

1 **DISENTANGLING SIGNATURES OF SELECTION BEFORE AND AFTER EUROPEAN COLONIZATION IN LATIN**
2 **AMERICANS**

3
4 Javier Mendoza-Revilla^{1,2,3*}, J. Camilo Chacón-Duque^{4,5}, Macarena Fuentes-Guajardo⁶, Louise Ormond¹, Ke
5 Wang⁷, Malena Hurtado³, Valeria Villegas³, Vanessa Granja³, Victor Acuña-Alonzo⁸, Claudia Jaramillo⁹, William
6 Arias⁹, Rodrigo Barquera Lozano^{7,8}, Jorge Gómez-Valdés⁸, Hugo Villamil-Ramírez^{10,11}, Caio C. Silva de
7 Cerqueira¹², Keyla M. Badillo Rivera¹³, Maria A. Nieves-Colón¹⁴, Christopher R. Gignoux¹⁵, Genevieve L. Wojcik¹⁶,
8 Andrés Moreno-Estrada¹⁷, Tábita Hunemeier¹², Virginia Ramallo^{12,18}, Lavinia Schuler-Faccini¹², Rolando
9 Gonzalez-José¹⁸, Maria-Cátira Bortolini¹², Samuel Canizales-Quinteros^{10,11}, Carla Gallo³, Giovanni Poletti³,
10 Gabriel Bedoya⁹, Francisco Rothhammer¹⁹, David Balding^{1,20}, Matteo Fumagalli²¹, Kaustubh Adhikari²², Andrés
11 Ruiz-Linares^{1,23,24*¶} and Garrett Hellenthal^{1*¶}

12
13 ¹ Department of Genetics, Evolution and Environment, and UCL Genetics Institute, University College London,
14 London, UK

15 ² Human Evolutionary Genetics Unit, Institut Pasteur, UMR2000, CNRS, Paris, France

16 ³ Laboratorios de Investigación y Desarrollo, Facultad de Ciencias y Filosofía, Universidad Peruana Cayetano
17 Heredia, Lima, Perú

18 ⁴ Centre for Palaeogenetics, Stockholm, Sweden

19 ⁵ Department of Archaeology and Classical Studies, Stockholm University, Stockholm, Sweden

20 ⁶ Departamento de Tecnología Médica, Facultad de Ciencias de la Salud, Universidad de Tarapacá, Arica, Chile.

21 ⁷ Department of Archaeogenetics, Max Planck Institute for the Science of Human History, Jena, Germany

22 ⁸ National Institute of Anthropology and History, Mexico City, Mexico

23 ⁹ GENMOL (Genética Molecular), Universidad de Antioquia, Medellín, Colombia

24 ¹⁰ Unidad de Genómica de Poblaciones Aplicada a la Salud, Facultad de Química, UNAM-Instituto Nacional de
25 Medicina Genómica, Mexico City, Mexico

26 ¹¹ Universidad Nacional Autónoma de México e Instituto Nacional de Medicina Genómica, Mexico City, Mexico

27 ¹² Departamento de Genética, Universidade Federal do Rio Grande do Sul, Porto Alegre, Brazil

28 ¹³ Department of Genetics, Stanford School of Medicine, Stanford, California, United States

29 ¹⁴ Department of Anthropology, University of Minnesota Twin Cities, Minneapolis, Minnesota, United States

30 ¹⁵ University of Colorado Anschutz Medical Campus, Aurora, Colorado, United States

31 ¹⁶ Bloomberg School of Public Health, John Hopkins University, Baltimore, Maryland, United States

1 ¹⁷ Laboratorio Nacional de Genómica para la Biodiversidad (UGA-LANGEBIO), CINVESTAV, Irapuato,
2 Guanajuato, Mexico
3 ¹⁸ Instituto Patagónico de Ciencias Sociales y Humanas-Centro Nacional Patagónico, CONICET, Puerto Madryn,
4 Argentina
5 ¹⁹ Instituto de Alta Investigación, Universidad de Tarapacá, Arica, Chile.
6 ²⁰ Schools of BioSciences and Mathematics & Statistics, University of Melbourne, Melbourne, Australia
7 ²¹ School of Biological and Behavioural Sciences, Queen Mary University of London, London, UK
8 ²² School of Mathematics and Statistics, Faculty of Science, Technology, Engineering and Mathematics, The
9 Open University, Milton Keynes, UK
10 ²³ Ministry of Education Key Laboratory of Contemporary Anthropology and Collaborative Innovation Center of
11 Genetics and Development, Fudan University, Shanghai, China
12 ²⁴ Aix-Marseille Université, CNRS, EFS, ADES, Marseille, France
13
14 [¶]These authors jointly supervised this work
15 * Correspondence to: javier.mendoza@upch.pe (J.M.R); andresruiz@fudan.edu.cn (A.R.L.);
16 g.hellenthal@ucl.ac.uk (G.H.)
17

1 **Abstract**

2 Throughout human evolutionary history, large-scale migrations have led to intermixing (i.e., admixture)
3 between previously separated human groups. While classical and recent work have shown that studying
4 admixture can yield novel historical insights, the extent to which this process contributed to adaptation
5 remains underexplored. Here, we introduce a novel statistical model, specific to admixed populations, that
6 identifies loci under selection while determining whether the selection likely occurred post-admixture or prior
7 to admixture in one of the ancestral source populations. Through extensive simulations we show that this
8 method is able to detect selection, even in recently formed admixed populations, and to accurately
9 differentiate between selection occurring in the ancestral or admixed population. We apply this method to
10 genome-wide SNP data of ~4,000 individuals in five admixed Latin American cohorts from Brazil, Chile,
11 Colombia, Mexico and Peru. Our approach replicates previous reports of selection in the HLA region that are
12 consistent with selection post-admixture. We also report novel signals of selection in genomic regions spanning
13 47 genes, reinforcing many of these signals with an alternative, commonly-used local-ancestry-inference
14 approach. These signals include several genes involved in immunity, which may reflect responses to endemic
15 pathogens of the Americas and to the challenge of infectious disease brought by European contact. In addition,
16 some of the strongest signals inferred to be under selection in the Native American ancestral groups of modern
17 Latin Americans overlap with genes implicated in energy metabolism phenotypes, plausibly reflecting
18 adaptations to novel dietary sources available in the Americas.

19
20

ACCEPTED MANUSCRIPT

1 Introduction

2 Admixed populations offer a unique opportunity to detect recent selection. In the human lineage, genomic
3 studies have demonstrated the pervasiveness of admixture events in the history of the vast majority of human
4 populations (Patterson et al. 2012; Hellenthal et al. 2014; Lazaridis et al. 2014). By inferring the ancestral
5 origins of particular genetic loci in the genomes of recently admixed individuals, recent studies have provided
6 evidence that such admixture has facilitated the spread of adaptative genetic mutations in humans. Notable
7 examples include the transfer of a protective allele in the Duffy blood group gene likely providing resistance to
8 *Plasmodium vivax* malaria in Malagasy and Cape Verdeans from sub-Saharan Africans (Hodgson et al. 2014;
9 Pierron et al. 2018; Hamid et al. 2021), and the transmission of the lactase persistence allele in the Fula
10 pastoralists from Western Eurasians (Vicente et al. 2019).

11 An ideal setting in which to test whether and how admixture contributed to genetic adaptation is Latin
12 America. The genetic make-up of present day Latin Americans stems mainly from three ancestral populations:
13 indigenous Native Americans, Europeans (mainly from the Iberian Peninsula), and Sub-Saharan Africans (Wang
14 et al. 2007; Moreno-Estrada et al. 2013; Moreno-Estrada et al. 2014; Homburger et al. 2015; Chacon-Duque et
15 al. 2018; Luisi et al. 2020) that were brought together starting ~500 years ago. The admixed genomes of Latin
16 Americans are thus the result of an intermixing process between human populations that had been evolving
17 independently for tens-of-thousands of years and that were suddenly brought together in a new environment.
18 In this new environment, the ancestral genomes were quickly subjected to novel pressures that were largely
19 unfamiliar from where they firstly evolved. Therefore, the genomes of Latin Americans potentially harbor
20 signals of recent adaptations attributable to beneficial variants, e.g. introduced from a particular ancestral
21 population, increasing rapidly in frequency post-admixture. Motivated by this, several studies have explored
22 the genomes of admixed Latin Americans for signatures of selection occurring since the admixture event (Tang
23 et al. 2007; Basu et al. 2008; Ettinger et al. 2009; Guan 2014; Rishishwar et al. 2015; Deng et al. 2016; Zhou et
24 al. 2016; Norris et al. 2020; Vicuna et al. 2020). These studies have relied on an approach similar to that of
25 admixture mapping, where the ancestry of a genomic region in each admixed individual is assigned to a
26 particular ancestral population, a technique known as local-ancestry-inference (LAI). Loci with significantly
27 more inferred ancestry inherited from one ancestral population are assumed to have evolved under some form
28 of selection (Tang et al. 2007).

29 In addition, the genetic make-up of Latin Americans offers the opportunity to detect selection in their ancestral
30 populations, as large cohorts of Latin Americans can be leveraged to reconstruct genetic variation patterns in
31 each source population. This is of particular use for exploring selection in Native Americans, since Native
32 American groups are currently underrepresented in genomic studies (Sirugo et al. 2019) and as a consequence
33 only a few studies have centered on detecting adaptive signals of indigenous groups from the Americas. Such
34 studies have identified strong selective signals at different genes, particularly at those related to immunity,
35 highlighting the selective pressures that Native Americans were subjected to after they entered the continent
36 (Lindo et al. 2018; Reynolds et al. 2019; Avila-Arcos et al. 2020).

37 With some exceptions (Cheng et al. 2021), these studies either limited their analyses to Latin Americans with
38 high Native American ancestry or used LAI to infer loci in individuals that derive from a Native American source.
39 However, such approaches may result in a reduction of statistical power due to removal of individuals with
40 non-Native American ancestry, inaccurate local ancestry estimation and/or through removing segments
41 challenging to assign.

42 Here we present a novel statistical model that identifies loci that have undergone selection before or after an
43 admixture event (which we refer to as pre- or post-admixture selection, respectively). In contrast to previous
44 methods, this approach is based on allele frequencies and does not require assignments of local ancestry along

1 the genome. We illustrate the utility of our new method by performing a selection scan in five Latin American
2 cohorts collected as part from the CANDELA Consortium (Ruiz-Linares et al. 2014). Our results suggest that
3 several loci have been subjected to natural selection in admixed Latin American populations, and in their
4 ancestral populations, replicating many of these signals using LAI. Many of the putative selected SNPs are
5 strongly associated to relevant phenotypes, or act as expression quantitative loci (eQTL) in relevant tissues,
6 providing further evidence of their functional effect. Overall, our analyses highlight the usefulness of our
7 method to detect signals of selection in admixed populations or their ancestral populations, and reveal novel
8 candidate genes implicated in the adaptive history of groups from the American continent.

9 **Results**

10 ***Overview of AdaptMix***

11 In part following Balding and Nichols (1995), and analogous to previous approaches (Long 1991; Mathieson et
12 al. 2015; Cheng et al. 2021), our model AdaptMix assumes that, under neutrality, the allele frequencies of an
13 admixed target population can be described using a beta-binomial model, with expected allele frequency equal
14 to a mixture of sampled allele frequencies from a set of groups that act as surrogates to the admixing sources
15 (fig. 1). In our case the admixed target population is a Latin American cohort, defined below, and we use three
16 surrogate groups to represent Native American, European, and African admixing source populations. The
17 mixture values are inferred a priori, e.g. using ADMIXTURE (Alexander et al. 2009) (fig. 1a), qpAdm (Haak et al.
18 2015) or SOURCEFIND (Chacon-Duque et al. 2018), as the average amount of ancestry that each admixed target
19 individual matches to a set of reference populations. (The reference populations used by these programs may
20 be the same as the surrogate populations, but they need not be as illustrated below.) We find the variance
21 parameter that maximises the likelihood of this beta-binomial model across all SNPs. This variance term aims
22 to limit the number of false-positives attributable to genetic drift in the target population following admixture
23 and/or the use of inaccurate surrogates for the ancestral populations. Then, at each SNP, we calculate the
24 probability of observing allele counts equal to or more extreme than those observed in the target population,
25 hence providing a *P*-value testing the null hypothesis that the SNP is neutral (see Methods).

