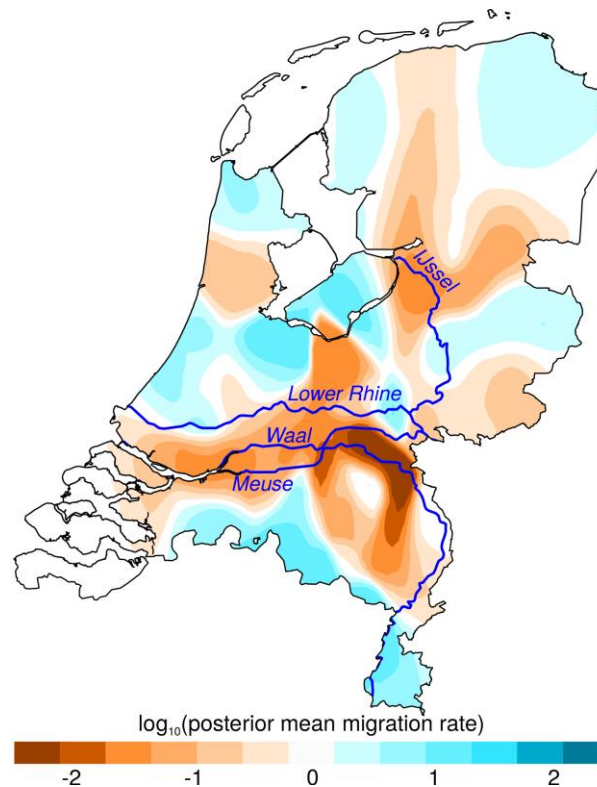


Figure 5 The effective migration surface of the Netherlands. Contour map shows the mean of 10 independent EEMS posterior migration rate estimates between 800 demes modelled over the land surface of the Netherlands. A value of 1 (blue) indicates a tenfold greater migration rate over the average; -1 (orange) indicates tenfold lower migration than average. The courses of major rivers are included to highlight their correlation with migrational cold spots. Map boundary data from the Database of Global Administrative Areas (GADM; <https://gadm.org>); river course data from Natural Earth (<https://www.naturalearthdata.com>).

161 are less confounded than statistics using SNP PCs, albeit with overlapping confidence intervals for the relatively
162 small Dutch sample. To more adequately represent the large-scale multi-population data typically used in modern
163 GWAS, we extended our analysis to the full ALS case-control dataset from which the Dutch data derive¹⁷, including
164 36,052 individuals from twelve European countries and the USA. For computational tractability, instead of
165 ChromoPainter we used PBWT-paint (<https://github.com/richarddurbin/pbwt/blob/master/pbwtPaint.c>), a scalable
166 approximate haplotype painting method based on the positional Burrows-Wheeler transform²⁶. When run on our
167 original Dutch dataset of 1,626 individuals, the structure rendered by PBWT-paint was almost identical to
168 ChromoPainter ($r_{PC1}^2 = 0.99$; $r_{PC2}^2 = 0.98$; Supplementary figure 3), indicating its suitability for this analysis.
169 PBWT-paint captured pervasive global and local structure in the multi-population GWAS data that both separated
170 and subdivided countries (Figure 6b). Top PCs of PBWT-paint coancestry explained substantially more variance in
171 phenotype than SNP PCs and GWAS statistics including PBWT-paint PCs as covariates were significantly less
172 confounded than statistics corrected by SNP PCA (Figure 6a, LD score regression intercepts).

173 Discussion

174 We have studied the Netherlands as a model population, harnessing information from shared haplotypes and recent
175 developments in spatial modelling to gain intricate insights into the geospatial distribution and likely origin of Dutch
176 population genetic structure. The structure identified through shared haplotypes is surprisingly strong; some Dutch

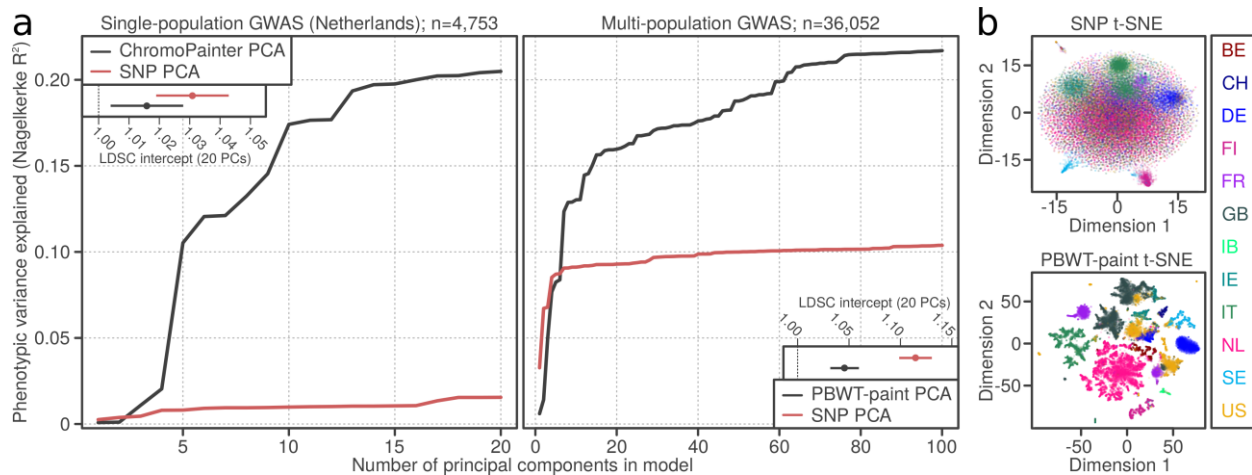


Figure 6 Fine-grained population structure and genome-wide association study (GWAS) confounding. (a) Variance in phenotype (amyotrophic lateral sclerosis) explained by principal components (PCs) for a single-population Dutch GWAS (left) and a multi-population GWAS (right). Insets show linkage disequilibrium score regression (LDSC) intercept terms (a summary estimate of GWAS confounding) when the first 20 single nucleotide polymorphism (SNP)-based PCs (SNP PCA) or the first 20 haplotype-based PCs (ChromoPainter/PBWT-paint PCA) are included as GWAS covariates. (b) Summary visualisations (t-distributed stochastic neighbour embedding, t-SNE) of local and global structure in the multi-population GWAS based on SNP genotypes (top) or haplotype sharing inferred using the scalable PBWT-paint chromosome painting algorithm (bottom). Individuals are coloured by country of origin; labels (right) follow ISO 3166-1 country codes, except IB, which was labelled Iberia (containing Spanish and Portuguese data) in the original GWAS dataset. PCA, principal component analysis; PBWT, positional Burrows-Wheeler transform.

177 genetic clusters identified this way are more mutually distinct (by F_{ST}) than whole European countries. We have also
 178 introduced a novel use of IBD sharing combined with PCA and Gaussian mixture model-based clustering to
 179 characterise changing population structure over time, revealing transient genetic structure layered over strong and
 180 stable north-south differentiation in the Netherlands. This is contextualised by somewhat distinct demographic
 181 histories between genetic groups in the Netherlands, with consistently lower N_e in the north than the south. A
 182 potential source of the north-south differentiation is impaired migration across the east-west courses of the Rhine,
 183 Meuse and Waal, which effectively separate southern Dutch populations from the north. The population structure
 184 observed in the Netherlands is especially remarkable when considered in terms of the country's size and extensive
 185 infrastructure; notably Denmark, which is roughly equal in geographical area, is genetically homogeneous, forming
 186 only a single cluster when interrogated using fineSTRUCTURE²⁷, despite its island-rich geography. Both the United
 187 Kingdom and Ireland also exhibit at least one large indivisible cluster constituting a large fraction of the
 188 population⁴⁻⁶, however no extraordinarily large clusters dominate the Dutch sample. Mean F_{ST} between Dutch
 189 clusters also greatly outmeasures that observed between Irish clusters, suggesting that the extent of population
 190 differentiation is higher in the Netherlands, despite Dutch land area being less than half that of the island of Ireland.

191 While coarse geographical trends in Dutch genetic structure have previously been described using single-marker
 192 PCA¹, our use of shared haplotypes reveals structure at a much higher resolution, differentiating subpopulations
 193 between, and sometimes within, provinces (Figure 1). As a striking example, individuals from the east and west of
 194 North Brabant (NBE and NBW in Figure 1) are mutually genetically distinguishable and are more distinct from
 195 clusters to their north than Limburg, despite being geographically closer. This deviation from haplotype sharing
 196 mirroring geography appears to be driven by strong genetic affinity to Belgium (Figure 2), reflecting a long history
 197 of demographic and sovereign overlap across a 100 km frontier spanning the modern Dutch-Belgian border. In

198 contrast, the majority of ancestral influence in Limburg, which also shares a substantial border with Belgium, is
199 equally split between Belgium to the west and Germany to the east. Notably, the Belgian border with the south of
200 Dutch Limburg is almost entirely described by the course of the Meuse, which may have acted as a historical
201 impediment to migration, thus distinguishing individuals in this region genetically. This is reflected in IBD
202 clustering, in particular the distinction of southern Limburgish individuals from the rest of the Netherlands in short
203 (1-3 cM) segments, which otherwise only describe coarse north-central-south structure (Figure 3). Future work
204 explicitly modelling Dutch-Belgian and Dutch-German frontiers using additional Belgian and German genetic data
205 with associated geography will resolve the historical and present-day role of the Meuse in distinguishing distinct
206 population clusters in the south of the Netherlands.

207 Similarly to North Brabant, groups of individuals in North and South Holland show significant genetic separation
208 despite mutual geographic proximity. While we have chosen to group the four South Holland clusters for visual
209 brevity in Figure 1, they are robustly distinct by TVD analysis (Figure 1a), indicating that significant population
210 differentiation exists even within South Holland. Migration and admixture in the highly urbanised *Randstad* has
211 been proposed as a driver of genetic diversity and loss of geographic structure in this region¹; the overlaid
212 geographical distribution of regional ancestry profiles (Figure 2) for this area lends support to this hypothesis.
213 However, the geographical ranges of the four South Holland clusters are somewhat independent (Supplementary
214 figure 4), indicating that some degree of genetic structure has survived this urbanisation. Previous studies have
215 highlighted the correlation between decreasing autozygosity and increased urbanisation²⁸; future work leveraging
216 the ChromoPainter/fineSTRUCTURE framework coupled with length-binned IBD and Gaussian mixture model-
217 based clustering will more explicitly delineate the interplay between urbanisation and population structure over time.
218 To this end, highly urbanised areas such as the *Randstad* will be particularly informative.