26 Assuming a pulse of admixture, this test is designed to detect selection occurring: (i) in the admixed population
27 following the admixture event (i.e. along the purple line “e” in fig. 1b), and/or (ii) in one (or more) of the
28 source/surrogate pairings (i.e. along the red and/or blue lines (a)-(d) in fig 1b). Note that scenario (ii) includes
29 selection occurring in any of the ancestral source populations (i.e. along the lines “c” or “d” in fig. 1b) and/or in
30 any of the surrogate populations (i.e. along the lines “a” or “b” in fig. 1b). At SNPs with evidence of selection
31 (i.e. low *P*-values), we distinguish between (i) and (ii) by exploring how genotype counts of admixed target
32 individuals relate to their inferred admixture proportions contributed by each surrogate. Under scenario (i), we
33 assume that selection affects all target individuals equally, regardless of their admixture proportions, which in
34 turn assumes all ancestries were present when selection occurred. In contrast, under scenario (ii), we expect
35 selection to more strongly affect one of the source/surrogate population pairings. Intuitively, if (ii) is true,
36 individuals with nearly 100% ancestry from the source/surrogate pair experiencing selection will have genotype
37 counts that deviate the most from expectations under the neutral model, while individuals with nearly 0%
38 ancestry from this pair will have counts that closely follow the neutral model (fig. 1c). If instead (i) is true, this
39 pattern is attenuated, though it can be challenging in practice to distinguish (ii) from (i) if allele frequencies
40 strongly differ between surrogate groups (fig. 1d). Assuming a multiplicative model of selection, which is
41 numerically close to an additive model, we find the selection coefficients that maximize the fit of the data to
42 model (i) and to model (ii) when separately treating each source/surrogate pair as the selected group. We
43 report ratios of likelihoods, equivalent here to using differences in Akaike Information Criterion (AIC), to
44 quantify our ability to distinguish among scenarios (i) and (ii).

1 In summary, for each tested SNP we infer (a) a P -value testing the null hypothesis of neutrality, (b) the relative
2 evidence (i.e. likelihood ratios) for whether selection occurred post-admixture or in one of the admixing
3 sources and (c) the selection strength summed across time.

4 **Simulations**

5 We tested our approach using simulations designed to resemble our Latin American cohort in terms of
6 sample size, inferred admixture proportions, and the extent to which our surrogates match the true
7 admixing sources. As post-admixture selection in recently admixed population is challenging to detect
8 unless selection is strong, we included selection coefficients (s) of large magnitude. We note that the
9 upper range values are consistent with those estimated in recently admixed populations, including Latin
10 Americans (Zhou et al 2016, Pierron et al 2019, Vicente et al 2019, Hamid et al 2021) (see Methods).
11 At a false-positive rate of 5×10^{-5} , these simulations indicate we have ~ 50 - 90% power to detect selection for
12 scenario (i) (i.e., post-admixture selection) with $s=0.15$ - 0.20 , with s defined as the selection strength per
13 generation in homozygotes carrying two copies of the selected allele, and selection occurring over 12
14 generations under various modes of selection (additive, dominant, multiplicative, recessive) (fig. 2a,
15 supplementary fig. S1). For scenario (ii), in the case of selection occurring in the Native American source, power
16 depends on the overall amount of Native American ancestry (fig. 2a). As an example, Brazil-like simulations
17 ($<15\%$ average Native American ancestry) show little power, Colombia-like simulations ($\sim 30\%$ average Native
18 American ancestry) typically exhibit $>50\%$ power, and other simulated populations (~ 50 - 70% average Native
19 American ancestry) exhibit $>75\%$ power under scenario (ii) assuming $s \geq 0.1$ over 50 generations, with similar
20 power if instead $s \sim 0.025$ over 150 generations (supplementary fig. S2). Simulations including a bottleneck in
21 the Native American source population (see Methods) showed reduced power, likely because the stronger
22 genetic drift both masks the selection signal (Refoyo-Martínez et al. 2019; Cuadros-Espinoza et al. 2021) and
23 makes the surrogate population more genetically differentiated from its corresponding source (supplementary
24 fig. S3). Detecting selection occurring in the European or African source depends on the overall amount of
25 European and African ancestry in a similar manner (e.g., fig. 2a, supplementary fig. S4-S5). For SNPs where we
26 detect selection, we mis-classify the type of selection $\leq 2\%$ of the time, e.g., concluding post-admixture
27 selection when the truth is selection in the Native American source $\sim 1\%$ of the time across all selection
28 coefficients (fig. 2b). However, our approach often fails to classify selection scenarios unless selection strengths
29 are large (e.g., $s > 0.1$).

30 We also compared the power of AdaptMix to that of Ohana, a recently developed maximum likelihood method
31 that infers selection by modeling ancestral admixture components, which has been shown to have similar or
32 higher power to other state-of-the-art methods (Cheng et al. 2021). Following Cheng et al. (2021), we
33 simulated a realistic demographic model relating four populations meant to represent African, East Asian,
34 European, and Native American sources. We also simulated an admixed population that descends from a 50-
35 50% mixture of the European and Native American sources, with selection occurring prior to admixture in only
36 the ancestral Native American source (see Methods). We then applied AdaptMix and Ohana to four sampled
37 populations that descend from the Africa, East Asian, European, and admixed populations. In these
38 simulations, AdaptMix has ~ 0.4 - 4.8% less power than Ohana if running Ohana with an ideal number of
39 ancestry components K , in this case $K=4$, that distinguishes the admixed population (supplementary fig. S6),
40 and if Ohana only tests for selection in the ancestry component most representative of the admixed population
41 (fig. 3a). However, AdaptMix has up to 5.12% more power than Ohana in this setting when using Ohana's more
42 general test that does not assume selection only in the admixed population. Furthermore, if using a suboptimal
43 K , e.g. $K=3$, Ohana's power is greatly reduced, since the Native American and East Asian sources are both

1 classified into the same ancestry component (supplementary fig. S6). We also performed simulations,
2 mimicking those in Cuadros-Espinoza et al. (2021), under which selection occurs post-admixture in the admixed
3 population, with admixture occurring 70 generations ago (see Methods). In these simulations, AdaptMix
4 outperformed Ohana even when using the ideal number of clusters $K=3$, presumably because Ohana does not
5 classify the admixed individuals into their own ancestry component (supplementary fig. S6), which should
6 maximize its power. In these post-admixture simulations, both AdaptMix and Ohana outperform two local
7 ancestry deviation (LAD) approaches (RFMix, ELAI) (Maples et al. 2013; Guan 2014), perhaps because the older
8 admixture time resulted in difficulties accurately assigning local ancestry segments to source populations.

9 **Applying AdaptMix to the five Latin American cohorts of CANDELA**

10 We divided Latin Americans into five cohorts based on country of origin: Brazil ($n=190$), Chile ($n=896$),
11 Colombia ($n=1125$), Mexico ($n=773$), and Peru ($n=834$), using individuals sampled as part of the CANDELA
12 Consortium (Ruiz-Linares et al. 2014), testing each cohort for selection separately (supplementary fig. S7).
13 Analyzing each cohort by country of origin results in a higher number of individuals, and thus increases the
14 statistical power to detect selection. As demonstrated in Chacon-Duque et al (2018), however, there is notable
15 population sub-structure within each country. To test for robustness of our selection signals to this sub-
16 structure, we supplemented each of these analyses by testing subsets of individuals within a country based on
17 their inferred ancestry matching to Native American reference groups from Chacon-Duque et al. (2018). This
18 gave six additional tested groups with sufficient ancestry represented: 'Mapuche' ($n=434$) in Chile, 'Chibcha
19 Paez' ($n=200$) in Colombia, 'Nahua' ($n=466$) and 'South Mexico' ($n=78$) in Mexico, and 'Andes Piedmont'
20 ($n=195$) and 'Quechua' ($n=147$) in Peru (supplementary fig. S8). To infer the proportion of African, European,
21 and Native American ancestry in each Latin American, we applied unsupervised ADMIXTURE with $K=3$ clusters
22 jointly to all CANDELA individuals and 553 Native American, Iberian, and West African reference individuals (fig.
23 1a).

24 Note that the choice of surrogate populations defines the selection test between each surrogate and its
25 corresponding ancestral source in scenario (ii). In this way, our test acts as an analogue to F_{ST} comparing two
26 populations, but while accounting for admixture in one of the populations. As an illustration, we tested the
27 Brazilian cohort for selection using northwest Europeans from England and Scotland (GBR) from the 1000
28 Genomes Project (1KGP) (The 1000 Genomes Project Consortium 2015) as a surrogate for the Brazilian cohort's
29 European ancestry source (supplementary fig. S9). Given the majority ($\sim 80\%$) of ancestry in the Brazilian cohort
30 is related to Iberian Europeans, this test is most-powered to detect selection acting along the branch
31 separating present-day northwest Europeans and descendants of Iberians who traveled to Brazil post-
32 Columbus. In this analysis, we infer strongest signals of selection at the *HERC2/OCA2* and *LCT/MCM6* genes.
33 This replicates previously reported selection signals when comparing northwest Europeans to present-day
34 Iberians (Poulter et al. 2003; Bersaglieri et al. 2004), and likely indicates selection for lighter skin pigmentation
35 and lactase persistence in northwest Europeans that is unrelated to any selection in the Americas.

36 As another example, we also tested each Latin American cohort separately while using Han Chinese from
37 Beijing (CHB) from the 1KGP as a surrogate for Native American ancestry (supplementary fig. S10). In this
38 analysis, SNPs that follow model (ii) indicate selection along the branch separating present-day Han Chinese
39 and Native American populations. For this test, we find the strongest signals of selection at previously reported
40 selected genes in East Asians, including those related to alcohol metabolism such as *ADH7* and *ADH1B*
41 (Galinsky et al. 2016; Gu et al. 2018) that both are classified as selection under model (ii). The strongest overall
42 signal in this scan, which was unclassified, overlapped the *POU2F3* gene, implicated in the regulation of viral
43 transcription, keratinocyte differentiation and other cellular events. Selection signals at this gene have been

1 reported to be under selection in Native American populations from throughout the Americas (Amorim et al.
2 2017), and also shows evidence for Neanderthal adaptive introgression in East Asians (Racimo et al. 2017).
3 For our main analyses, we use 205 Iberians (from 1KGP and Chacon-Duque et al. (2018)) to represent European
4 ancestry surrogates. Therefore, given the likely short split time between present-day Iberians and Europeans
5 that migrated to the Americas during the colonial era, we are underpowered to detect selection in the
6 European source only (see simulations). We use 206 West Africans from the 1KGP to represent the African
7 ancestry source, which has been reported as a good proxy to the African genetic sources (from Chacon-Duque
8 et al. (2018)). For this reason, we should similarly have low power to find selection occurring only in the African
9 source/surrogate. At any rate we do not test for selection related to African ancestry, because the Latin
10 American cohort here have ~6% African ancestry on average, limiting power further (see supplementary fig.
11 S5). We combined 142 individuals with <1% non-Native American inferred ancestry from 19 Native American
12 groups (supplementary table S1) to represent the Native American surrogate. By using individuals sampled
13 from geographically spread Native American groups as the Native American ancestry surrogate, we aim to
14 identify regional selection signals experienced by some Native American groups but not others. We also expect
15 to have the highest power when testing for selection type (ii) in Native Americans, as there is likely to be the
16 most time separating this ‘average’ Native American surrogate and the admixing source of each regional Latin
17 American cohort. To avoid confounding our inference, we excluded individuals with >1% inferred ancestry
18 matching to surrogates other than Native Americans, Iberian Europeans, and West Africans using SOURCEFIND
19 (Chacon-Duque et al. 2018). Also, since the time since admixture among these groups is relatively short in the
20 CANDELA cohort (likely <15 generations ago), detecting selection post-admixture can only identify relatively
21 strong selection signals (see simulations).

22 ***AdaptMix identifies 47 regions of putative selection***

23 For each Latin American cohort, we considered SNPs under selection as those having P -values less than the
24 5×10^{-5} false-positive threshold in the population-matched neutral simulations, which corresponds to a model-
25 based P -value of 6.75×10^{-6} – 1.07×10^{-7} (supplementary table S2). For Chile, Colombia, Mexico and Peru, we
26 report loci that pass these criteria both in the analysis of all individuals from that country and in at least one of
27 three alternative analyses for that country that are designed to test for robustness to latent population
28 structure (supplementary fig. S11). The first of these alternative analyses consisted of identifying signals of
29 selection using AdaptMix on each of the six Native American subsets defined above (e.g., in either the ‘Andes
30 Piedmont’ or ‘Quechua’ subset when testing for selection in Peruvians) (supplementary table S3). The other
31 two alternative analyses were based on LAI. In particular we used ELAI (Guan 2014) to assign each genomic
32 region of an admixed individual to a Native American, European, or African ancestral source. For the second
33 alternative analysis, designed to test for post-admixture selection, we assessed whether the proportion of
34 ancestry inferred from one of these three sources in a local region deviated substantially from the genome-
35 wide average (supplementary table S4). For the third alternative analysis, designed to test for selection in the
36 Native American source, we instead used the Population Branch Statistic (PBS) (Yi et al. 2010) to test for
37 selection in one of the six Native American subset groups defined above, using allele frequencies computed
38 from LAI-inferred Native American segments from the subset of individuals representing that Native American
39 group (see Methods) (supplementary fig. S8 and supplementary table S5).

40 Overall, we find 51 candidate regions to have evidence of positive or purifying selection passing the
41 criteria above, 47 of which target protein-coding genes (supplementary table S6 and fig. 4). Four of
42 these 47 candidate gene regions contain at least one SNP exhibiting strong evidence (likelihood ratio
43 >1,000) of selection affecting all admixed individuals regardless of ancestry proportions, which we

1 assume reflects post-admixture selection. Furthermore, 18 of these 47 regions exhibit strong evidence
2 of selection containing at least one SNP (likelihood ratio $>1,000$) in the Native American source only.
3 The 25 remaining candidate gene regions are unclassified into either type of selection (likelihood ratio
4 $\leq 1,000$).

5 To prioritize candidate casual genes, we annotated the protein-coding gene that had the highest overall
6 Variant-to-Gene (V2G) scores (Ghousaini et al. 2021; Ochoa et al. 2021) for the SNPs showing the strongest
7 evidence of selection in each candidate gene region. The overall V2G score aggregates differentially weighted
8 evidence of variant-gene association from several sources, including cis-QTL data, chromatin interaction
9 experiments, *in silico* function predictions (e.g., Variant Effect Predictor from Ensembl), and distance between
10 the variant and each gene's canonical transcription starting site. For each of these candidate genes we then
11 annotated the phenotype with the highest overall association score based on the Open Targets Platform
12 (Koscielny et al. 2017; Ochoa et al. 2021).

13 While most of these associated phenotypes represent genetic disorders, syndromes, or different types of
14 measurements (medically or non-medically-related), many are also related to immune response and diet – two
15 major selective forces previously reported to shape the human genome (Karlsson et al. 2014; Fan et al. 2016).
16 We therefore organize the description of our candidate selection signals into two main sections below that
17 cover only these two features, with signals of all other hits in supplementary table S6. For brevity, below we
18 only highlight putatively selected regions where at least one significant SNP had an associated GWAS or eQTL
19 signal. For our significant SNPs related to immune-response genes, GWAS signals included SNPs associated to
20 white blood cell counts in a large multi-continental cohort (including Latin American individuals) (Chen et al.
21 2020), and eQTL signals included cis-associated SNPs to gene expression in 15 immune-related cell types from
22 the DICE project (Schmiedel et al. 2018). For our significant SNPs related to diet, GWAS signals included
23 metabolic, anthropometric, and lipid levels from the UK Biobank cohort (Loh et al. 2018), and eQTL signals
24 included cis-associated SNPs to gene expression in adipose, muscle, and liver tissue from the GTEx Project
25 (Lonsdale et al. 2013).

26 **Signals at immune-related genes**

27 Fifteen of the 47 candidate gene regions contained at least one protein-coding gene either related to the
28 development or regulation of the immune system or that has been previously associated to the quantification
29 of immune cell types, susceptibility progression to infectious diseases, or autoimmune disorders. For example,
30 we replicate a well-known signal encompassing several immune-related genes at 6p21 that are part of the
31 human leukocyte antigen (HLA) system (fig. 4 and supplementary fig. S12-S14). These included SNPs (AdaptMix
32 P -value $< 5.00 \times 10^{-7}$) near several MHC class I genes (*HLA-G*, *HLA-H*, *HLA-A*, and *HLA-J*) in each of the Chilean,
33 Colombian, Mexican and Peruvian cohorts, with the Colombian cohort containing several SNPs classified as
34 being selected post-admixture (likelihood ratio $> 1,000$). Encouragingly, we inferred African ancestry enrichment
35 (Z -score > 2.5) in each cohort ~ 60 kb downstream from our top AdaptMix signals using LAI, with maximum Z -
36 score > 9 (one-sided P -value $< 4.09 \times 10^{-21}$) in the Chilean cohort (fig. 5). In addition, other signals were inferred
37 upstream in the Chilean cohort at a 5' UTR SNP of the *ZBTB12* gene (rs2844455, AdaptMix P -value $= 5.45 \times 10^{-8}$),
38 the Mexican cohort at an intronic SNP of *HLA-DMA* (rs28724903, AdaptMix P -value $= 3.87 \times 10^{-8}$), and the
39 Peruvian cohort at an intronic SNP of the MHC class III gene *STK19* (rs6941112, AdaptMix P -value $= 7.57 \times 10^{-9}$).
40 Many of these HLA genes have been previously characterized as subject to be under selection post-admixture
41 in different Latin American populations by showing an excess of African ancestry at the HLA locus (Tang et al.
42 2007; Basu et al. 2008; Ettinger et al. 2009; Guan 2014; Rishishwar et al. 2015; Deng et al. 2016; Zhou et al.
43 2016; Norris et al. 2020; Vicuna et al. 2020).

1 In addition to HLA, we infer previously unreported selection signals in four candidate gene regions that each
 2 harbor genes with well-established roles in the immune system, with each region containing at least one SNP
 3 significantly associated (P -value $<5\times 10^{-8}$) to white blood cell counts or the expression of an immune-related
 4 gene in immune cells (P -value $<10^{-5}$) (see Methods). Among these, one signal at 1p13 in the Chilean cohort
 5 encompasses the *CD101* gene (fig. 6a), which belongs to a family of cell-surface immunoglobulins superfamily
 6 proteins and plays a role as an inhibitor of T-cell proliferation (Soares et al. 1998; Bouloc et al. 2000). Within
 7 this region five SNPs are classified as being selected post-admixture and show also an increase of LAI-inferred
 8 European ancestry (maximum Z-score=3.40, one-sided P -value=3.36 $\times 10^{-4}$). Strikingly, the region contains a
 9 synonymous SNP (Ile588, CADD score of 9.23) (rs3736907, AdaptMix P -value=1.05 $\times 10^{-9}$) that strongly affects
 10 *CD101* expression in T cells (eQTL P -value $< 2.42\times 10^{-10}$) and is associated with neutrophil (GWAS P -
 11 value=2.08 $\times 10^{-10}$) and total white cell count (GWAS P -value=3.61 $\times 10^{-9}$) (fig. 6a).

12 The second signal, at 18p11 also in Chileans, encompasses the *PTPN2* gene, a tyrosine-specific phosphatase
 13 involved in the Janus kinase (JAK)-signal transducer and activator of transcription (STAT) signaling pathway (fig.
 14 6b). The JAK-STAT pathway has an important role in the control of immune responses, and dysregulation of this
 15 pathway is associated with various immune disorders (Shuai and Liu 2003). Several SNPs with low AdaptMix P -
 16 values (P -value $<1.69\times 10^{-7}$) in the 18p11 region are also associated with eosinophil counts (GWAS P -
 17 value $<1.13\times 10^{-10}$) and the expression of *PTPN2* in natural killer (NK) cells (eQTL P -value $<1.14\times 10^{-9}$) (fig. 6b).

18 The other two novel signals, both in the Peruvian cohort, are consistent with selection in Native Americans only
 19 (likelihood-ratio $>1,000$). The first, at 17q25, contains the *CD300LF* gene that encodes for a membrane
 20 glycoprotein that contains an immunoglobulin domain, and which plays an important role in the maintenance
 21 of immune homeostasis by promoting macrophage-mediated efferocytosis (Borrego 2013). Notably, a 3'UTR
 22 SNP (rs9913698, AdaptMix P -value=3.11 $\times 10^{-9}$) is strongly associated with monocyte count (GWAS P -
 23 value=1.00 $\times 10^{-33}$), total white cell count (GWAS P -value=5.96 $\times 10^{-24}$), lymphocyte count (GWAS P -
 24 value=2.50 $\times 10^{-19}$), and neutrophil count (GWAS P -value=1.30 $\times 10^{-9}$) (supplementary fig. S15). The second signal
 25 is at 22q11 adjacent to the *MIF* gene (fig. 6c), which is implicated in macrophage function in host defense
 26 through the suppression of anti-inflammatory effects of glucocorticoids (Calandra and Roger 2003). Variants
 27 within *MIF* have been recently associated to rheumatoid arthritis in southern Mexican patients (Santoscoy-
 28 Ascencio et al. 2020). The SNP rs2330635 (AdaptMix P -value=7.06 $\times 10^{-8}$) is strongly associated to the expression
 29 of *MIF* in T-cells (eQTL P -value $<8.63\times 10^{-5}$) and NK cells (eQTL P -value=5.77 $\times 10^{-9}$) and is also marginally
 30 associated to neutrophil counts (GWAS P -value=2.46 $\times 10^{-6}$) (fig. 6c).

31 Overall, these findings suggest that some of the clearest signals of adaptation in the Americas can be ascribed
 32 to immune-related selective pressures. These plausibly resulted from both the introduction of novel pathogens
 33 after European colonization and the endemic pathogens encountered by the first Native Americans during the
 34 initial peopling of the continent.

35 **Signals at genes related to diet**

36 Among the 47 candidate regions, nine regions contained at least one protein-coding gene potentially related to
 37 dietary practices through their association with metabolism-related phenotypes or anthropometric-related
 38 measurements (supplementary table S6). Among these, we infer three previously unreported signals where at
 39 least one of the selected SNPs was associated to metabolic- or anthropometric-related phenotypes, or to the
 40 expression of the candidate gene in adipose, muscle, or liver tissue (see Methods). One of these three hits
 41 (rs4636058, AdaptMix P -value=5.70 $\times 10^{-10}$), at 6p22 in the Chilean cohort, is classified as being selected post-
 42 admixture and shows an increase of LAI-inferred European ancestry (Z-score=3.78, one-sided P -value=7.82 $\times 10^{-4}$).
 43 It is located at 6q22 and encompasses the *SLC35F1* gene, whose function is not known, though several
 44 studies have associated this gene with different measurements of cardiac function (Hoffmann et al. 2017;

1 Warren et al. 2017; Giri et al. 2019). Notably, SNP rs4636058 is marginally associated to cholesterol levels
2 (UKBB GWAS P -value= 3.8×10^{-4}) and body fat percentage (UKBB GWAS P -value= 4.29×10^{-4}). Another of these
3 three hits, at 1q31 in the Mexican cohort, is consistent with selection in Native Americans (likelihood-
4 ratio $>1,000$) (fig. 7a). The 1q31 signal includes an intronic SNP (rs1171148, AdaptMix P -value= 2.31×10^{-8}) of
5 *BRINP3*, a gene associated to body mass index in studies across different human groups (Pulit et al. 2019; Zhu
6 et al. 2020). Within this region, various SNPs are associated to different metabolic-related phenotypes,
7 including the SNP rs1171148 that is associated with hip circumference (UKBB GWAS P -value= 4.96×10^{-8}) and
8 marginally associated with body mass index (UKBB GWAS P -value= 5.51×10^{-5}) (fig. 7a).

9 Finally, the third hit (rs5030938, AdaptMix P -value= 3.79×10^{-15}), which had the highest overall AdaptMix score,
10 is inferred in the Peruvian cohort at 10q22 and indicates selection in Native Americans (likelihood-ratio $>1,000$)
11 (fig. 7b). This SNP is associated with the expression of *HKDC1* in liver (eQTL P -value= 2.19×10^{-5}), adipose visceral
12 (eQTL P -value= 1.46×10^{-5}), and adipose subcutaneous tissue (eQTL P -value= 1.36×10^{-4}) (fig. 6b). *HKDC1* encodes
13 and hexokinase that catalyzes the rate-limiting and first obligatory step of glucose metabolism (Ludvik et al.
14 2016), and several studies have associated variants within this gene with glucose levels in pregnant women
15 (Hayes et al. 2013; Guo et al. 2015; Kanthimathi et al. 2016; Tan et al. 2019) and with weight at birth
16 (Warrington et al. 2019).