219 The principal fineSTRUCTURE split in the Netherlands describes north-south genetic differentiation (Figure 1) that
220 is strong and persistent over time (Figure 3). We hypothesised that this reflects partially independent demographic
221 histories so we estimated ancestral N_e for northern (NNL) and southern (SNL) Dutch fineSTRUCTURE
222 populations, revealing superexponential growth in both populations with a sudden increase in rate following the 17th
223 century Dutch Golden Age (Figure 4a). Historical N_e follows the same approximate trajectory for both populations
224 but is consistently lower for the northern cluster, corroborating previous observations of increased homozygosity in
225 northern Dutch populations¹ and consistent with a model of northerners representing a founder isolate from
226 southerners (although a more complex demographic model may better explain these observations)^{1,2}. The apparent
227 absence of N_e decline in 14th-century Netherlands initially hints at the possibility that the Black Death had a weaker
228 impact in the region than elsewhere in Europe; although this agrees with the views of some historians, it is hotly
229 debated by others²⁹. Per province, however, most N_e estimates display a prominent dip at this time (Figure 4b),
230 suggesting that merging non-randomly mating subpopulations into a countrywide group (Figure 4a) artificially
231 inflates diversity, thus smoothing over any population crash following the Black Death. Population structure is thus
232 important when estimating N_e and trends countrywide and in NNL and SNL clusters (Figure 4a) should be
233 interpreted carefully: it is possible that a substantial population crash brought about by the Black Death might have
234 had only a marginal impact on the overall effective size of the breeding population in these merged groups. Indeed,
235 the rate of exponential growth in countrywide N_e (Figure 4a) is marginally shallower in the 10 generations following

236 the Black Death (0.024; 95% c.i. 0.0235-0.0251) compared to the 10 generations prior (0.017; 95% c.i. 0.016-
237 0.018), indicating enduring strain on the overall Dutch population prior to its recovery in the Dutch Golden Age.

238 Previous works have hinted that north-south genetic differentiation in the Netherlands may have been facilitated by
239 cultural division between the predominantly Catholic south and the Protestant north¹. Given that the north-south
240 structure observed in 1-3 cM IBD bins (expected time depth ~700 BCE) greatly precedes different forms of
241 Christianity (Figure 3), our data support a model in which the Protestant Reformation of the 16th and 17th centuries
242 exploited pre-existing demographic subdivisions, leading to correlation between distinct cultural affinities and
243 clusters of genetic similarity. Geographical modelling supports the role of migrational boundaries in establishing and
244 maintaining this population substructure, especially rivers (Figure 5). A substantial belt of low inferred migration
245 runs across the Netherlands, corresponding closely to the roughly parallel east-west courses of the Lower Rhine,
246 Waal and Meuse rivers and correlating with the geographical boundary of the principal north-south
247 fineSTRUCTURE split. Absolute assignment of causality to these geographical correlates is, however, not possible
248 and, given the dense network of waterways in the Netherlands, could be misleading. For example, a strong
249 migrational cold spot in the east of the Netherlands runs parallel to the IJssel (Figure 5), but could potentially be
250 better explained by the course of the Apeldoorn Canal, a politically fraught waterway constructed in the early 19th
251 Century. Similarly, a cold spot in the northwest directly overlays the North Sea Canal (completed in 1876). As both
252 of these are human-made waterways, it is not certain whether their courses are consequences or determinants of low
253 movement of people across their paths.

254 As well as internal geography, outside populations have also played an important and significant role in the
255 establishment of population structure in the Netherlands (Figure 2; Table 1); however the variety and extent of
256 demographic upheaval and mobility of European populations over history obscure the likely historical provenance
257 of many inferred admixture signals. As an important exception, however, ancestry profiles show a small but
258 significant contribution of Danish haplotypes in the north and west of the Netherlands, a possible vestige of Viking
259 raids in coastal areas in the 9th and 10th centuries. This is corroborated by an inferred GLOBETROTTER single-date
260 admixture event in the NHFG (North Holland, Friesland and Groningen) cluster (Figure 1) between 759 and 1290
261 CE with Danish haplotypes as a major admixing source (Table 1). The demographic legacy of more than a century
262 of Danish Viking raids and settlement in the Netherlands has been the subject of some debate; from our data, it
263 appears that the modern Dutch genome has indeed been partially shaped by historical Viking admixture. This
264 Danish Viking contact is contemporaneous with a critical period in the establishment of the modern Dutch genome
265 from other outside sources (1004-1111 CE; Table 1), although the precise historical correlates of the admixture
266 events detected in the remaining Dutch regions are less obvious. Future densely sampled ancient DNA datasets from
267 informative time depths in the Netherlands and northwest Europe will enable direct estimation of ancestral
268 population structure, admixture, demographic affinities and effective population sizes, improving precision over the
269 current study which depends on proxy patterns of haplotype sharing between modern individuals. Similarly, regional
270 ancestry and admixture inference are limited by the use of modern proxy populations in place of true ancestral
271 sources; nevertheless, there are ample advantages to the use of modern data, including large sample size and
272 relevance to research on modern human health and disease. In particular, as in our previous work in Ireland⁶,
273 samples in the current Dutch dataset were not specifically selected to have pure ancestry in each geographical area

274 (eg all grandparents from the same region⁴) meaning the degree of structure observed is not idealised or exaggerated
275 by sampling, but instead representative of the structure expected in any GWAS that includes Dutch data.

276 We therefore explored the impact of fine-scale genetic structure described in this study and others⁴⁻¹¹ on GWAS
277 statistics, using the ALS study from which the Dutch data derive as an exemplar trait. Generally, population-based
278 PCs should not predict case/control status; if they do, this indicates that (sub)populations are stratified between cases
279 and controls, introducing bias that artificially inflates GWAS statistics. In both Dutch-only and multi-population
280 analyses, fine-scale genetic structure detected by haplotype sharing (ChromoPainter or PBWT-paint) explained
281 substantially more variance in phenotype (ALS case/control status) than standard SNP-only PCA (Figure 6a). This
282 demonstrates the power of shared haplotypes to simultaneously capture subtle genetic structure within single
283 countries (that is potentially invisible to standard single-marker PCA), broader structure between countries and
284 potential cryptic technical artefacts such as platform- or imputation-derived bias. We found that shared haplotypes
285 are effective for controlling GWAS inflation: statistics calculated using haplotype-based PCs as covariates showed
286 lower overall confounding than single marker-based covariates, as measured by LD score regression intercepts
287 (Figure 6a). In the age of large-scale, single-country and cross-population biobanks, the additional power of
288 haplotype sharing methods to detect fine-scale local population structure will be crucial for ensuring robust GWAS
289 results unconfounded by ancestry. For example, a recent study of latent structure in the UK Biobank demonstrated
290 that a GWAS for birth location returned significant hits even after correction for 40 single-marker PCs³⁰, suggesting
291 that residual fine-grained population structure may influence other GWAS from this cohort. Ongoing developments
292 in scalable haplotype sharing algorithms such as PBWT-paint will help to address this problem by facilitating the
293 creation of biobank-scale haplotype sharing resources, simultaneously improving studies of human health and
294 disease and enabling large-scale, fine-grained population genetic studies of human demography.

295 Methods

296 Data and quality control

297 We mapped fine-grained genetic structure in the Netherlands using a population-based Dutch ALS case-control
298 dataset (n=1,626; subset of stratum sNL3 from a genome-wide association study (GWAS) for amyotrophic lateral
299 sclerosis¹⁷) and a European reference dataset subsampled from a GWAS for multiple sclerosis¹⁹ (MS; n = 4,514;
300 EGA accession ID EGAD00000000120). 1,422 Dutch individuals had associated residential data (hometown at time
301 of sampling) which were used for geographical analyses. For estimating GWAS confounding, we separately
302 analysed the Netherlands on its own using a larger ALS case/control dataset (n = 4,753; strata sNL1, sNL3 and
303 sNL4 from reference 17) and the complete multi-population GWAS dataset¹⁷ (n = 36,052) from which this Dutch
304 subset was derived. Data handling for estimating confounding is further described under “Estimating GWAS
305 confounding” below. For population structure analyses, we applied quality control (QC) using PLINK v1.9³¹; briefly
306 we removed samples with high missingness (>10%), high heterozygosity (>3 median absolute deviations from
307 median) and single-marker PCA outliers (>5 standard deviations from mean for PCs 1-20). We also filtered out A/T
308 and G/C SNPs and SNPs with minor allele frequency <0.05, high missingness (>2%) or in Hardy Weinberg
309 disequilibrium ($p < 1 \times 10^{-6}$). Before running Chromopainter/fineSTRUCTURE we retained only one individual from

310 any pair or group that exhibited greater than 7.5% genomic relatedness ($\hat{\pi}$) and removed SNPs with any missing
311 genotypes as the algorithm does not tolerate missingness or relatedness well. For European reference data we also
312 removed individuals suggested by the QC of the source study¹⁹ and we extracted individuals only of European
313 descent. As this European dataset included MS patients, we filtered out SNPs in a 15 Mb region surrounding the
314 strongly associated HLA locus (GRCh37 position chr6:22,915,594–37,945,593) to avoid bias generated from this
315 association, following previous works. The final Dutch and European reference datasets contained 374,629 SNPs
316 and 363,396 SNPs respectively at zero missingness. The merge of these datasets contained 147,097 SNPs at zero
317 missingness. Data were phased per chromosome with the 1000 Genomes Project phase 3 reference panel³² using
318 SHAPEIT v2³³ (for ChromoPainter/fineSTRUCTURE) and Beagle v4.1 (for IBD estimation). Both programmes
319 were run with default settings; allele concordance was checked prior to phasing (SHAPEIT: --check; Beagle:
320 conform-gt utility).

321 fineSTRUCTURE analysis

322 We used ChromoPainter/fineSTRUCTURE¹⁸ to detect fine-grained population structure using default settings. In
323 brief, each individual was painted using all other individuals (-a 0 0), first estimating N_e and μ (switch rate and
324 mutation rate) with 10 expectation-maximization iterations, then the model was finally run using these parameter
325 estimates. The fineSTRUCTURE Markov chain Monte Carlo (MCMC) model was then run on the resulting
326 coancestry matrix with two chains for 3,000,000 burnin and 1,000,000 sampling iterations, sampling every 10,000
327 iterations. We extracted the state with the maximum posterior probability and performed an additional 10,000 burnin
328 iterations before inferring the final trees using both the climbtree and maximum concordance methods. For all
329 subsequent analyses the maximum concordance tree was used.

330 Cluster robustness

331 To assess the robustness of clustering in the Dutch data we calculated TVD⁴ and F_{ST} . TVD is a distance metric for
332 assessing the distinctness of pairs of clusters, calculated from the ChromoPainter chunklength matrix. TVD is
333 calculated as the sum of the absolute differences between copying vectors for all pairs of clusters, where the copying
334 vector for a given cluster A is a vector of the average lengths of DNA donated to individuals in A by all clusters.
335 Intuitively, the TVD of two clusters reflects distance between those clusters in terms of haplotype sharing amongst
336 all clusters, and is a meaningful method for assessing the effectiveness of fineSTRUCTURE clustering. To assess
337 whether the observed clustering performed better than chance we permuted individuals between cluster pairs
338 (maintaining cluster size) and calculated the number of permutations that exceeded our original TVD score for that
339 pairing of clusters. We used 1,000 permutations where possible, and otherwise used the maximum number of unique
340 permutations. P-values were calculated from the number of permutations greater than or equal to the observed TVD
341 divided by the total permutations; all p-values were less than 0.001, indicating robust clustering. Finally we
342 generated a TVD tree for $k=16$ by merging pairs of clusters with the lowest TVD successively using methods
343 described previously⁸ (Supplementary figure 5). The tree was built in $k-1$ steps, with TVD recalculated at each step
344 from the remaining populations. Branch lengths were scaled proportional to the TVD value of the corresponding
345 pair of populations using adapted code from the original paper. To assess cluster differentiation independently of the

346 ChromoPainter model, F_{ST} was calculated between Dutch clusters using PLINK 1.9. For comparison, we calculated
347 F_{ST} between European countries present in reference 19.