17 Overall, these results support previous hypothesis that genes related to energy metabolism were probably
18 critical in the establishment of stable human populations in distinct ecoregions (Hancock et al. 2010), including
19 those of the Americas (Amorim et al. 2017; Reynolds et al. 2019).

20 Discussion

21 Analytical considerations

22 Here we present AdaptMix, a novel statistical model that identifies loci under selection in admixed populations.
23 Our model is based on the principle that allele frequencies in an admixed population can be modeled as a
24 linear combination of the allele frequencies in the ancestral populations proportional to their admixing
25 contributions, and that deviations from the expectation can be a product of selection. This selection test is
26 related to the work of Long (1991) and Mathieson et al. (2015). One difference is that our approach directly
27 infers and models the variance of the predicted allele frequencies in the admixed population given the set of
28 surrogates used for ancestral sources. This parameter can help control for large deviations in allele frequency
29 arising solely from genetic drift experienced in the admixed population (Long 1991; Bhatia et al. 2014) and/or
30 from using inaccurate proxies for one or more of the source populations. In some applications here, e.g. the
31 Brazilian cohort, AdaptMix gives P -values with a median near 0.5 as expected under the null hypothesis of
32 neutrality (supplementary fig. S16). However, simulations under neutrality that follow a slightly different model
33 than our inference approach (see Methods), shows AdaptMix gives both an excess of high and low P -values
34 relative to the uniform distribution expected under neutrality (supplementary fig. S17). This suggests our P -
35 values are not well-calibrated, perhaps reflecting deviations from the underlying model and necessitating
36 caution when choosing thresholds for significance. One potential issue is that SNPs with low minor-allele-
37 frequency (MAF) likely well-fit their expected frequencies under the neutral model, given their lower expected
38 variance in sampling frequency. Therefore, datasets with a high proportion of such SNPs may decrease the
39 inferred variance parameter to an undesirably low value. Binning SNPs by MAF and inferring a separate
40 variance parameter for each bin may help. Here we based our significance thresholds on neutral simulations
41 tailored to each cohort, including matching for genome-wide sampled allele frequencies, and focus only on the
42 strongest association signals that resulted in low false-positive rates based on simulated neutral SNPs.
43 However, we caution that necessarily simulations are over-simplifications of complex latent demographic
44 processes, and more work is required to verify these signals.

1 Another important contribution of our test is that it can infer whether selection disproportionately affects one
2 source/surrogate pairing or affects all ancestry backgrounds equally. We assume selection affecting all ancestry
3 backgrounds indicates selection occurring post-admixture, which is more parsimonious than an alternative
4 explanation of independent selection events differentiating allele frequencies between each admixing source
5 and its surrogate. For inferred selection in a source/surrogate pairing, this can reflect selection occurring in that
6 source and/or its surrogate, possibly even following the admixture event. Post-admixture selection affecting
7 only one source may be possible in cases of selection only occurring in a particular environment that is
8 correlated with admixture fractions. For example, selection we infer to occur in Native Americans may be
9 attributable to Europeans introducing a new environmental pressure (e.g. infectious disease) that
10 disproportionately affected fitness in indigenous Americans. However, the split time between the true Native
11 American ancestral source and our Native American surrogate is likely much longer than the time since colonial
12 era admixture, suggesting selection pre-admixture as a more plausible explanation given the longer time to act.
13 Supporting this, our inferred selection coefficients (which are summed over time) in cases where we conclude
14 selection in Native Americans are typically greater than 2 (supplementary table S6). If selection had occurred
15 post-admixture continuously over the last 12 generations (corresponding to an admixture date of ~1650CE),
16 this value approximately corresponds to a per generation selection coefficient ~0.16, which is strong relative to
17 previous reports of recent selection in human populations (e.g. Hamid et al. (2021)). In contrast, our four
18 signals concluding post-admixture selection infer a per generation selection coefficient <0.1, which falls more
19 in line with previous inference of selection strengths.

20 For 18 genomic regions where we conclude selection in the Native American source (supplementary table S6),
21 it is possible this is capturing selection in (some subset of) groups that comprise the Native American surrogate
22 group we use here, rather than in the (more localized) Native American source of the admixed population. The
23 lack of overlap in selection signals when analysing the five CANDELA cohorts, and lack of concordance of our
24 signals with those from PBS testing for selection in this combined Native American surrogate (supplementary
25 fig. S18), suggests our signals are not being driven by selection in this combined population in practice. Another
26 potential concern is that our likelihood ratio test may preferentially conclude selection in the Native American
27 source if the combined Native American surrogate generally represents a poor match to the true source.
28 Encouragingly, when using PBS to test for selection in LAI-inferred Native American segments from individuals
29 with high degrees of ancestry recently related to the tested Native American source, an analysis that does not
30 use the allele frequency of the combined Native American surrogate, PBS scores for SNPs in 6 of these 18
31 regions fall into the top 99.99th percentile (supplementary fig. S19-24), with the remaining 13 regions
32 containing SNPs in the top 99th percentile. However, relative to our approach, LAI-based selection scans (e.g.,
33 Avila-Arcos et al. (2020)) may be more robust to using combined data from multiple populations to represent
34 one surrogate, since it only requires matching to a subset of individual's haplotype patterns in the reference
35 panel.

36 We also checked whether the top signals recently reported to be under selection in the Native American
37 ancestry component of an admixed Mexican population using Ohana (Cheng et al. 2021) showed evidence of
38 selection in our scan of a different Mexican cohort. Notably, we found that 7 out of the top 10 candidate genes
39 reported in Cheng et al. (2021) contained at least one nearby SNP (i.e., within 50kb from the reported gene)
40 with AdaptMix selection scores above the 95th percentile in the Mexican cohort, including 4 SNPs with scores
41 above the 99th percentile, and one SNP with a score above the 99.9th percentile. We also found that among the
42 18 SNPs classified as being selected in the Native American ancestors of the Peruvian cohort, 12 of these were
43 found at higher frequencies in ancient DNA (aDNA) from >700-year-old populations sampled in Peru relative to
44 any other aDNA data sampled elsewhere in the Americas (supplementary fig. S25).

1 In general our approach has decreased power to distinguish whether selection occurred post-admixture versus
2 in one of the ancestral sources, if reference population allele frequencies are very different and/or selection is
3 weak (fig. 1c). Inferring excess ancestry matching using LAI would likely better classify whether selection was
4 post-admixture in such cases, e.g. a scenario where one population that is fixed (or nearly-fixed) for the
5 protective allele intermixes with a population nearly-fixed for the non-protective allele, with the admixed
6 population subsequently undergoing selection. An example of this is a recently reported excess of African
7 ancestry, likely attributable to post-admixture selection, on the Duffy-null allele in inhabitants of Santiago
8 Island in Cape Verde (Hamid et al. 2021). However, our test to detect whether *any* type of selection occurred
9 should not be affected by these scenarios. In addition, our approach may identify post-admixture selection in
10 scenarios that excess-ancestry LAI-based would miss by design, such as cases where the selected allele is at a
11 similar frequency in all reference populations. Perhaps the most important contrast to LAI and other
12 approaches detecting selection in admixed populations (Cheng et al. 2021), is that in principle our approach
13 can be applied to populations that descend from the mixture of genetically similar groups, e.g. if using
14 haplotype-based approaches (e.g. SOURCEFIND) to infer ancestry proportions. Future work should assess the
15 power of this technique under such admixture settings.

16 While our method assumes a single pulse of admixture, theoretically our ability to diagnose and classify
17 selection occurring in only one source should not be affected by multiple instances of (or continuous)
18 admixture from that or any other source. This is because the signal of allele frequency deviation due to
19 selection in such cases is entirely determined by the amount of ancestry inherited from that source, and not by
20 the number of admixture pulses. In contrast, if an admixed population experiences selection and then receives
21 new migrants from one of the original admixing sources that are unaffected by this selection, e.g. later
22 European migrants to the Americas, in theory this may attenuate our ability to determine that selection
23 occurred post-admixture. However, in a simple scenario of one such additional admixture pulse, contributing
24 10-50% of DNA, the correct post-admixture selection theoretical model fits as well or better to the theoretical
25 truth than does the incorrect model concluding selection in the source that did not contribute new migrants
26 (supplementary fig. S26).

27 As noted above, and consistent with other tests comparing populations (Mathieson 2020), the choice of
28 surrogate group can make a difference in the inferred selection signals. For example, our largest signal of
29 Native American selection, at 10q22 and most strongly signalled in the 'Andes Piedmont' Peruvian subgroup,
30 disappears if replacing the 'combined Native American' surrogate group with Han Chinese (CHB from the 1KGP)
31 (supplementary fig. S10). In this case, the frequency of the putatively selected allele (rs5030938) is 67% in LAI-
32 inferred Native haplotypes in the Peruvian 'Andes Piedmont' subgroup, which is notably higher than the 38-
33 54% observed in LAI-inferred Native American haplotypes in four non-Peruvian sub-groups, and thus consistent
34 with selection (supplementary table S7). However, it is lower than that of CHB (~76%), which explains the lack
35 of signal when using CHB as a surrogate. The frequency in Yakut, a Siberian group that perhaps better
36 represents ancestral Native Americans than CHB does (Wang et al. 2007), is closer to that of frequency
37 estimates across non-Peruvian Native American groups (0.46-0.5). In general, there is a trade-off between
38 using surrogates more distantly related to the source, which may decrease power to find regional adaptation
39 signals, versus choosing a more closely related surrogate, which may also decrease power by masking
40 adaptation signatures that it shares with the target source (e.g. using Iberians as a surrogate for European
41 ancestry of Latin Americans). Our inferred variance parameter can be used to investigate how well a given
42 surrogate captures genetic variation in the target population, with for example the inferred variance using CHB
43 as a surrogate ~5-10-fold higher relative to using the combined Native American surrogate.

44 ***Selection signals detected in the CANDELA cohort***

1 The candidate genes we infer to be affected by selection in Latin Americans and their Native American
2 ancestors are best viewed in the context of other previously reported signals. Reynolds et al. (2019) recently
3 performed a selection scan in three Native North American populations and identified some of the strongest
4 signals at immune-related genes including the interleukin 1 receptor Type 1 (*IL1R1*) gene in a sample from
5 several closely related communities in the southeastern United States, and the mucin 19 (*MUC19*) gene in a
6 central Mexican population. We do not replicate the MUC19 signal in the CANDELA Mexican cohort, which
7 could indicate that the Native American component in this cohort is not closely related to that of the central
8 Mexican Native American group. Nonetheless, we found some of our strongest signals of selection at several
9 loci encompassing genes involved in the immune response, including *CD300LF* and *MIF*, detected as being
10 selected in the Native American ancestors of Peruvians. Interestingly, *CD300LF* promotes macrophage-
11 mediated efferocytosis, while *MIF* play a role regulating macrophage function through the suppression of
12 glucocorticoids. These observations suggest that these two genes might have perhaps evolved in a coordinated
13 manner, possibly due to their phagocytic-related role against the novel pathogens encountered in the
14 Americas.