348 Ancestry profiles

349 We assessed the ancestral profile of Dutch samples in terms of a European reference made up of 4,514 European
350 individuals¹⁹ from Belgium, Denmark, Finland, France, Germany, Italy, Norway, Poland, Spain and Sweden.
351 European samples were first assigned to homogeneous genetic clusters using fineSTRUCTURE as in previous
352 work⁶ to reduce noise in painting profiles. We then modelled each Dutch individual's genome as a linear mixture of
353 the European donor groups using ChromoPainter, and applied ancestry profile estimation as described previously⁴
354 and implemented in GLOBETROTTER¹² (num.mixing.iterations: 0). This method estimates the proportion of DNA
355 which is most closely shared with each individual from each donor group calculated from a normalised
356 ChromoPainter chunklength output matrix, and then implements a multiple linear regression of the form

$$357 \quad Y_p = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_g X_g$$

358 to correct for noise caused by similarities between donor populations. Here, Y_p is a vector of the proportion of DNA
359 that individual p copies from each donor group, and X_g is the vector describing the average proportion of DNA that
360 individuals in donor group g copy from other donor groups including their own. The coefficients of this equation
361 $\beta_1 \dots \beta_g$ are thus interpreted as the "cleaned" proportions of the genome that target individual p copies from each
362 donor group, hence the ancestral contribution of each donor group to that individual. The equation is solved using a
363 non-negative-least squares function such that $\beta_g \geq 0$ and the sum of proportions across groups equals 1. We
364 discarded European groups that contributed less than 5% total to any individual, and refit to eliminate noise. We
365 then aggregated sharing proportions across donor groups (genetically homogenous clusters) from the same country
366 to estimate total sharing between an individual and a given country to investigate the regional distribution of sharing
367 profiles. Autocorrelation of ancestry profiles was assessed by Moran's I and Mantel's test (10,000 permutations) in
368 R version 3.2.3. Geographical directions of ancestry gradients were determined by rotating the plane of latitude-
369 longitude between 0° and 360° in 1° steps and finding the axis Y that maximised the coefficient of determination for
370 the linear regression $Y \sim A_c$, where A_c is the aggregated ancestry proportion for country c .

371 Identity-by-descent analyses

372 IBD segments were called in phased data using RefinedIBD²⁰ (default settings) to generate pairwise matrices of total
373 length of IBD shared between individuals for bins of different segment lengths. To identify population structure
374 captured by IBD sharing patterns we performed PCA on these matrices using the prcomp function in R version
375 3.2.3³⁴ and clustered the IBD matrices using a Gaussian mixture model implemented in the R package mclust³⁵. We
376 note that while previous work²¹ has shown that IBD matrices underperform the linked ChromoPainter matrix in
377 identifying population structure, they are arguably more interpretable for visualising temporal change as they can be
378 subdivided into cM bins corresponding to different time periods, a feature leveraged by emerging work on local
379 population structure²³. Patterns in IBD sharing that identify population subgroups in older (shorter) cM bins which
380 are preserved in more recent (longer) bins are interpreted as persistent population structure that has been influenced

381 by mating patterns in old and recent generations. Structure which emerges in a specific cM bin and is lost is likely to
382 reflect transient changes in panmixia that have not necessarily persisted. We approximated the age of segments in a
383 given cM bin using equation s19 from reference 23, under the assumption that the population is sufficiently large.

384 Inferring admixture events

385 To infer and date admixture events from European sources we ran GLOBETROTTER¹² with the Netherlands dataset
386 as a whole and in individual cluster groups defined from the Dutch fineSTRUCTURE maximum concordance tree
387 (Figure 1). To define European donor groups we used the European fineSTRUCTURE maximum concordance tree,
388 as with previous work⁶ to ensure genetically homogenous donor populations. We used ChromoPainter v2 to paint
389 Dutch and European individuals using European clusters as donor groups. This generated a copying matrix
390 (chunklengths file) and 10 painting samples for each Dutch individual. GLOBETROTTER was run for 5 mixing
391 iterations twice: once using the null.ind:1 setting to test for evidence of admixture accounting for unusual linkage
392 disequilibrium (LD) patterns and once using null.ind:0 to finally infer dates and sources. We further ran 100
393 bootstraps for the admixture date and calculated the probability of no admixture as the proportion of nonsensical
394 inferred dates (<1 or >400 generations). Confidence intervals were calculated from the bootstraps from the standard
395 model (null.ind:0) using the empirical bootstrap method, and a generation time of 28 years.

396 ADMIXTURE analysis

397 We performed ADMIXTURE analysis³⁶ on the combined Dutch and European samples to explore single marker-
398 based population structure in a set of 41,675 SNPs (LD-pruned using PLINK 1.9: $r^2 > 0.1$; sliding window 50 SNPs
399 advancing 10 SNPs at a time) SNPs. ADMIXTURE was run for $k=1-10$ populations, using 5 EM iterations at each k
400 value. The k value with the lowest cross validation error was selected for further analysis. We analysed the
401 distribution of proportions for each ADMIXTURE cluster across the Dutch dataset, and its relationship with
402 geography.

403 Estimating mean pairwise IBD sharing within and between groups

404 We compared IBD sharing within and between both clusters and provinces (Supplementary figure 2) using the mean
405 number of segments within a given length range (eg 1-2cM) shared between individuals. To calculate this mean for
406 a single group of size N with itself the denominator was $(N^2 - N)/2$; when comparing two groups of sizes N and M
407 the denominator was NM .

408 Estimating recent changes in population sizes

409 We used IBDNe²⁴ to estimate historical changes in N_e . IBDNe leverages information from the length distribution of
410 IBD segments to accurately estimate effective population size over recent generations, with a resolution limit of
411 about 50 generations for SNP data. We followed the authors' protocol and detected IBD segments using IBDseq
412 version r1206³⁷ with default settings and ran IBDNe on the resulting output with default settings, removing IBD
413 segments shorter than 4cM (minibd=4, the recommended threshold for genotype data). We compared estimated N_e
414 with recorded census size (<https://opendata.cbs.nl/statline/#/CBS/nl/dataset/37296ned/table?ts=1520261958200>) for

415 approximately equivalent dates (starting at 1946 CE for generation 0 and assuming 1 generation is 28 years) and
416 found that for generations 0 - 3 our N_e estimates were approximately $\frac{1}{3}$ of the census population (Supplementary
417 figure 6), which follows expectation if lifespan is $3\times$ the generation time. The slope of the ratios for the three
418 generations is near zero suggesting that our model tracks well with the census population; this is consistent with
419 reported expectation²⁴.

420 Estimating effective migration surfaces

421 To model geographic barriers to geneflow in the Netherlands we ran EEMS¹⁴. This software provides a visualisation
422 of hot and coldspots for geneflow across a habitat using a geocoded genetic dataset. To run EEMS, we generated an
423 average pairwise genetic dissimilarity matrix from our genotype data using the bed2diffs utility provided with the
424 software. We initially ran the EEMS model with 10 randomly initialised MCMC chains for a short run of 100,000
425 burn-in and 200,000 sampling iterations, thinning every 999 iterations, to find a suitable starting point. For these
426 runs we placed the data in 800 demes and used default settings with the following adjustments to the proposal
427 variances: $qEffctProposals2 = 0.000088888888$; $qSeedsProposals2 = 0.7$; $mEffctProposals2 = 0.7$. The resulting
428 chain with the highest log-likelihood was then used as the starting point for a further ten chains for 1,000,000 burn-
429 in iterations and 2,000,000 sampling iterations, thinning every 9,999 iterations. The model was run with the
430 following adjustments to the proposal variances: $qEffctProposals2 = 0.000088888888$; $qSeedsProposals2 = 0.7$;
431 $mEffctProposals2 = 0.7$. We plotted the results of our analysis using the rEEMSplot package in R and modified the
432 resulting vector graphics using Inkscape v0.91 to remove display artefacts caused by non-overlapping polygons.
433 MCMC convergence was assessed by inspecting the log-posterior traces (Supplementary figure 7).

434 Estimating GWAS confounding

435 To examine the contribution of observed fine-grained population structure to GWAS confounding, we estimated
436 how well phenotype could be predicted by principal components of haplotype sharing matrices in a 2016 GWAS for
437 ALS¹⁷, comparing our results to those obtained using standard single marker PCA. We separately analysed
438 1,060,224 zero-missingness Hapmap3 SNPs that passed QC in the original GWAS for Dutch data alone (1,971
439 cases, 2,782 controls) and for the complete multi-population GWAS (12,577 cases, 23,475 controls). Haplotypes for
440 unrelated individuals ($\hat{\pi} < 0.1$) were phased using SHAPEIT v2³³ and painted in terms of one another using
441 ChromoPainter v2¹⁸ for the Dutch dataset (estimating N_e and μ using the weighted average of 10 EM iterations on
442 chromosomes 1, 8, 15 and 20), and PBWT-paint (<https://github.com/richarddurbin/pbwt/blob/master/pbwtPaint.c>)
443 for the considerably larger multi-population GWAS dataset. PBWT-paint is a fast approximate implementation of
444 ChromoPainter suitable for large datasets. PCs of the resulting coancestry matrices were calculated using the
445 fineSTRUCTURE R tools (<https://www.paintmychromosomes.com>). For comparison we also calculated PCs on
446 independent markers from the SNP datasets using Plink v1.9, first removing long range LD regions³⁸
447 ([https://genome.sph.umich.edu/wiki/Regions_of_high_linkage_disequilibrium_\(LD\)](https://genome.sph.umich.edu/wiki/Regions_of_high_linkage_disequilibrium_(LD))) and pruning for LD
448 (--indep-pairwise 500 50 0.8). Variance in ALS phenotype explained by ChromoPainter/PBWT-paint PCs and SNP
449 PCs (Nagelkerke R^2) was estimated using the glm() function and fmsb package³⁹ in R version 3.2.3. To estimate
450 confounding in GWAS inflation, we implemented a logistic regression model GWAS (--logistic) in PLINK v1.9 for
451 each dataset using 20 ChromoPainter/PBWT-paint PCs or 20 SNP PCs as covariates and ran LD score regression⁴⁰

452 on the resulting summary statistics using recommended settings. Structure evident in the PBWT-paint matrix was
453 visualised and contrasted with corresponding SNP data in 2 dimensions using t-distributed stochastic neighbour
454 embedding (t-SNE)⁴¹ implemented in the Rtsne package in R version 3.2.3 (5,000 iterations; perplexity 30; top 100
455 PCs provided as initial dimensions).

456 Acknowledgements

457 This work has been supported by Science Foundation Ireland (17/CDA/4737), the Motor Neurone Disease
458 Association of England, Wales and Northern Ireland (957-799) and the European Research Council (ERC) under the
459 European Union's Horizon 2020 research and innovation programme (grant agreement n° 772376 – EScORIAL).
460 The collaboration project is co-funded by the PPP Allowance made available by Health~Holland, Top Sector Life
461 Sciences & Health, to stimulate public-private partnerships.

462 Author contributions

463 R.P.B. and R.L.McL. conceived the study. R.P.B, W.v.R., J.H.V and R.L.McL. contributed to study design. R.P.B.
464 and R.L.McL. conducted the analyses. R.P.B. and R.L.McL. drafted the manuscript. W.v.R., L.H.v.d.B. and J.H.V.
465 provided data and critical revision of the manuscript.