15 Regarding signals of selection post-admixture, several studies have consistently shown adaptive signals in
16 different Latin American populations at HLA by showing an excess of matching to African reference haplotypes
17 using LAI (Tang et al. 2007; Basu et al. 2008; Ettinger et al. 2009; Guan 2014; Rishishwar et al. 2015; Deng et al.
18 2016; Zhou et al. 2016; Norris et al. 2020; Vicuna et al. 2020). Given that African ancestry was enriched at this
19 region, the authors suggested that certain African alleles could have conferred a selective advantage to certain
20 infectious diseases most likely brought by Europeans. While AdaptMix is only able to classify selection in one
21 cohort (Colombia) out of our four HLA signals, we also replicated this excess of African ancestry in each of the
22 CANDELA cohorts (supplementary fig. S12). There is some debate as to whether these signals are genuine or
23 attributable to confounders such as inaccurate LAI inference (Pasaniuc et al. 2013). To illustrate the validity of
24 these concerns, people with entirely Northwest European ancestry from Britain infer excess ancestry related to
25 Africa in HLA, which – though perhaps influenced by genuine selection at HLA in Northwest Europeans –
26 presumably does not reflect genuine recent African ancestry (supplementary fig. S27). Instead this is more
27 likely attributable to the relatively high degree of genetic diversity in HLA mimicking African genetic diversity,
28 illustrating how these LAI-based tests can give false-positive signals when testing for post-admixture selection.
29 This may explain why AdaptMix does not replicate the moderate amount of excess African ancestry inferred by
30 LAI at HLA in the Brazilian cohort (supplementary fig. S12), which is predominantly of European ancestry.
31 Indeed regions under selection in admixed populations may be particularly difficult to classify accurately using
32 LAI, e.g. with the HLA region here having the lowest overall LAI classification probability (supplementary fig.
33 S28), especially in cases where the reference population have not experienced similar selection and hence may
34 have poorly matching genetic variation patterns. As our approach does not require LAI, it is robust to these
35 issues. While our model is not able to classify selection as post-admixture at most of our HLA signals, allele
36 frequency patterns in the admixed cohorts are consistent with post-admixture selection and often show allele
37 frequencies drifting away from those expected under our neutral model and towards those of the African or
38 European reference population (supplementary fig. S29). This is most evident in the Colombian cohort,
39 consistent with Africans contributing protective alleles as previously suggested (Tang et al. 2007; Basu et al.
40 2008; Ettinger et al. 2009; Guan 2014; Rishishwar et al. 2015; Deng et al. 2016; Zhou et al. 2016; Norris et al.
41 2020; Vicuna et al. 2020). In addition to HLA, we also identified a novel post-admixture selection signal in the
42 Chilean cohort that was accompanied by a significant increase of European ancestry at the *CD101* locus, again,
43 suggesting that protective alleles from Europeans might have also been adaptive to counter Old World-borne
44 diseases brought to the Americas.

1 The signals encompassing genes related to metabolic and anthropometric-related phenotypes are consistent
2 with novel dietary practices in the Americas driving adaptation, with many signals with an effect on relevant
3 phenotypes and/or tissues, classified as being selected in the Native American source. Previous studies have
4 shown evidence of adaptation at genes related to metabolic-related phenotypes and attributed the adaptation
5 to dietary pressures in Native Americans. Avila-Arcos et al. (2020) recently reported strong signals of selection
6 in the Mexican Huichol at several genes associated to lipid metabolism, including *APOA5* and *ABCG5*. We do
7 not replicate these signals in the CANDELA Mexican cohort, which could indicate that the Native American
8 component in this cohort is not closely related to that of the Huichol. The signals at *APOA5* and *ABCG5* are in
9 line with a previous finding of a strong selection signal at another ATP-binding cassette transporter A1 (*ABCA1*)
10 gene that has been associated with low high-density lipoprotein cholesterol in Latin Americans (Villarreal-
11 Molina et al. 2008; Acuña-Alonzo et al. 2010). As the *ABCA1* protein carrying the putative selected allele shows
12 a decrease cholesterol efflux, the authors suggest that this variant could have favored intracellular cholesterol
13 and energy storage, which in turn might have beneficially influenced the ability to accommodate fluctuations in
14 energy supply during severe famines and during the regulation of reproductive function (Acuña-Alonzo et al.
15 2010). Lindo et al. (2018) used a genomic transect of Andean highlanders from northern Peru, and found the
16 strongest signals of selection at *MGAM*, a gene related to starch digestion. The authors attributed this finding
17 to a dietary-related selective pressure perhaps brought by the transition to agriculture in this region. AdaptMix
18 shows evidence in the CANDELA Peruvian cohort within *MGAM* (rs7810984, AdaptMix P -value= 1.79×10^{-8} ,
19 above 99.9th percentile) only when using CHB as a surrogate for Native American ancestry. This again illustrates
20 how the choice of surrogate populations defines the selection test between each surrogate and its
21 corresponding ancestral source. It is possible that by including Andean Native Americans in our Native
22 American source population (supplementary table S1) we are affecting the power to detect selection in the
23 Andean Native American ancestors of the CANDELA Peruvian cohort, analogous to how Lindo et al. (2018) no
24 longer detect selection at *MGAM* if using PBS to compare ancient and present-day (Aymara) Andean groups.
25 Studies have also reported signals of selection in Native Americans groups shared with Siberian populations,
26 which the authors interpreted as an adaptation to polyunsaturated-rich diets prior or close to the peopling of
27 the Americas, likely in the Arctic Beringia. These included a signal overlapping the *WARS2* and *TBX15* genes,
28 previously associated to body fat distribution and adipose tissue differentiation (Fumagalli et al. 2015; Racimo,
29 Gokhman, et al. 2017), and the fatty acid desaturase (*FADS*) gene cluster that modulates the manufacture of
30 polyunsaturated fatty acids (Amorim et al. 2017; Harris et al. 2019) (but see Mathieson (2020) for an
31 alternative explanation of the *FADS* signal). Again, we inferred moderate selection evidence at these regions in
32 the CANDELA Peruvian cohort only when using CHB as surrogate for Native American ancestry (SNP rs2361028
33 near *TBX15*, AdaptMix P -value= 1.8×10^{-7} , above 99.5th percentile; SNP rs174576 within *FADS2*, AdaptMix P -
34 value= 3.8×10^{-8} , above 99.5th percentile). It is thus tempting to suggest that the three novel signals of selection
35 AdaptMix classifies as being under selection in Native Americans might be related to dietary pressures
36 experienced prior or during the peopling of the Americas (e.g., the *BRINP3* signal detected in Mexicans), or as a
37 product for a greater reliance of domesticated crops including potatoes (3400–1,600 CE) (Rumold and
38 Aldenderfer 2016) (e.g., the *HKDC1* signal detected in Peruvians). However, it is important to note that other
39 factors may also be attributable for some of these selection signals.

40 Of potential adaptive interest is the *STOX1* gene detected in the Peruvian cohort close to our highest overall
41 selection signal within *HKDC1* at 10q22 (fig. 6b). Mutations within *STOX1* have been associated to preeclampsia
42 (Van Dijk et al. 2005; van Dijk and Oudejans 2011), a pathology of pregnancy characterized by high blood
43 pressure and signs of damage to other organ system that can be lethal for the mother and for the fetus (Sibai
44 2003). Interestingly, in the single linkage study on preeclampsia conducted in Andean Peruvian families to date,

1 SNPs within *STOX1* show marginal association (P -value=0.004678) (supplementary fig. S30) (Badillo Rivera et al.
2 2021). Given that high elevation is linked to an increased incidence of preeclampsia (Zamudio 2007), it is
3 possible that natural selection has acted on genes related to this condition. Furthermore, the fact that variants
4 within *HKDC1* are associated with glucose levels in pregnant women (Hayes et al. 2013; Guo et al. 2015;
5 Kanthimathi et al. 2016; Tan et al. 2019) and considering the relationship between abnormal glucose levels and
6 preeclampsia (Joffe et al. 1998; Weissgerber and Mudd 2015), it is also possible that natural selection has
7 targeted variants at *HKDC1* due to its role in glucose metabolism.

8 Lastly, other environmental factors may also be attributable for some of these selection signals, such as
9 infectious diseases. There is growing evidence of a link between metabolic diseases and innate immunity or
10 inflammation (Pickup and Crook 1998; Kominsky et al. 2010; Lumeng and Saltiel 2011; Robbins et al. 2014). For
11 instance, it has been shown that cholesterol plays a key role in various infectious processes such as the entry
12 and replication of flaviviral infection (Osuna-Ramos et al. 2018). Additional studies in ancient and present-day
13 indigenous American populations will be needed to disentangle the putative selective pressures at these loci.

14 **Conclusion**

15 We have presented a novel allele frequency-based method that identifies loci under selection in admixed
16 populations, while determining whether the selection affected all ancestral sources equally, indicating
17 selection following admixture, or in only one of the sources. The novel candidate genes under selection provide
18 new insights into the adaptive traits necessary for the early habitation of the Americas and to respond to the
19 challenge of infectious pathogens corresponding to European contact. Future functional investigations will
20 allow a more detailed understanding of the consequences of selective pressures experienced in the American
21 continent, including its effect on present-day health outcomes.

22 **Materials and Methods**

23 **Genomic datasets**

24 The Latin American individual samples analyzed here were part of CANDELA Consortium (Ruiz-Linares et al.
25 2014). The CANDELA Consortium samples (<http://www.ucl.ac.uk/silva/candela>) have been described in detail
26 in previous publications (Ruiz Linares et al 2014; Chacon-Duque et al., 2018). These data include a total of 6,630
27 volunteers from five Latin American countries (Brazil, Chile, Colombia, Mexico and Peru). This dataset was
28 genotyped on the Illumina HumanOmniExpress chip platform including 730,525 SNPs. We also collated
29 reference populations from regions that have contributed to the admixture in Latin America. For Native
30 American samples we used individuals previously genotyped by Chacon-Duque et al. (2018). This dataset
31 comprises 19 Native American populations from throughout the Americas with genotype data
32 (supplementary table S1). For all the analyses described, we have only retained Native American individuals
33 that showed more than 99% Native American ancestry as estimated by ADMIXTURE (see below). For European
34 samples, we used genotype data from Portuguese and Spanish, individuals previously genotyped by Chacon-
35 Duque et al. (2018) and Spanish (IBS; Iberian Population in Spain) from the 1000 Genomes Project study (The
36 1000 Genomes Project Consortium 2015). For Sub-Saharan Africans, we used genotype data from Yoruba (YRI;
37 Yoruba in Ibadan, Nigeria), and Luhya (LWK; Luhya in Webuye, Kenya) individuals from the 1KGP. The reference
38 samples from Chacon-Duque et al. (2018) are described in more detail in the Supplementary Table 1 from the
39 mentioned publication. For some of our analysis we also included the 103 Han Chinese from Beijing (CHB) and
40 85 Europeans from England and Scotland (GBR) from the 1KGP as a surrogate for the Native American and
41 European source, respectively. Genotype data of the individuals from the 1KGP was downloaded from the 1000
42 Genomes Project FTP site available at <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/>.

1 **Data curation**

2 We used PLINK v1.9 (Chang et al. 2015) to exclude SNPs and individuals with more than 5% missing data or that
3 showed evidence of genetic relatedness as in Chacon-Duque et al. (2018). Due to the admixed nature of the
4 Latin American samples, there is an inflation in Hardy-Weinberg P -values, and therefore we did not exclude
5 SNPs based on Hardy-Weinberg deviation. After applying these filters, 625,787 autosomal SNPs and 7,986
6 individuals were retained for further analysis.

7 **Selecting admixed Latin American and reference individuals**

8 In order to select admixed Latin American individuals (i.e. individuals with varying degrees of Native American,
9 European and African ancestry), we conducted an unsupervised ADMIXTURE analysis at $K=3$ using a set of
10 103,426 LD-pruned SNPs including Native Americans, Iberian Europeans and West Africans. We then removed
11 non-admixed Latin American individuals that we define as having less than 10% or more than 90% Native
12 American genome-wide ancestry. To avoid confounding our selection inference due to underlying population
13 structure, we further excluded individuals with $>1\%$ inferred ancestry matching to surrogates other than Native
14 Americans, Iberian Europeans, and West Africans using SOURCEFIND estimates obtained for the same
15 individuals in Chacon-Duque et al. (2018). As expected, we observe a strong correlation between the
16 ADMIXTURE and SOURCEFIND estimates (average $r>0.99$) demonstrating the validity of this filtering approach
17 and demonstrating that most of the ancestry of the admixed Latin American individuals can be appropriately
18 modelled by a three-way admixture model. After this filtering procedure, the five Latin American populations
19 consisted of 190 Brazilians (BRA), 1125 Colombians (COL), 896 Chileans (CHL), 773 Mexicans (MEX) and 834
20 Peruvians (PER). From our Native American, European, and Sub-Saharan African reference populations, we also
21 removed individuals that contained more than 1% of ancestry from another group based on the ADMIXTURE
22 analysis described above. After this extra filter our final reference dataset was composed of 142 Native
23 Americans, 205 Europeans, and 206 Sub-Saharan Africans.

24 **Change in allele frequency under Wright-Fisher with multiplicative model of selection**

25 Assuming a multiplicative model of selection and random mating, the frequency of the three
26 genotypes in generation 1 at a biallelic locus with alleles A and a at frequencies p and $1 - p$,
27 respectively, in the previous generation is:

AA	Aa	aa
$(1 + s_1)^2 p^2 / c_1$	$(1 + s_1) 2p(1 - p) / c_1$	$(1 - p)^2 / c_1$

28 where $s_1 \in [-1, \infty]$ is the selection coefficient in generation 1 and $c_1 = (1 + s_1)^2 p^2 + (1 + s_1) 2p(1 - p) +$
29 $(1 - p)^2$. Note that each copy of the A allele changes fitness by a factor of $(1 + s_1)$.

30 Under the above, the allele frequency of (p_1) of allele A in generation 1 is:

$$31 \quad p_1 = \frac{(1 + s_1)^2 p^2 + (1 + s_1) p(1 - p)}{(1 + s_1)^2 p^2 + (1 + s_1) 2p(1 - p) + (1 - p)^2} \quad (1)$$

$$= \frac{(1 + s_1) p}{1 + s_1 p}$$

32 For generation 2, again assuming a multiplicative selection, the frequencies of the three genotypes are:

AA	Aa	aa
$(1 + s_2)^2 p_1^2 / c_2$	$(1 + s_2) 2p_1(1 - p_1) / c_2$	$(1 - p_1)^2 / c_2$

33 Where $s_2 \in [-1, \infty]$ is the selection coefficient in generation 2 and $c_2 = (1 + s_2)^2 p_1^2 + (1 + s_2) 2p_1(1 -$
34 $p_1) + (1 - p_1)^2$. Note that each copy of the allele A changes fitness by a factor of $(1 + s_2)$ in this generation.
35

1 The allele frequency (p_2) of allele A in generation 2 is:

$$\begin{aligned} p_2 &= \frac{(1 + s_2)p_1}{1 + s_2 p_1} & (2) \\ &= \frac{(1 + s_2) \left[\frac{(1 + s_1)p}{1 + s_1 p} \right]}{1 + s_2 \left[\frac{(1 + s_1)p}{1 + s_1 p} \right]} \\ &= \frac{(1 + s_2^*)p}{1 + s_2^* p}, \end{aligned}$$

2

3 where $s_2^* \equiv (s_1 + s_2 + s_1 s_2)$.

4 More generally, the allele frequency p_g of allele A in generation g is:

$$p_g = \frac{(1 + s)p}{1 + sp}, \quad (3)$$

5

6 where

$$s = [\sum_{i=1}^g s_i] + [\sum_{j=1}^{g-1} (s_j \sum_{i=j+1}^g s_i)] + \sum_{i=3}^g \Pi_i \approx \sum_{i=1}^g s_i, \quad (4)$$

7 with s_i the selection coefficient at generation i and Π_i the sum of the products of all $\binom{g}{i}$ combinations of
8 $\{s_1, \dots, s_i\}$ values. The approximation in equation (4) assumes the s_i are small, which should be a reasonable
9 approximation based on e.g. estimated selection coefficients in humans.

10 **Testing for evidence of selection at a SNP**

11 To assess the evidence of selection at a SNP, we employ a model inspired by that used in Mathieson et al.
12 (2015) and based on the Balding-Nichols formulation (Balding and Nichols 1995). In particular for the allele
13 count X_j at SNP j in the target population, we assume:

$$Pr(X_j = x_j | M, p_j, D) = \text{Beta-Binomial} \left(x_j; 2M, \frac{1-D}{D} p_j, \frac{1-D}{D} (1-p_j) \right), \quad (5)$$

14

15 where M is the number of target individuals and D is a variance parameter that is measuring the degree of
16 uncertainty about p_j . More generally, D can be thought of as genetic drift parameter. The above model
17 implicitly assumes that the frequency of the allele in the target population follows a $\text{Beta}(\text{mean} =$
18 $p_j, \text{variance} = D p_j (1 - p_j))$. Under neutrality, we assume

$$p_j = \frac{1}{M} \sum_{k=1}^K \left(\left[\sum_{i=1}^M \alpha_k(i) \right] f_{jk} \right), \quad (6)$$

19

20 where f_{jk} is the sampled frequency of the allele in the surrogate population at SNP j for source k , and $\alpha_k(i)$ is
21 the inferred admixture proportion from population k in individual i . We first find \hat{D} as the value of D that
22 maximizes $\prod_{j=1}^J [Pr(X_j | M, p_j, D)]$, using the optim function in R with the 'Nelder-Mead' algorithm. Then, fixing
23 $D = \hat{D}$ in equation (5), for each SNP we find the two-sided P -value testing the null hypothesis that the
24 observed allele counts follow this neutral model.

25 The variance under (5) is small for SNPs with very high or very low p_j , so such SNPs tend to reject this null
26 model even in cases where the observed target population allele frequency does not deviate notably from its
27 neutral expectation p_j in (6). Therefore, we used an alternative parameterisation where we assumed the
28 frequency of the allele in the target population follows a $\text{Beta}(\text{mean} = p_j, \text{variance} = V)$. This was achieved

1 by substituting D in equation (5) at SNP j with $\min\left[\frac{V}{p_j(1-p_j)}, 0.99999\right]$, necessary to ensure numerical
 2 stability, and finding \hat{V} . In practice this means that SNPs with minor allele frequency $< (1.00001 \times V)$ had
 3 variance $(0.99999 p_j(1-p_j))$ rather than V . While our use of V achieved the desired property of mitigating
 4 false-positives at SNPs with low MAF, one potential drawback is that datasets containing a high proportion of
 5 low-MAF SNPs may drive the inferred V to be small relative to the variance expected at high-MAF SNPs under
 6 neutrality. In other words, under neutrality it is possible that $V > D p_j(1-p_j)$ at low-MAF SNPs, yet
 7 $V < D p_j(1-p_j)$ at high-MAF SNPs. If so, high-MAF SNPs may reject the neutral model more frequently than it
 8 should under neutrality. Indeed, this seems to be the case: in some of our neutral simulations described below,
 9 SNPs with *AdaptMix* P -value < 0.05 are 1.7-fold enriched for SNPs with MAF > 0.3 relative to all tested SNPs. We
 10 reiterate this is partially by design since we use our formulation with V precisely to avoid inferring selection at
 11 low-MAF SNPs. Future work, for example inferring V separately for SNPs binned by MAF, may lead to better P -
 12 value calibration under neutrality.

13 **Determining whether selection occurred pre or post-admixture**

14 Consider the scenario in fig. 1b, where sampled population C descends from an admixture of unsampled
 15 populations A^* and B^* , who are represented by sampled surrogate populations A and B , respectively. Our
 16 test aims to distinguish whether selection occurred post-admixture along branch (e) versus along any of
 17 branches (a)-(d). Let f_C be the allele frequency of a sample from population C . At a neutral SNP:

$$18 \quad E[f_C] = \alpha f_{A^*} + (1 - \alpha) f_{B^*}, \quad (7)$$

19 where f_{A^*} and f_{B^*} are true allele frequencies of A^* and B^* at the SNP, respectively, and α is the admixture
 20 proportion from A^* . Letting f_k be the sampled allele frequency for population k serving as surrogate to the
 21 true admixing population k^* , it seems reasonable to assume:

$$22 \quad E[f_C] = \alpha f_A + (1 - \alpha) f_B. \quad (8)$$

23 Note that this also holds under selection along branch (f) in fig. 1b, which we ignore here (but which can be
 24 tested by comparing allele frequencies in A and B). Equation (8) assumes that f_A and f_B are equally good
 25 proxies for the admixing populations' frequencies f_{A^*} and f_{B^*} , respectively, at the SNP, which may not be true.
 26 We test the effect of this using simulations, described below, in which surrogates vary in how well they reflect
 27 their respective true admixing
 28 sources.

29 In the case of a multiplicative model of selection along branch (e) in fig. 1b at this SNP, using equation (3) we
 30 assume:

$$31 \quad E[f_C] = \frac{(1+s)[\alpha f_A + (1-\alpha)f_B]}{1+s[\alpha f_A + (1-\alpha)f_B]} \equiv E_c[f_C], \quad (9)$$

32 where s is the selection strength (i.e. equation [4]) along branch (e).
 33 Alternatively, under a multiplicative model for selection along branches (a) and/or (c) in fig. 1b, with analogous
 34 results for selection along branches (d) and/or (b), instead we assume:

$$35 \quad E[f_C] = \alpha \left[\frac{(1+s_A)f_A}{1+s_A f_A} \right] + (1-\alpha)f_B = f_B + \alpha \left[\frac{(1+s_A)f_A}{1+s_A f_A} - f_B \right] \equiv E_A[f_C], \quad (10)$$

1 where s_A is the selection strength along branches (a) and/or (c). Importantly, $E_A[f_c]$ is linear in α , while $E_C[f_c]$,
 2 is not, which we aim to exploit to distinguish between these two scenarios.
 3 Here, assuming CANDELA population T can be described as a mixture of K sources, we assume the genotype g_i
 4 of individual $i \in [1, \dots, M]$ from T follows:

$$g_i \sim \text{Binomial}(2, f_T(i)). \quad (11)$$

5
 6 Under neutrality, we set $f_T(i)$ in (11) to:

$$f_T^N(i) = \sum_{k=1}^K [\alpha_k(i) f_k], \quad (12)$$

7
 8 where f_k is the sampled allele frequency at the given SNP for the surrogate population to the source
 9 contributing $\alpha_k(i)$ admixture to individual i .

10 In the case of selection in T post-admixture, we generalize equation (9) and set $f_T(i)$ in (11) to:

$$f_T^P(i|s) = \frac{(1+s)[\sum_{k=1}^K \alpha_k(i) f_k]}{1+s[\sum_{k=1}^K \alpha_k(i) f_k]}. \quad (13)$$

11
 12 For the alternative case of selection along the branches separating source A and its sampled surrogate A^* , we
 13 generalize equation (10) and replace $f_T(i)$ in (11) with:

$$f_T^A(i|s_A) = \left[\sum_{k \neq A}^K \alpha_A(i) f_k \right] + \alpha_A(i) \left[\frac{(1+s_A) f_A}{1+s_A f_A} \right]. \quad (14)$$

14
 15 In practice, we fix $\alpha_A(i)$ to be the proportion of DNA that each target individual i matches to surrogate k as
 16 inferred by ADMIXTURE. We define:

$$L^P(s) \equiv \prod_{i=1}^M \left[f_T^P(i|s)^{g_i} (1 - f_T^P(i|s))^{2-g_i} \right], \quad (15)$$

17
 18 where g_i is the genotype for target individual i . We use the optim function in R with the ‘Nelder-Mead’
 19 algorithm to find the maximum-likelihood estimate (MLE) \hat{s} , which is the value of s that maximizes equation
 20 (15).

21 Similarly we define:

$$L^A(s_A) \equiv \prod_{i=1}^M \left[f_T^A(i|s_A)^{g_i} (1 - f_T^A(i|s_A))^{2-g_i} \right], \quad (16)$$

22
 23 again finding \hat{s}_A , as the MLE for s_A .

24 We note that $[2 - 2\log(L^P(\hat{s}))]$ and $[2 - 2\log(L^A(\hat{s}_A))]$ are analogous to AIC values for these respective
 25 models. Following AIC theory, we calculate:

$$I = \frac{\min[L^P(\hat{s}), L^A(\hat{s}_A)]}{\max[L^P(\hat{s}), L^A(\hat{s}_A)]} \leq 1, \quad (17)$$

26
 27 where, relative to the model with higher likelihood out of (15) and (16), the model with smaller likelihood is I
 28 times as probable to minimise the loss of information when used to represent the unknown true model (Akaike
 29 1974).

1 Note we could analogously calculate the likelihood under the neutral model, i.e., using equation (12). Then, as
 2 an alternative to the selection testing approach described in Section ‘Testing for evidence of selection at a
 3 SNP’, we could use a likelihood-ratio-statistic approach to test for selection using either (15) or (16) as the
 4 alternative model likelihood. We explored this alternative testing approach, but do not use it here because it
 5 gave lower P -values when simulating under neutrality. This observation may in part be alleviated if we
 6 estimated f_{k^*} under both the neutral and alternative models rather than fixing $f_{k^*} = f_k$. However, estimating
 7 f_{k^*} is confounded with estimating s or s_A under the alternative models.