466 Conflict of interest statement

467 All authors have nothing to declare.

468 **References**

- 469 1. Abdellaoui, A. *et al.* Population structure, migration, and diversifying selection in the Netherlands. *Eur J Hum*
470 *Genet* **21**, 1277–1285 (2013).
- 471 2. Genome of the Netherlands Consortium. Whole-genome sequence variation, population structure and
472 demographic history of the Dutch population. *Nat. Genet.* **46**, 818–825 (2014).
- 473 3. Lawson, D. J. *et al.* Is population structure in the genetic biobank era irrelevant, a challenge, or an opportunity?
474 *Hum. Genet.* (2019) doi:10.1007/s00439-019-02014-8.
- 475 4. Leslie, S. *et al.* The fine-scale genetic structure of the British population. *Nature* **519**, 309–314 (2015).
- 476 5. Gilbert, E. *et al.* The Irish DNA Atlas: Revealing Fine-Scale Population Structure and History within Ireland. *Sci*
477 *Rep* **7**, 17199 (2017).
- 478 6. Byrne, R. P. *et al.* Insular Celtic population structure and genomic footprints of migration. *PLoS Genet.* **14**,
479 e1007152 (2018).
- 480 7. Gilbert, E. *et al.* The genetic landscape of Scotland and the Isles. *Proc Natl Acad Sci U A* **116**, 19064–19070
481 (2019).
- 482 8. Kerminen, S. *et al.* Fine-Scale Genetic Structure in Finland. *G3* **7**, 3459–3468 (2017).
- 483 9. Takeuchi, F. *et al.* The fine-scale genetic structure and evolution of the Japanese population. *PLoS One* **12**,
484 e0185487 (2017).
- 485 10. Raveane, A. *et al.* Population structure of modern-day Italians reveals patterns of ancient and archaic ancestries
486 in Southern Europe. *Sci Adv* **5**, eaaw3492 (2019).
- 487 11. Bycroft, C. *et al.* Patterns of genetic differentiation and the footprints of historical migrations in the Iberian
488 Peninsula. *Nat Commun* **10**, 551 (2019).
- 489 12. Hellenthal, G. *et al.* A genetic atlas of human admixture history. *Science* **343**, 747–751 (2014).
- 490 13. Novembre, J. & Peter, B. M. Recent advances in the study of fine-scale population structure in humans. *Curr*
491 *Opin Genet Dev* **41**, 98–105 (2016).
- 492 14. Petkova, D., Novembre, J. & Stephens, M. Visualizing spatial population structure with estimated effective
493 migration surfaces. *Nat Genet* **48**, 94–100 (2016).
- 494 15. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted
495 arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
- 496 16. Mathieson, I. & McVean, G. Differential confounding of rare and common variants in spatially structured
497 populations. *Nat. Genet.* **44**, 243–246 (2012).

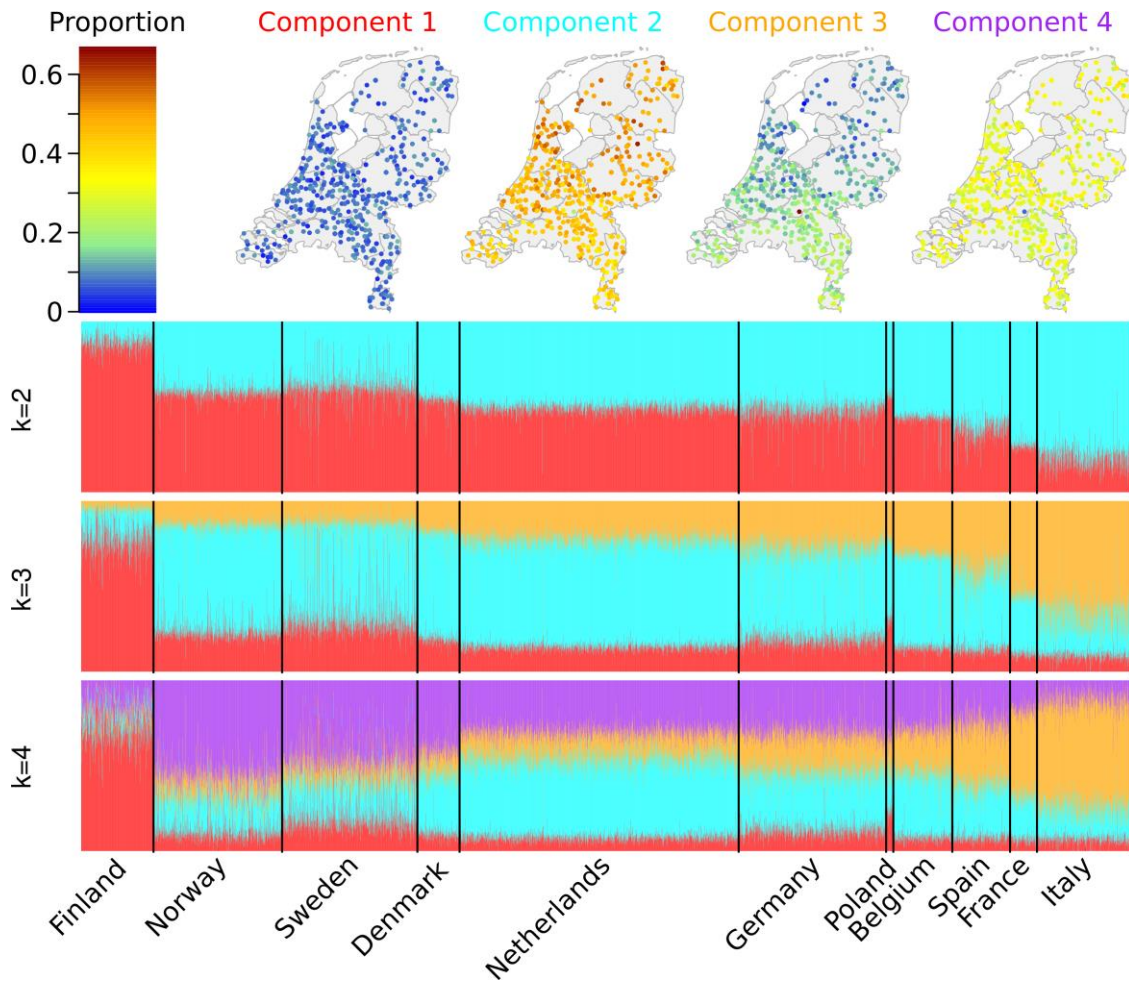
- 498 17. van Rheenen, W. *et al.* Genome-wide association analyses identify new risk variants and the genetic architecture
499 of amyotrophic lateral sclerosis. *Nat. Genet.* **48**, 1043–1048 (2016).
- 500 18. Lawson, D. J., Hellenthal, G., Myers, S. & Falush, D. Inference of population structure using dense haplotype
501 data. *PLoS Genet.* **8**, e1002453 (2012).
- 502 19. Sawcer, S. *et al.* Genetic risk and a primary role for cell-mediated immune mechanisms in multiple sclerosis.
503 *Nature* **476**, 214–219 (2011).
- 504 20. Browning, B. L. & Browning, S. R. Improving the accuracy and efficiency of identity-by-descent detection in
505 population data. *Genetics* **194**, 459–471 (2013).
- 506 21. Lawson, D. J. & Falush, D. Population identification using genetic data. *Annu Rev Genomics Hum Genet* **13**,
507 337–361 (2012).
- 508 22. Palamara, P. F. Population genetics of identity by descent. (2014).
- 509 23. Al-Asadi, H., Petkova, D., Stephens, M. & Novembre, J. Estimating recent migration and population-size
510 surfaces. *PLoS Genet* **15**, e1007908 (2019).
- 511 24. Browning, S. R. & Browning, B. L. Accurate Non-parametric Estimation of Recent Effective Population Size
512 from Segments of Identity by Descent. *Am J Hum Genet* **97**, 404–418 (2015).
- 513 25. Herlihy, D. *The Black Death and the Transformation of the West*. (Harvard University Press, 1997).
- 514 26. Durbin, R. Efficient haplotype matching and storage using the positional Burrows-Wheeler transform (PBWT).
515 *Bioinforma. Oxf. Engl.* **30**, 1266–1272 (2014).
- 516 27. Athanasiadis, G. *et al.* Nationwide Genomic Study in Denmark Reveals Remarkable Population Homogeneity.
517 *Genetics* **204**, 711–722 (2016).
- 518 28. Nalls, M. A. *et al.* Measures of autozygosity in decline: globalization, urbanization, and its implications for
519 medical genetics. *PLoS Genet.* **5**, e1000415 (2009).
- 520 29. Roosen, J. & Curtis, D. R. The ‘light touch’ of the Black Death in the Southern Netherlands: an urban trick? *Econ*
521 *Hist Rev* **72**, 32–56 (2019).
- 522 30. Haworth, S. *et al.* Apparent latent structure within the UK Biobank sample has implications for epidemiological
523 analysis. *Nat. Commun.* **10**, 333 (2019).
- 524 31. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience*
525 **4**, 7 (2015).
- 526 32. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74
527 (2015).

- 528 33. Delaneau, O., Marchini, J. & Zagury, J.-F. cois. A linear complexity phasing method for thousands of genomes.
529 *Nat Methods* **9**, 179–181 (2011).
- 530 34. CoreTeam, R. *R: A Language and Environment for Statistical Computing. Vienna, Austria: R Foundation for*
531 *Statistical Computing; 2015.* (2015).
- 532 35. Scrucca, L., Fop, M., Murphy, T. B. & Raftery, A. E. mclust 5: Clustering, Classification and Density Estimation
533 Using Gaussian Finite Mixture Models. *R J* **8**, 289–317 (2016).
- 534 36. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated individuals.
535 *Genome Res.* **19**, 1655–1664 (2009).
- 536 37. Browning, B. L. & Browning, S. R. Detecting identity by descent and estimating genotype error rates in
537 sequence data. *Am J Hum Genet* **93**, 840–851 (2013).
- 538 38. Price, A. L. *et al.* Long-range LD can confound genome scans in admixed populations. *Am J Hum Genet* **83**,
539 132–5; author reply 135–9 (2008).
- 540 39. Nakazawa, M. fmsb: Functions for medical statistics book with some demographic data, 2014. *R Package*
541 (2018).
- 542 40. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide
543 association studies. *Nat. Genet.* **47**, 291–295 (2015).
- 544 41. van der Maaten, L. & Hinton, G. Visualizing data using t-SNE. *J. Mach. Learn. Res.* **9**, 2579–2605 (2008).