8 Simulations

9 *Estimating how well each surrogate reflects its corresponding true admixing source*

10 We aimed to generate simulations that mimic our real data. To do so, we first generate a measure of how well
 11 a sampled surrogate population k reflects its corresponding true (unknown) source population. In particular,
 12 we estimate a drift parameter d_k in the following manner. First, at each SNP j we use `nlminb` in R to find the
 13 estimated values $\{\tilde{f}_1^j, \dots, \tilde{f}_K^j\}$ for $\{f_{1^*}, \dots, f_{K^*}\}$, respectively, that minimize:

$$\sum_{i=1}^M \left(x_i^j - \sum_{k=1}^K \alpha_k(i) f_{k^*} \right)^2, \quad (18)$$

14
 15 Where $x_i^j \in \{0,1,2\}$ is the allele count for the admixed target individual $i \in [1, \dots, M]$ at the SNP and each
 16 $\tilde{f}_k^j \in [0,1]$. Then, for each source k , with observed allele counts G_k^j and total counts M_k^j at SNP j in the
 17 surrogate population, following Balding-Nichols (Balding and Nichols 1995) we assume:

$$G_k^j \text{Beta} \sim \text{Binomial} \left(M_k^j \frac{d_k}{1-d_k} \tilde{f}_k^j, \frac{d_k}{1-d_k} (1 - \tilde{f}_k^j) \right). \quad (19)$$

18
 19 We then used the ‘Nelder-Mead’ algorithm in the `optim` function in R to find the $d_k \in [0,1]$ that maximized the
 20 product of (19) across all SNPs. This gave the values reported in Table 1.

21 Large estimated d_k (>0.1) correspond to cases where there is little admixture from that source in our sampled
 22 individuals from that country, i.e. for African admixture in most countries and Native American admixture in
 23 Brazil. As values inferred using such little data are presumably unreliable, we cap them at 0.05 for the
 24 simulations below. While these values are a guide, in practice we adjusted these values by a multiple of 2-7 to
 25 generate neutral simulations that had the same inferred drift \hat{D} , described in section ‘Testing for evidence of

Target	Native American	European	African
Brazil	0.173	0.007	0.102
Chile	0.02	0.011	0.226
Colombia	0.044	0.012	0.044
Mexico	0.024	0.007	0.223
Peru	0.015	0.009	0.119

26 selection at a SNP’, as that observed in the real data.

27 **Table 1.** Inferred d_k measuring how well the sampled surrogate (column) reflect the true admixing sources for
 28 each target population (row).

29 *Generating simulated allele frequencies*

1 We simulated admixed individuals who had experienced selection, with genome-wide admixture proportions
 2 $\alpha_k(i)$ from source populations $k \in [1, \dots, K]$ for simulated individuals $i \in [1, \dots, M]$ matching those inferred by
 3 ADMIXTURE in the real data. To do so, for each simulation we repeated the following procedure:

- 4 1. For each source k , at each SNP we sample starting allele frequencies f_k^* from a
 5 $Beta\left(\frac{d_k}{1-d_k} f_k, \frac{d_k}{1-d_k} (1 - f_k)\right)$, where f_k is the sampled frequency of the respective surrogate
 6 population and d_k are defined in Table 1 (but capped at 0.05).
- 7 2. We randomly select SNPs to undergo selection. If selection is occurring in source population k
 8 prior to admixture, we randomly sample from among SNPs for which $f_k^* < 0.5$. If selection is
 9 occurring post-admixture, we instead randomly sample from among SNPs for which
 10 $\sum_{i=1}^M (\sum_{k=1}^K f_k^* \alpha_k(i)) / M < 0.5$.
- 11 3. We randomly select neutral SNPs from among all remaining SNPs, i.e., those not among the
 12 SNPs chosen in (2), in the real data.
- 13 4. To simulate selection:
 - 14 • If selection is occurring prior to admixture, we simulate selection in the relevant source
 15 population for g generations under a specified model of selection (additive, dominant,
 16 multiplicative, recessive) using Wright-Fisher with a population size of N_e individuals.
 - 17 • If selection is occurring after admixture, we simulate selection separately in each of the
 18 source populations for g generations, under a specified model of selection using
 19 Wright-Fisher with a population size of N_e individuals per population.
- 20 5. At each SNP, we sample allele counts for each individual i from a $Binomial(2, p_i)$ with
 21 $p_i = \sum_{k=1}^K [f_k^g \alpha_k(i)]$, where:
 - 22 • $f_k^g = f_k^*$ for neutral SNPs
 - 23 • $f_k^g = f_k^*$ at selected SNPs for source populations k not undergoing selection (i.e., in
 24 cases where selection is pre-admixture)
 - 25 • f_k^g is the sampled final frequencies in step (4) after g generations, at selected SNPs for
 26 source population k undergoing selection

27
 28 We then analyse data from the simulated target population individuals using the real sampled data from the
 29 surrogate populations. For simulations here, we use $N_e = 10000$ for the African, European, and Native
 30 American source groups.

31 Our procedure in steps (4)-(5) to simulate selection and admixture ensures the admixed individuals have
 32 variable admixture proportions while remaining computationally tractable. An alternative to this would be to
 33 generate M admixed populations using observed f_k values, with the admixture proportions for population i
 34 equal to $\alpha_1(i), \dots, \alpha_K(i)$, and then simulate each admixed population for g generations using Wright-Fisher,
 35 either with or without selection. Such simulations would match the approach used by our model to classify
 36 selection as type (i) or type (ii) (Section ‘Determining whether selection occurred pre- or post-admixture’).
 37 However, we chose the above for reasons of computational efficiency, as we have many individuals (i.e.,
 38 $M > 1000$). Note also that our selection test (Section ‘Determining whether selection occurred pre- or post-
 39 admixture’) is different from this simulation procedure, in that our test models the combined allele frequency
 40 across all admixed individuals, using the mean admixture contributions across target individuals to calculate
 41 the expected frequency. This, in addition to the way we infer the variance term that describes the distribution
 42 of each SNP’s sampled allele frequency (see ‘Testing for evidence of selection at a SNP’ above), may explain
 43 why our model exhibits an excess of SNPs with small P -values even when simulating no selection. This is

1 despite using all SNPs to infer the model's variance parameter, which is designed to make more SNPs fit the
2 model (likely explaining the excess of high P -values we also see, e.g., in supplementary fig. S17). While
3 including this variance parameter does somewhat control P -values by e.g., giving in some cases a median P -
4 value near 0.5, as expected under neutrality, our no-selection simulations suggest caution in directly using our
5 model's P -values for assessing selection evidence. This suggests some degree of plausible simulations would be
6 helpful to calibrate the model's reported P -values.

7 **Forward simulations**

8 To explore the effect of the effective population size (N_e) on the population undergoing selection, we
9 conducted additional simulations using the forward simulator SLiM 3 (Haller and Messer 2019). We used the
10 demographic model of an admixed Mexican (MEX) population recently presented in Cheng et al. (2021). The
11 model is based on parameter estimates from Gravel et al. (2011) and Gutenkunst et al. (2009) of a three-
12 population demography including an African (AFR), European (EUR), and Asian population (ASN). The main
13 difference is the inclusion of an additional Native American population that splits from the ASN population. The
14 MEX population is modelled as a 50%/50% admixture between EUR and the Native American population. We
15 consider five different N_e 's for the Native American population ($N_e=800, 1000, 200, 5000, \text{ and } 10000$). The
16 selection occurs only in the ancestral Native American population with no ongoing selection in MEX. In the
17 original model selection lasts for 500 generations, which resulted in the allele being fixed before the admixture
18 event in simulations with high N_e , particularly when testing for high selection coefficients. To avoid this fixation
19 which might result in a bias when estimating the power, we modelled selection in the Native American
20 population for 300 generations. All other parameters were the same as in the original model. We simulate a
21 region of 4Mb with a mutation rate of 10^{-8} and a recombination rate of 10^{-8} base pairs per generation and
22 sample 20 diploid individuals from each population. We simulate a single selected site under an additive model
23 within a ± 10 kb window of the center of the simulated region. As in our previous simulations we consider 10
24 different selection coefficients ($s=0.01$ to 0.1 in steps of 0.01 , with s defined here as the increased fitness when
25 carrying one copy of the advantageous allele) with a starting frequency for the selected site being equal to or
26 higher than 0.01 but lower than 0.1 . Following Haller and Messer (2019), we scale times down by a factor of 10 ,
27 and scale up the mutation rate, recombination rate and selection coefficient by a factor of 10 . We conducted a
28 total of 500 independent regions to estimate the statistical power for each combination of N_e and selection
29 coefficient.

30 We additionally simulated an 80 Mb region under neutrality (i.e., $s=0$) using the same settings as previously
31 described. For AdaptMix, admixture proportions were estimated by applying supervised ADMIXTURE with $K=3$
32 to this neutral region, setting AFR, EUR, and ASN as the reference populations. Note that, as the MEX
33 population does not have AFR ancestry, this simulation setting is also assessing the power under a model
34 misspecification, which might be more realistic for most real-world applications. The 80 Mb neutral region was
35 then used to generate a null distribution of P -values. The power of AdaptMix was based on a P -value cutoff
36 that resulted in a false-positive rate of 5×10^{-5} of this null distribution.

37 **Comparison of AdaptMix against other selection approaches**

38 We compared the performance of AdaptMix to two different approaches under the two scenarios: (i) selection
39 in one of the source populations and (ii) selection in the admixed population following the admixture event.

40 To assess the power under scenario (i), we compared AdaptMix against Ohana, a maximum likelihood method
41 for finding regions under positive selection by modeling ancestry components (Cheng et al 2021). Importantly,
42 Ohana has been shown to retain similar or higher power compared to other state-of-the-art methods. We
43 compared AdaptMix and Ohana under the demographic setting previously described, but simulating selection

1 for 500 generations and fixing the N_e of the Native American population undergoing selection to 800, as in the
2 original publication.

3 To assess the power under scenario (ii), we compared AdaptMix against Ohana and to a local ancestry
4 deviation (LAD) approach. A LAD approach here consists of evaluating whether a genomic region is enriched for
5 a particular ancestry compared to their genome-wide average, and relies on local ancestry inference. LAD
6 approaches have been extensively used to detect signals of selection following an admixture event in several
7 recently admixed populations, including Latin Americans (Tang et al. 2007; Basu et al. 2008; Ettinger et al.
8 2009; Guan 2014; Rishishwar et al. 2015; Deng et al. 2016; Zhou et al. 2016; Norris et al. 2020; Vicuna et al.
9 2020). For this scenario we used the demographic model recently presented in Cuadros-Espinoza et al. (2021),
10 which involves a simple two-way admixture model. Briefly, the demographic model consists of two populations
11 that split from their common ancestor 2080 generations ago, and then intermix 70 generations ago to produce
12 a third admixed population. The admixture proportions are set to 50%/50%, and selection occurs only in the
13 admixed population for 70 generations until the present. All other parameters are set to those presented in the
14 publication, except for the removal of background selection. We sample 50 diploid individuals from each of the
15 three populations, i.e. the admixed population and the two intermixing populations X and Y, at the end of the
16 simulation, as in the original publication.

17 As in our previous simulations, we additionally simulated an 80 Mb region under neutrality (i.e., $s=0$) using the
18 same settings as previously described for each scenario. In the case of AdaptMix the 80 Mb neutral region was
19 used to estimate admixture proportions, based on a supervised ADMIXTURE analysis with $K=2$, using X and Y as
20 surrogates, and to generate a null distribution of P -values. In the case of Ohana, we used the 80 Mb neutral
21 region to estimate the ancestral component proportions and the covariance matrix, and to generate a null
22 distribution of log-likelihood ratios from its selection scan. The maximum number of iterations to estimate the
23 ancestral component proportions and the covariance matrix was set to 20. For Ohana, we considered both the
24 global hypothesis testing whether any ancestry component has an outlying score in the covariance matrix, and
25 a population-specific hypothesis testing whether a specific ancestry component has an outlying score. For the
26 population-specific hypothesis, in scenario (i) we tested the ancestry component most representative of the
27 Native American component in MEX and in scenario (ii) we tested the ancestry component most prevalent in
28 the admixed population. In the case of LAD, which was only used for scenario (ii), we performed local ancestry
29 inference using both RFMix (Maples et al. 2013) and ELAI (Guan 2014). We ran RFMix with the phase
30 correction feature enabled and performed two rounds of the EM algorithm to improve local ancestry calls. In
31 the case of ELAI, we performed 20 rounds of EM iterations. To obtain local ancestry assignment probabilities,
32 we conducted 10 independent runs and averaged probabilities across all runs, as recommended in the ELAI
33 manual. All other parameters for both methods were set to the default except for the time of admixture, which
34 was set to the true generation time. We performed local ancestry inference on the 80 Mb neutral region to
35 generate a null distribution of Z-scores.