546 **Supplementary material**

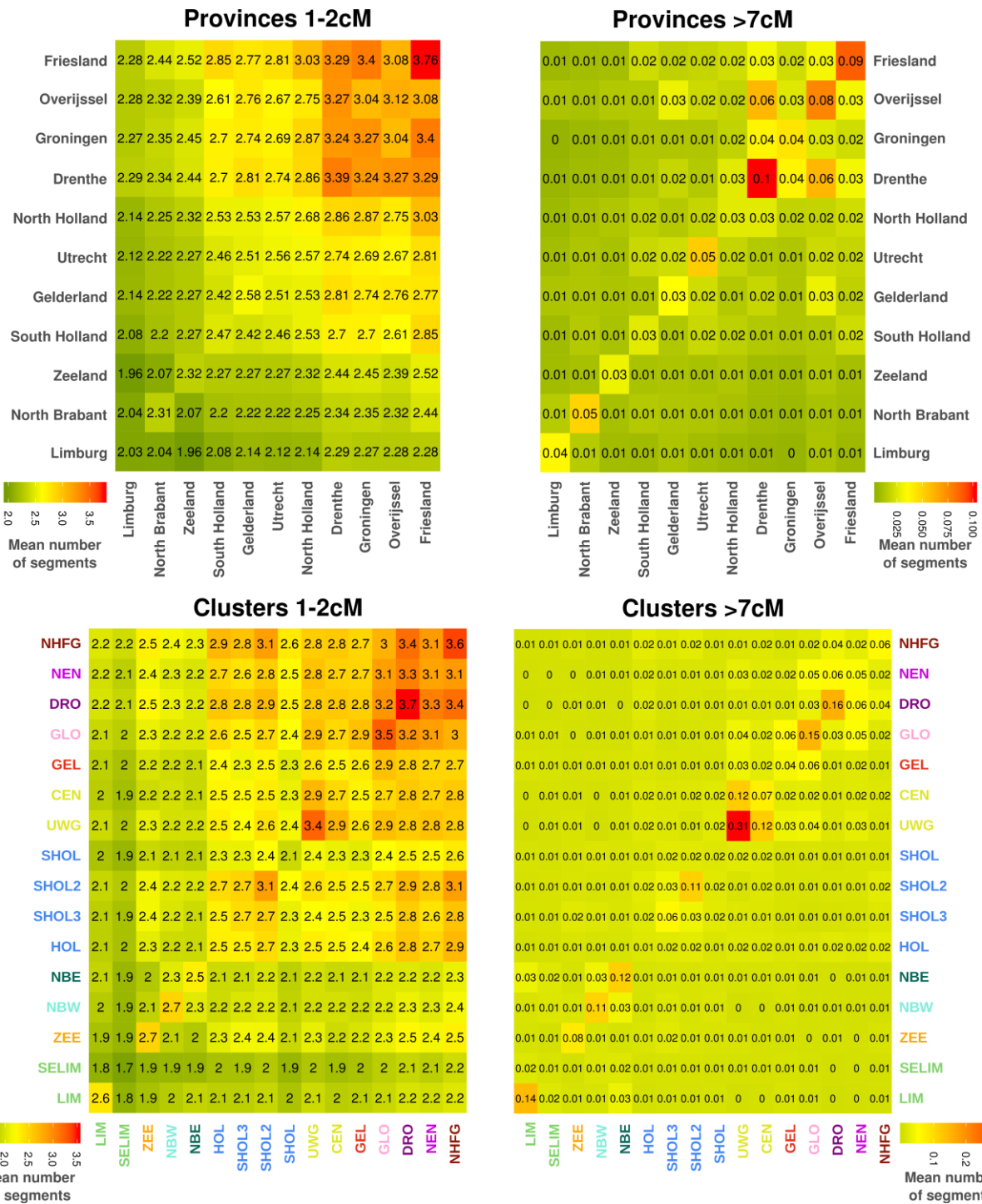
547 **Supplementary table 1 Mean pairwise F_{ST} ($\times 10^{-3}$) for European groups from Sawcer *et al.*¹⁹**

| | Finland | Sweden | Norway | Germany | Italy | Denmark | Belgium | Spain | Poland | France |
|----------------|----------------|---------------|---------------|----------------|--------------|----------------|----------------|--------------|---------------|---------------|
| Finland | 0 | 3.93 | 5.10 | 5.85 | 11.1 | 5.56 | 6.71 | 9.99 | 5.81 | 7.59 |
| Sweden | - | 0 | 0.362 | 0.702 | 4.87 | 0.362 | 1.07 | 3.76 | 1.84 | 1.81 |
| Norway | - | - | 0 | 0.899 | 4.96 | 0.407 | 1.07 | 3.73 | 2.56 | 1.76 |
| Germany | - | - | - | 0 | 2.77 | 0.399 | 0.289 | 2.11 | 1.26 | 0.678 |
| Italy | - | - | - | - | 0 | 4.08 | 2.29 | 1.21 | 5.42 | 1.55 |
| Denmark | - | - | - | - | - | 0 | 0.526 | 3.01 | 2.13 | 1.22 |
| Belgium | - | - | - | - | - | - | 0 | 1.53 | 2.43 | 0.367 |
| Spain | - | - | - | - | - | - | - | 0 | 4.56 | 0.660 |
| Poland | - | - | - | - | - | - | - | - | 0 | 2.84 |
| France | - | - | - | - | - | - | - | - | - | 0 |



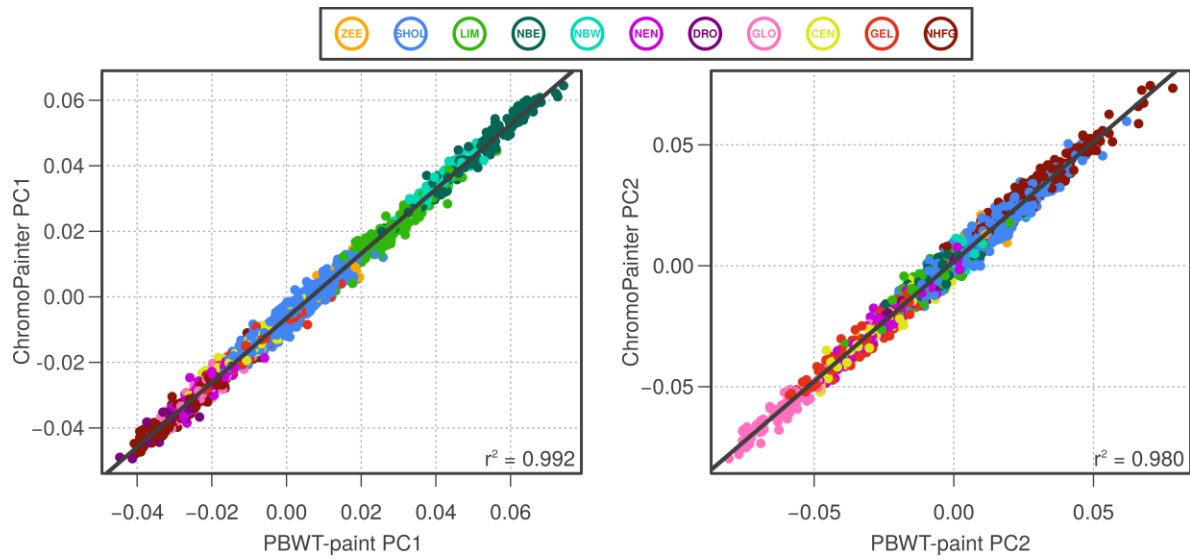
549

550 **Supplementary figure 1 ADMIXTURE modelling for Dutch and European samples.** Maps depict the regional breakdown of
551 ADMIXTURE components for k=4 split. Dutch samples have a high value for admixture component 2, which is next highest in
552 Germany and Belgium. Components 2 and 3 show opposing north-south gradients in the Netherlands, with component 2 highest
553 in the north and component 3 highest in the south. Component 3 is best represented in southern European countries such as Italy.



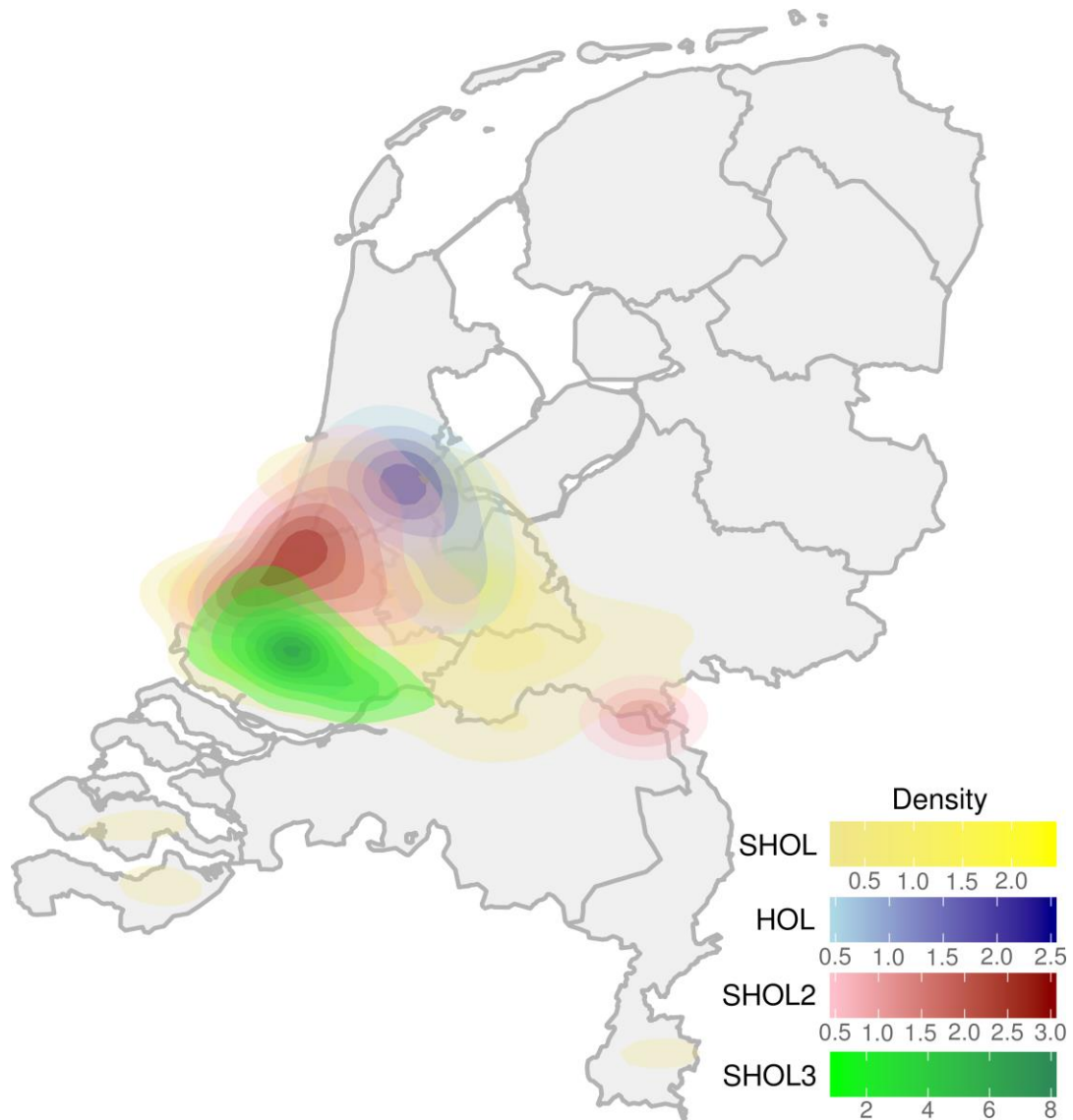
554

555 **Supplementary figure 2 Old (left) and recent (right) IBD sharing per province (top) and per cluster (bottom).** Average
 556 sharing of old (short) segments is enriched in northern provinces and clusters. Average sharing of recent (long) segments is
 557 higher on average within clusters than within provinces, indicating haplotypic clustering captures marginally more recent
 558 ancestry.



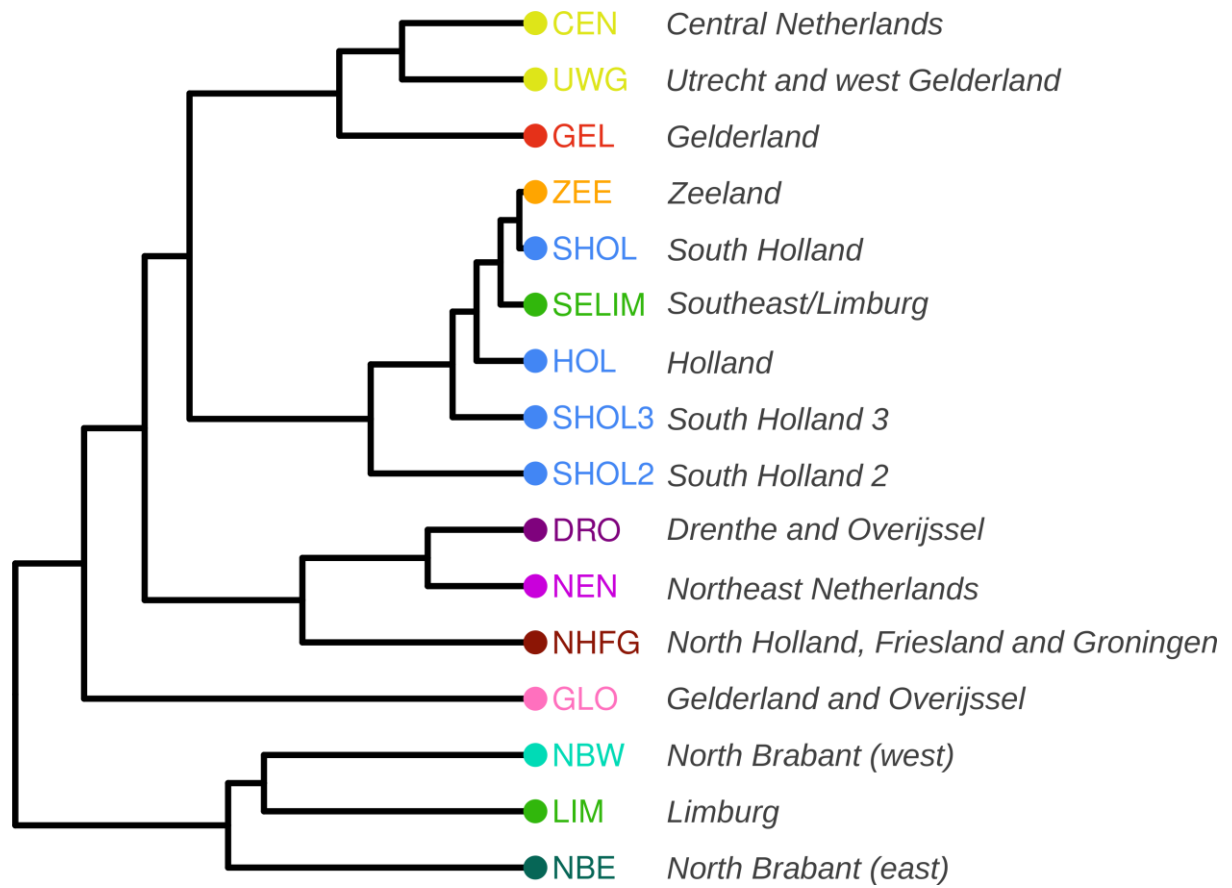
559

560 **Supplementary figure 3 Benchmark of PBWT-paint vs ChromoPainter.** Scatterplots comparing the first two principal
561 components (PCs) of the coancestry matrices produced by ChromoPainter and PBWT-paint, showing strong correlation. Points
562 are coloured by cluster groups defined in Figure 1. For all pairwise comparisons in the two coancestry matrices, Pearson's $\rho =$
563 0.82 ($0.82-0.821$; $p < 2 \times 10^{-16}$).



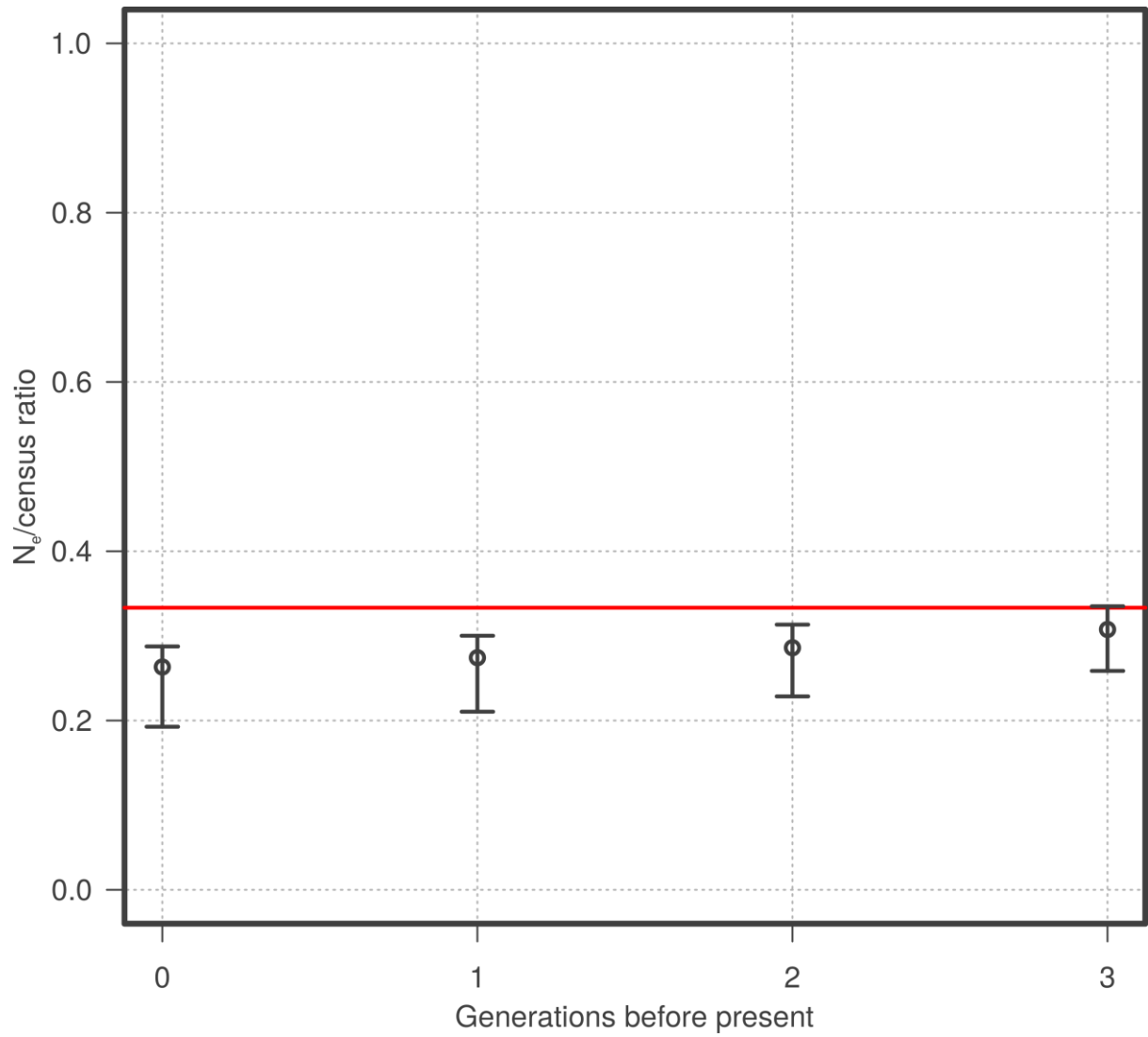
564

565 **Supplementary figure 4 Geographic distribution of South Holland clusters from the SHOL cluster group.** 2D kernel
566 density estimates are shown for the geographic spread of samples from clusters SHOL (yellow), HOL (blue), SHOL2 (red), and
567 SHOL3 (green) which form the SHOL cluster group in Figure 1. Kernel density estimates were calculated using the
568 `stat_density2d` function in `ggplot2` (R version 3.2.3) with default settings. >80% of samples are contained within plotted polygons
569 for each cluster. Notably, although overlapping, three of the four clusters show quite distinct geographic ranges.



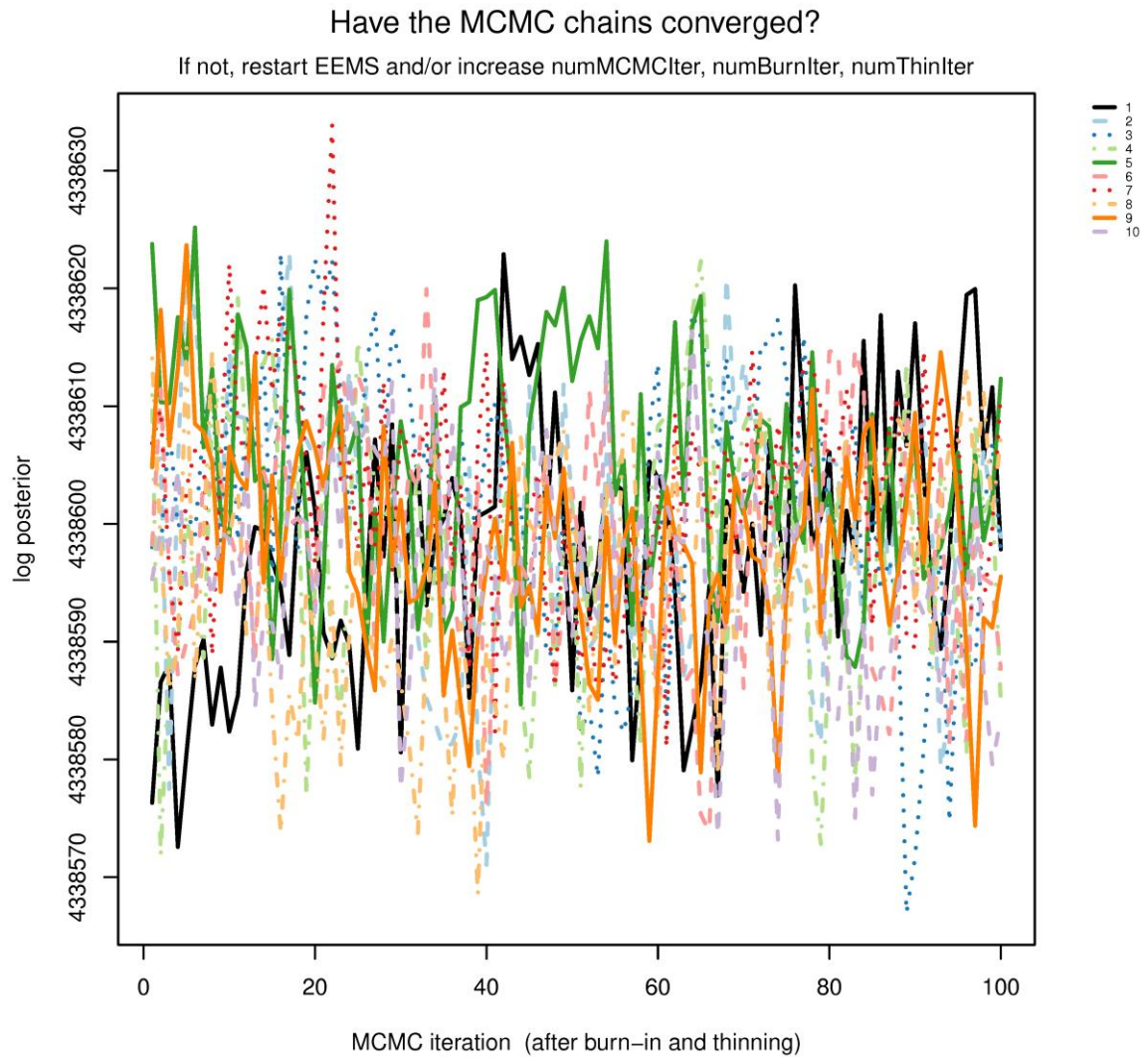
570

571 **Supplementary figure 5 Total variation distance (TVD) tree for k=16 split in the Netherlands.** Clusters are coloured and
572 labelled according to scheme in Figure 1.



573

574 **Supplementary figure 6 Ratio of estimated N_e /Census is stable over the past 3 generations.** The red line at 0.33 corresponds
575 to the expected ratio of N_e to census if lifespan is 3 times the generation time.



576

577 **Supplementary figure 7 Convergence of MCMC chains for EEMS run in Netherlands.** 10 independently seeded MCMC
578 chains reach approximate convergence.

579 Supplementary note 1

580 Project MinE ALS GWAS Consortium authors

581 Wouter van Rheenen¹, Aleksey Shatunov², Russell L. McLaughlin³, Rick A.A. van der Spek¹, Alfredo Iacoangeli^{2,4},
582 Kevin P. Kenna¹, Kristel R. van Eijk¹, Nicola Ticozzi^{5,6}, Boris Rogelj^{7,8}, Katarina Vrabec⁹, Metka Ravnik-Glavač^{9,10},
583 Blaž Koritnik¹¹, Janez Zidar¹¹, Lea Leonardis¹¹, Leja Dolenc Grošelj¹¹, Stéphanie Millecamps¹², François
584 Salachas^{12,13,14}, Vincent Meininger^{15,16}, Mamede de Carvalho^{17,18}, Susana Pinto¹⁷, Marta Gromicho¹⁷, Ana Pronto-
585 Laborinho¹⁷, Jesus S. Mora¹⁹, Ricardo Rojas-García^{20,21}, Meraida Polak^{22,23}, Siddharthan Chandran^{24,25}, Shuna
586 Colville²⁴, Robert Swingler²⁴, Karen E. Morrison²⁶, Pamela J. Shaw²⁷, John Hardy²⁸, Richard W. Orrell²⁹, Alan
587 Pittman^{28,30}, Katie Sidle²⁹, Pietro Fratta³¹, Andrea Malaspina^{32,33}, Simon Topp², Susanne Petri³⁴, Susanna Abdulla³⁵,
588 Carsten Drepper³⁶, Michael Sendtner³⁶, Thomas Meyer³⁷, Roel A. Ophoff^{38,39}, Kim A. Staats³⁹, Martina Wiedau-
589 Pazos⁴⁰, Catherine Lomen-Hoerth⁴¹, Vivianna M. Van Deerlin⁴², John Q. Trojanowski⁴², Lauren Elman⁴³, Leo
590 McCluskey⁴³, A. Nazli Basak⁴⁴, Thomas Meitinger⁴⁵, Peter Lichtner⁴⁵, Milena Blagojevic-Radivojkov⁴⁵, Christian R.
591 Andres⁴⁶, Gilbert Bensimon^{47,48,49}, Bernhard Landwehrmeyer⁵⁰, Alexis Brice^{51,52,53,54,55}, Christine A.M. Payan^{47,49},
592 Safaa Saker-Delye⁵⁶, Alexandra Dürr⁵⁷, Nicholas W. Wood⁵⁸, Lukas Tittmann⁵⁹, Wolfgang Lieb⁵⁹, Andre Franke⁶⁰,
593 Marcella Rietschel⁶¹, Sven Cichon^{62,63,64,65,66}, Markus M. Nöthen^{62,63}, Philippe Amouyel⁶⁷, Jean-François
594 Dartigues⁶⁸, Andre G. Uitterlinden^{69,70}, Fernando Rivadeneira^{69,70}, Karol Estrada⁶⁹, Albert Hofman^{70,71}, Charles
595 Curtis^{72,73}, Anneke J. van der Kooi⁷⁴, Markus Weber⁷⁵, Christopher E. Shaw², Bradley N. Smith², Daisy Sproviero⁷⁶,
596 Cristina Cereda⁷⁶, Mauro Ceroni⁷⁷, Luca Diamanti⁷⁷, Roberto Del Bo⁷⁸, Stefania Corti⁷⁸, Giacomo P. Comi⁷⁸, Sandra
597 D'Alfonso⁷⁹, Lucia Corrado⁷⁹, Cinzia Bertolin⁸⁰, Gianni Sorarù⁸⁰, Letizia Mazzini⁸¹, Viviana Pensato⁸², Cinzia
598 Gellera⁸², Cinzia Tiloca⁵, Antonia Ratti^{5,6}, Andrea Calvo^{83,84}, Cristina Moglia^{83,84}, Maura Brunetti^{83,84}, Rosa
599 Capozzo⁸⁵, Chiara Zecca⁸⁵, Christian Lunetta⁸⁶, Silvana Penco⁸⁷, Nilo Riva⁸⁸, Alessandro Padovani⁸⁹, Massimiliano
600 Filosto⁹⁰, PARALS registry⁹¹, SLALOM group⁹¹, SLAP registry⁹¹, SLAGEN Consortium⁹¹, NNIPPS Study Group⁹¹,
601 Ian Blair⁹², Garth A. Nicholson^{92,93}, Dominic B. Rowe⁹², Roger Pamphlett⁹⁴, Matthew C. Kiernan⁹⁵, Julian
602 Grosskreutz⁹⁶, Otto W. Witte⁹⁶, Robert Steinbach⁹⁶, Tino Prell⁹⁶, Beatrice Stubendorff⁹⁶, Ingo Kurth^{97,98}, Christian
603 A. Hübner⁹⁷, P. Nigel Leigh⁹⁹, Federico Casale⁸³, Adriano Chio^{83,84}, Ettore Beghi¹⁰⁰, Elisabetta Pupillo¹⁰⁰, Rosanna
604 Tortelli¹⁰¹, Giancarlo Logroscino^{102,103}, John Powell², Albert C. Ludolph⁵⁰, Jochen H. Weishaupt⁵⁰, Wim
605 Robberecht^{104,105,106}, Philip Van Damme^{104,105,106}, Robert H. Brown¹⁰⁷, Jonathan D. Glass^{22,23}, John E. Landers¹⁰⁷,
606 Orla Hardiman^{108,109}, Peter M. Andersen¹¹⁰, Philippe Corcia^{46,111,112}, Patrick Vourc'h⁴⁶, Vincenzo Silani^{5,6}, Michael
607 A. van Es¹, R. Jeroen Pasterkamp¹¹³, Cathryn M. Lewis^{72,114}, Gerome Breen^{72,73}, Ammar Al-Chalabi², Leonard H.
608 van den Berg¹, Jan H. Veldink¹

- 609 1. Department of Neurology, UMC Utrecht Brain Center, University Medical Center Utrecht, Utrecht University,
610 Utrecht, The Netherlands.
- 611 2. Maurice Wohl Clinical Neuroscience Institute, King's College London, Department of Basic and Clinical
612 Neuroscience, London, UK.
- 613 3. Complex Trait Genomics Laboratory, Smurfit Institute of Genetics, Trinity College Dublin, Dublin, Republic of
614 Ireland.
- 615 4. Department of Biostatistics and Health Informatics, Institute of Psychiatry, Psychology and Neuroscience,
616 King's College London, London, UK.

- 617 5. Department of Neurology and Laboratory of Neuroscience, IRCCS Istituto Auxologico Italiano, Milano, Italy.
- 618 6. Department of Pathophysiology and Transplantation, 'Dino Ferrari' Center, Università degli Studi di Milano,
619 Milano, Italy.
- 620 7. Department of Biotechnology, Jožef Stefan Institute, Ljubljana, Slovenia.
- 621 8. Biomedical Research Institute BRIS, Ljubljana, Slovenia.
- 622 9. Department of Molecular Genetics, Institute of Pathology, Faculty of Medicine, University of Ljubljana, SI-
623 1000 Ljubljana, Slovenia.
- 624 10. Institute of Biochemistry, Faculty of Medicine, University of Ljubljana, SI-1000 Ljubljana, Slovenia.
- 625 11. Ljubljana ALS Centre, Institute of Clinical Neurophysiology, University Medical Centre Ljubljana, SI-1000
626 Ljubljana, Slovenia.
- 627 12. Institut du Cerveau et de la Moelle épinière, Inserm U1127, CNRS UMR 7225, Sorbonne Universités, UPMC
628 Univ Paris 06 UMRS1127, Paris, France.
- 629 13. Centre de Référence Maladies Rares SLA Ile de France, Département de Neurologie, Hôpital de la Pitié-
630 Salpêtrière, Paris, France.
- 631 14. GRC-UPMC SLA et maladies du Motoneurone, France.
- 632 15. Ramsay Generale de Santé, Hopital Peupliers, Paris, France.
- 633 16. Réseau SLA Ile de France.
- 634 17. Instituto de Fisiologia, Instituto de Medicina Molecular, Faculdade de Medicina, Universidade de Lisboa,
635 Lisbon, Portugal
- 636 18. Department of Neurosciences, Hospital de Santa Maria-CHLN, Lisbon, Portugal.
- 637 19. ALS Unit, Hospital San Rafael, Madrid, Spain
- 638 20. Neurology Department, Hospital de la Santa Creu i Sant Pau de Barcelona, Autonomous University of
639 Barcelona, Barcelona, Spain.
- 640 21. Centro de Investigación en red en Enfermedades Raras (CIBERER), Spain.
- 641 22. Department Neurology, Emory University School of Medicine, Atlanta, GA, USA.
- 642 23. Emory ALS Center, Emory University School of Medicine, Atlanta, GA, USA.
- 643 24. Euan MacDonald Centre for Motor Neurone Disease Research, Edinburgh, UK.
- 644 25. Centre for Neuroregeneration and Medical Research Council Centre for Regenerative Medicine, University of
645 Edinburgh, Edinburgh, UK.
- 646 26. School of Medicine, Dentistry and Biomedical Sciences, Queen's University Belfast, UK.
- 647 27. Sheffield Institute for Translational Neuroscience (SITraN), University of Sheffield, Sheffield, UK.
- 648 28. Department of Molecular Neuroscience, Institute of Neurology, University College London, UK.
- 649 29. Department of Clinical Neuroscience, Institute of Neurology, University College London, UK.
- 650 30. Reta Lila Weston Institute, Institute of Neurology, University College London, UK.
- 651 31. Department of Neuromuscular Diseases, UCL Queen Square Institute of Neurology.
- 652 32. Centre for Neuroscience and Trauma, Blizard Institute, Queen Mary University of London, London, UK.
- 653 33. North-East London and Essex Regional Motor Neuron Disease Care Centre, London, UK.
- 654 34. Department of Neurology, Hannover Medical School, Hannover, Germany.
- 655 35. Department of Neurology, Otto-von-Guericke University Magdeburg, Magdeburg, Germany.
- 656 36. Institute of Clinical Neurobiology, University Hospital Wuerzburg, Germany.

- 657 37. Charité – Universitätsmedizin, Berlin, Germany.
- 658 38. Department of Human Genetics, David Geffen School of Medicine, University of California, Los Angeles, CA,
659 USA.
- 660 39. Center for Neurobehavioral Genetics, Semel Institute for Neuroscience and Human Behavior, University of
661 California, Los Angeles, CA, USA.
- 662 40. Department of Neurology, David Geffen School of Medicine, University of California, Los Angeles, CA, USA.
- 663 41. Department of Neurology, University of California, San Francisco, CA, USA.
- 664 42. Center for Neurodegenerative Disease Research, Perelman School of Medicine at the University of
665 Pennsylvania, Philadelphia, PA, USA.
- 666 43. Department of Neurology, Perelman School of Medicine at the University of Pennsylvania, PA USA.
- 667 44. Koç University, School of Medicine, KUTTAM-NDAL, Istanbul Turkey.
- 668 45. Institute of Human Genetics, Helmholtz Zentrum München, Neuherberg, Germany.
- 669 46. Centre SLA, CHRU de Tours, Tours, France; UMR 1253, iBrain, Université de Tours, Inserm, Tours, France.
- 670 47. APHP, Département de Pharmacologie Clinique, Hôpital de la Pitié-Salpêtrière, France.
- 671 48. Université Pierre & Marie Curie, Pharmacologie, Paris VI, Paris, France.
- 672 49. BESPIM, CHU-Nîmes, Nîmes, France.
- 673 50. Department of Neurology, Ulm University, Ulm, Germany.
- 674 51. INSERM U 1127, Hôpital de la Pitié-Salpêtrière, 75013 Paris, France.
- 675 52. CNRS UMR 7225, Hôpital de la Pitié-Salpêtrière, 75013 Paris, France.
- 676 53. Sorbonne Universités, Université Pierre et Marie Curie Paris 06 UMRS 1127, Hôpital de la Pitié-Salpêtrière,
677 75013 Paris, France.
- 678 54. Institut du Cerveau et de la Moelle épinière, Hôpital de la Pitié-Salpêtrière, 75013 Paris, France.
- 679 55. APHP, Département de Génétique, Hôpital de la Pitié-Salpêtrière, 75013 Paris, France.
- 680 56. Genethon, CNRS UMR 8587 Evry, France.
- 681 57. Department of Medical Genetics, l'Institut du Cerveau et de la Moelle Épinière, Hoptial Salpêtrière, 75013
682 Paris, France.
- 683 58. Department of Neurogenetics, Institute of Neurology, University College London, UK.
- 684 59. PopGen Biobank and Institute of Epidemiology, Christian Albrechts-University Kiel, Kiel, Germany.
- 685 60. Institute of Clinical Molecular Biology, Kiel University, Kiel, Germany.
- 686 61. Department of Genetic Epidemiology in Psychiatry, Central Institute of Mental Health, Faculty of Medicine
687 Mannheim, University of Heidelberg, Germany
- 688 62. Institute of Human Genetics, University of Bonn, Bonn, Germany.
- 689 63. Department of Genomics, Life and Brain Center, Bonn, Germany.
- 690 64. Division of Medical Genetics, University Hospital Basel, University of Basel, Basel, Switzerland.
- 691 65. Department of Biomedicine, University of Basel, Basel, Switzerland.
- 692 66. Institute of Neuroscience and Medicine INM-1, Research Center Juelich, Juelich, Germany.
- 693 67. University of Lille, Inserm, CHU Lille, Institut Pasteur de Lille, U1167 - RID-AGE - Risk Factor and molecular
694 determinants of aging diseases, Labex Distalz, F-59000 Lille, France.
- 695 68. Bordeaux University, ISPED, Centre INSERM U1219-Epidemiologie-Biostatistique & CIC-1401, CHU de
696 Bordeaux, Pole de Sante Publique, Bordeaux, France.

- 697 69. Department of Internal Medicine, Genetics Laboratory, Erasmus Medical Center Rotterdam, Rotterdam, The
698 Netherlands.
- 699 70. Department of Epidemiology, Erasmus Medical Center Rotterdam, Rotterdam, The Netherlands.
- 700 71. Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA, USA.
- 701 72. Social, Genetic & Developmental Psychiatry Centre, Institute of Psychiatry, Psychology & Neuroscience,
702 King's College London, London, UK.
- 703 73. NIHR Maudsley Biomedical Research Centre, Maudsley Hospital and Institute of Psychiatry, Psychology &
704 Neuroscience, King's College London, London, UK.
- 705 74. Amsterdam UMC, department of Neurology, University of Amsterdam, Neuroscience, Amsterdam
- 706 75. Neuromuscular Diseases Unit/ALS Clinic, Kantonsspital St. Gallen, 9007, St. Gallen, Switzerland.
- 707 76. Genomic and post-Genomic Center, IRCCS Mondino Foundation, Pavia, Italy.
- 708 77. General Neurology, IRCCS Mondino Foundation, Pavia, Italy
- 709 78. Neurologic Unit, IRCCS Foundation Ca' Granda Ospedale Maggiore Policlinico, Milan, Italy.
- 710 79. Department of Health Sciences, Interdisciplinary Research Center of Autoimmune Diseases, UPO, Università
711 del Piemonte Orientale, Novara, Italy.
- 712 80. Department of Neurosciences, University of Padova, Padova, Italy.
- 713 81. ALS Centre Department of Neurology Maggiore della Carità University Hospital, Novara.
- 714 82. Unit of Genetics of Neurodegenerative and Metabolic Diseases, Fondazione IRCCS Istituto Neurologico 'Carlo
715 Besta', Milano, Italy.
- 716 83. "Rita Levi Montalcini" Department of Neuroscience, ALS Centre, University of Torino, Turin, Italy.
- 717 84. Azienda Ospedaliera Città della Salute e della Scienza, Torino, Italy.
- 718 85. Department of Clinical research in Neurology, University of Bari "A. Moro", at Pia Fondazione "Card. G.
719 Panico", Tricase (LE), Italy.
- 720 86. NEMO Clinical Center, Serena Onlus Foundation, Niguarda Ca' Granda Hospital, Milan, Italy.
- 721 87. Medical Genetics Unit, Department of Laboratory Medicine, Niguarda Ca' Granda Hospital, Milan, Italy.
- 722 88. Department of Neurology, Institute of Experimental Neurology (INSPE), Division of Neuroscience, San
723 Raffaele Scientific Institute, Milan, Italy.
- 724 89. Neurology Unit, Department of Clinical and Experimental Sciences, University of Brescia, Italy.
- 725 90. Neurology Unit, Department of Neuroscience and Vision, Spedali Civili Hospital, Brescia, Italy.
- 726 91. A list of members and affiliations appears at the end of this Supplementary Note.
- 727 92. Department of Biomedical Sciences, Faculty of Medicine and Health Sciences, Macquarie University, Sydney,
728 New South Wales, Australia.
- 729 93. University of Sydney, ANZAC Research Institute, Concord Hospital, Sydney, New South Wales, Australia.
- 730 94. Discipline of Pathology, Sydney Medical School, Brain and Mind Centre, The University of Sydney, New
731 South Wales 2050, Australia.
- 732 95. Brain and Mind Centre, The University of Sydney, New South Wales 2050, Australia.
- 733 96. Hans-Berger Department of Neurology, Jena University Hospital, Jena, Germany.
- 734 97. Institute of Human Genetics, Jena University Hospital, Jena, Germany.
- 735 98. Institute of Human Genetics, Medical Faculty, RWTH Aachen University, Aachen, Germany

- 736 99. Department of Neurology, Brighton and Sussex Medical School Trafford Centre for Biomedical Research,
737 University of Sussex, Falmer, East Sussex, UK.
- 738 100.Laboratory of Neurological Diseases, Department of Neuroscience, IRCCS Istituto di Ricerche Farmacologiche
739 Mario Negri, Milano, Italy.
- 740 101.Institute of Neurology, University College of London (UCL), London, UK
- 741 102.Department of Basic Medical Sciences, Neuroscience and Sense Organs, University of Bari 'Aldo Moro', Bari,
742 Italy.
- 743 103.Unit of Neurodegenerative Diseases, Department of Clinical Research in Neurology, University of Bari 'Aldo
744 Moro', at Pia Fondazione Cardinale G. Panico, Tricase, Lecce, Italy.
- 745 104.KU Leuven - University of Leuven, Department of Neurosciences
- 746 105.VIB, Center for Brain & Disease Research, Laboratory of Neurobiology, Leuven, Belgium.
- 747 106.University Hospitals Leuven, Department of Neurology, Leuven, Belgium.
- 748 107.Department of Neurology, University of Massachusetts Medical School, Worcester, MA, USA.
- 749 108.Academic Unit of Neurology, Trinity College Dublin, Trinity Biomedical Sciences Institute, Dublin, Republic
750 of Ireland.
- 751 109.Department of Neurology, Beaumont Hospital, Dublin, Republic of Ireland.
- 752 110.Department of Clinical Science, Neurosciences, Umeå University, Sweden.
- 753 111.Federation des Centres SLA Tours and Limoges, LITORALS, Tours, France.
- 754 112.INSERM U1253, "iBrain", Université François-Rabelais de Tours, Faculté de Médecine 10, Bd Tonnelé,
755 37032 Tours Cedex 1, France.
- 756 113.Department of Translational Neuroscience, UMC Utrecht Brain Center, University Medical Center Utrecht,
757 Utrecht University, Utrecht, The Netherlands.
- 758 114.Department of Medical and Molecular Genetics, King's College London, London, UK.

759 *Italian Consortium for the Genetics of ALS (SLAGEN) members*

760 Daniela Calini, Isabella Fogh, Antonia Ratti, Vincenzo Silani, Nicola Ticozzi, Cinzia Tiloca, Barbara Castellotti,
761 Cinzia Gellera, Viviana Pensato, Franco Taroni, Cristina Cereda, Mauro Ceroni, Stella Gagliardi, Giacomo Comi,
762 Stefania Corti, Roberto Del Bo, Lucia Corrado, Sandra D'Alfonso, Letizia Mazzini, Elena Pegoraro, Giorgia Querin,
763 Massimiliano Filosto and Gianni Sorarù

764 *Registro Lombardo Sclerosi Laterale Amyotrofica (SLALOM) group members*

765 Francesca Gerardi, Fabrizio Rinaldi, Maria Sofia Cotelli, Luca Chiveri, Maria Cristina Guaita, Patrizia Perrone,
766 Giancarlo Comi, Carlo Ferrarese, Lucio Tremolizzo, Marialuisa Delodovici, Massimiliano Filosto and Giorgio Bono

767 *Piemonte and Valle d'Aosta Registry for Amyotrophic Lateral Sclerosis (PARALS) group members*

768 Stefania Cammarosano, Antonio Canosa, Dario Cocito, Leonardo Lopiano, Luca Durelli, Bruno Ferrero, Antonio
769 Bertolotto, Alessandro Mauro, Luca Pradotto, Roberto Cantello, Enrica Bersano, Dario Giobbe, Maurizio Gionco,
770 Daniela Leotta, Lucia Appendino, Roberto Cavallo, Enrico Odddenino, Claudio Geda, Fabio Poglio, Paola
771 Santimaria, Umberto Massazza, Antonio Villani, Roberto Conti, Fabrizio Pisano, Mario Palermo, Franco Vergnano,

772 Paolo Provera, Maria Teresa Penza, Marco Aguggia, Nicoletta Di Vito, Piero Meineri, Ilaria Pastore, Paolo
773 Ghiglione, Danilo Seliak, Nicola Launaro, Giovanni Astegiano and Bottacchi Edo

774 *Sclerosi Laterale Amyotrofica-Puglia (SLAP) registry members*

775 Isabella Laura Simone, Stefano Zoccolella, Michele Zarrelli and Franco Apollo

776 *Neuroprotection and Natural History in Parkinson Plus Syndromes (NNIPPS) Study group members*

777 William Camu, Jean Sebastien Hulot, Francois Viallet, Philippe Couratier, David Maltete, Christine Tranchant,
778 Marie Vidailhet.