36 To estimate and compare the power between the different approaches we simulated a total of 500
37 independent regions under each scenario and for each selection coefficient tested. Each independent
38 simulation consisted of a 2Mb region with a mutation rate of 10^{-8} and a recombination rate of 10^{-8} base pairs
39 per generation. We simulate a single selected site under an additive model near the center of the simulated
40 region and consider 5 different selection coefficients ($s=0.01$ to 0.05 in steps of 0.01 , with s defined here as the
41 increased fitness when carrying one copy of the advantageous allele). The power of each method was based on
42 a P -value cutoff that resulted in a false-positive rate of 0.05 of the respective null distribution.

43 Finally, we also compared the execution time of AdaptMix and Ohana (supplementary table S8). We find that
44 Ohana was much faster when running on a single node, for example taking 80 seconds to run on 150

1 individuals at >200,000 SNPs using 5 iterations, compared to running ADMIXTURE and AdaptMix taking ~5,700
2 seconds in the same population.

3 **Estimation of allele frequencies in ancient Native Americans**

4 To estimate allele frequencies in ancient Native Americans, we queried the Allen Ancient DNA Resource (AADR)
5 available at: <https://reich.hms.harvard.edu/allen-ancient-dna-resource-aadr-downloadable-genotypes-present-day-and-ancient-dna-data>. We downloaded the '1240K' dataset v50.0, which contains ancient and
6 present-day individuals (from either shotgun sequencing data or in-solution target capture, with a range of
7 coverages) at 1,233,013 sites. In order to obtain data for ancient Native Americans without non-Native
8 American ancestry, we kept only individuals with a reported date of more than 500 years BP from countries in
9 the Americas and the Caribbean that passed the quality control filters as defined in the database. After that we
10 selected populations with a minimum of 10 non-missing allelic counts when estimating allele frequencies.

12 **Local ancestry analysis in the CANDELA cohort**

13 Local ancestry assignment was conducted using the HMM approach implemented in ELAI (Guan 2014). The
14 phased genotype data needed as input was obtained by using SHAPEIT2 (Delaneau et al. 2012) with default
15 parameter settings. Genetic distances were obtained from the HapMap Phase II genetic map build GRCh37
16 (Gibbs et al. 2003). As reference continental panels, we used the same Native American, European, and African
17 individuals as in our AdaptMix analysis. ELAI was run setting the admixture generation parameter to 20, and
18 with 20 rounds of EM iterations. To obtain local ancestry assignment probabilities, we conducted 10
19 independent runs and averaged probabilities across all runs as recommended in the ELAI manual. To test for
20 LAD we estimated Z-scores for each ancestry across each locus, and obtained the corresponding one-sided *P*-
21 values testing for a positive deviation.

22 **Population Branch Statistic (PBS) analysis in the CANDELA cohort**

23 We first selected Latin American individuals carrying a specific Native American ancestry component based on
24 the inferred Native American ancestry proportions previously estimated by Chacon-Duque et al 2018 in the
25 CANDELA sample. Specifically, for each Native American ancestry component, we selected CANDELA individuals
26 with >10% inferred ancestry from that particular Native American ancestry component, and with <1%
27 combined inferred ancestry combined across all other Native American components. Thus, each group of
28 admixed Latin Americans was composed primarily of Native American ancestry from a particular Native
29 American component, plus European and African ancestry. We then estimated allele frequencies for each
30 Native American component by considering only alleles (i.e. haplotypes) that were considered of Native
31 American origin with local-ancestry posterior probability >0.9. We only computed allele frequencies for a
32 Native American component if all SNPs genome-wide had >100 alleles (haplotypes) assigned to Native
33 American origin. This resulted in allele frequency estimates for six Native American components, including
34 'Quechua', 'Andes Piedmont', 'Chibcha Paez', 'Nahua1', 'South Mexico', and 'Mapuche' ancestral components
35 (see Chacon-Duque et al. (2018) for a detail description of the inferred components). Pairwise F_{ST} were then
36 estimated using Hudson's estimator as in equation 9 of Bhatia et al. (2013). The branch length (T) between two
37 populations was computed as $T = -\log_{10}(1 - F_{ST})$ (Cavalli-Sforza 1969). The Population Branch Statistic (PBS)
38 (Yi et al. 2010) combines the pairwise branch lengths between three populations, which was computed as:

$$39 \text{ PBS}_{Target} = \frac{T^{Target,Control} + T^{Target,Outgroup} + T^{Control,Outgroup}}{2}.$$

40 PBS values were computed for each Native American component, using all possible pairwise combinations of
41 the other Native components as the control and outgroup populations. The rationale of this analysis was to try
42 to find signals of selection exclusive to a given Native American group (i.e. that likely occurred after the

1 divergence between Native American lineages). For some of our analysis we also used the CHB population from
2 the 1000 Genomes Project, the European reference population, or the African reference population, as control
3 and outgroup populations.

4 **Summary statistics for GWAS and eQTL data**

5 To assess the biological consequence of selected variants, we queried summary statistics from GWASs of
6 relevant phenotypes, and gene-expression data (i.e expression quantitative locus [eQTL] studies) from relevant
7 cell or tissues. For our GWAS query, we retrieved data from immune and metabolic-related phenotypes, as
8 these traits are known to have been subjected to strong selective pressures across several human groups (Fan
9 et al. 2016). Immune-related phenotypes included (i) total white cell count, neutrophil count, lymphocyte
10 count, monocyte count, basophil count, and eosinophil count from the Chen et al. (2020) GWAS study
11 conducted across five continental ancestry groups. Metabolic-related phenotypes included body mass index
12 (BMI), body fat percentage, type II diabetes status, hip circumference, waist circumference, HDL levels, LDL
13 levels, cholesterol levels, and triglycerides levels (Loh et al. 2018). Summary statistics from these GWAS
14 analyses were based on the UK BioBank cohort available at: <http://www.nealelab.is/uk-biobank>. For our eQTL
15 query, we retrieved cis-associations summary statistics of 15 human immune cell types from the DICE
16 (Database of Immune Cell Expression, Expression quantitative trait loci [eQTLs] and Epigenomics) project
17 (Schmiedel et al. 2018), available at: <https://dice-database.org/downloads>. We also retrieved cis-association
18 summary statistics from adipose (subcutaneous, and visceral omentum), muscle (skeletal), and liver tissue from
19 the GTEx Project v7 (Lonsdale et al. 2013) available at: <https://gtexportal.org/home/datasets>.

20 **Acknowledgements**

21 We thank the volunteers for their enthusiastic support for this research. We also thank Alvaro Alvarado,
22 Mónica Ballesteros Romero, Ricardo Cebrecos, Miguel Ángel Contreras Sieck, Francisco de Ávila Becerril, Joyce
23 De la Piedra, María Teresa Del Solar, Paola Everardo Martínez, William Flores, Martha Granados Riveros,
24 Rosilene Paim, Ricardo Gunski, Sergeant João Felisberto Menezes Cavalheiro, Major Eugênio Correa de Souza
25 Junior, Wendy Hart, Ilich Jafet Moreno, Paola León-Mimila, Francisco Quispealaya, Diana Rogel Diaz, Ruth
26 Rojas, and Vanessa Sarabia, for assistance with volunteer recruitment, sample processing and data entry. We
27 also thank Francois Balloux, Aida Andres, Mark McCarthy, Etienne Patin, and Sebastian Cuadros Espinoza for
28 helpful discussion and critical comments on earlier versions of the manuscript. We are very grateful to the
29 institutions that allowed the use of their facilities for the assessment of volunteers, including: Escuela Nacional
30 de Antropología e Historia and Universidad Nacional Autónoma de México (México); Universidade Federal do
31 Rio Grande do Sul (Brazil); 13ª Companhia de Comunicações Mecanizada do Exército Brasileiro (Brazil);
32 Pontificia Universidad Católica del Perú, Universidad de Lima and Universidad Nacional Mayor de San Marcos
33 (Perú). Work leading to this publication received funded from: Leverhulme Trust (RPG-2018-208 to MF), the
34 National Natural Science Foundation of China (#31771393 to ARL), the Scientific and Technology Committee of
35 Shanghai Municipality (18490750300 to ARL), Ministry of Science and Technology of China (2020YFE0201600 to
36 ARL), Shanghai Municipal Science and Technology Major Project (2017SHZDZX01 to ARL) and the 111 Project
37 (B13016 to ARL), the Leverhulme Trust (F/07 134/DF to ARL), BBSRC (BB/I021213/1 to ARL), the Excellence
38 Initiative of Aix-Marseille University - A*MIDEX (a French “Investissements d’Avenir” programme to ARL),
39 Wellcome Trust/Royal Society (098386/Z/12/Z to GH), the National Institute for Health Research University
40 College London Hospitals Biomedical Research Centre, BBSRC (BB/R01356X/1), Universidad de Antioquia (CODI
41 sostenibilidad de grupos 2013- 2014 and MASO 2013-2014), Conselho Nacional de Desenvolvimento Científico
42 e Tecnológico, Fundação de Amparo à Pesquisa do Estado do Rio Grande do Sul (Apoio a Núcleos de Excelência

1 Program) and Fundação de Aperfeiçoamento de Pessoal de Nível Superior. JM-R was supported by a doctoral
2 scholarship from CONCYTEC-PERU (224-2014-FONDECYT).

3 **Data availability**

4 This project only analyses data that has been previously reported in other publications. Raw genotype data for
5 reference populations can be accessed as described previously (The 1000 Genomes Project Consortium 2015;
6 Chacon-Duque et al. 2018). Raw genotype data from CANDELA cannot be made available due to restrictions
7 imposed by ethical approval. Summary statistics from the selection analysis will be deposited in a public
8 repository upon publication.

9 **Software availability**

10 Scripts for selection analyses will be uploaded to a software developer public repository upon publication. The
11 current version of AdaptMix presented in this study is available upon request from g.hellenthal@ucl.ac.uk.
12

1 **Main Figure legends**

2 **Fig. 1. Schematic and intuition of the AdaptMix model. (a)** For each CANDELA individual (columns),
 3 ADMIXTURE-inferred proportions of ancestry related to Native American, European, and African reference
 4 individuals. **(b)** Assuming only two admixing sources in this illustration for simplicity, the model assumes
 5 ancestral populations (A^* and B^*) contribute ancestry proportions α_A and α_B , respectively, to an admixed
 6 population (C') that is ancestral to the tested population (C). Assuming neutrality, the expected allele
 7 frequency (p_0) of C' is estimated using these proportions and the allele frequencies surrogate populations A
 8 and B related to A^* and B^* , respectively. The sampled allele frequency (p) of C is compared to p_0 , with large
 9 deviations indicative of selection (shown with an asterisk in the distribution). **(c and d)** The relationship
 10 between p_0 , the expected allele frequency in the admixed population under neutrality or selection, and α_B , the
 11 ancestry proportion contributed from ancestral population B^* . If selection occurred prior to admixture during
 12 the split between populations B^* and its surrogate B (i.e. along the blue branch in **(b)**), this relationship
 13 increases linearly (blue lines), becoming more differentiated from neutrality (grey line) as the admixture from
 14 B^* increases. In contrast, under selection post-admixture (i.e. along the purple branch in **(b)**), the expected
 15 allele frequency (purple lines) can deviate from neutrality even when the admixture from B^* is near 0. The
 16 difference between the post-admixture and pre-admixture lines is more clear when allele frequencies in
 17 populations A and B are similar (top plot). Solid blue and red lines indicate the allele frequencies in the
 18 surrogate populations A and B , which are used to calculate p_0 .

19 **Fig 2. Performance of AdaptMix to detect and classify selection in simulated Latin American populations. (a)**
 20 Power to detect selection post-admixture, selection in Native Americans, or selection in Europeans in
 21 simulated populations mimicking the Latin American cohorts. Power is based on a P -value cutoff that resulted
 22 in a false-positive rate of 5×10^{-5} in neutral simulations. The power estimated for a given selection coefficient is
 23 based on combining simulations using four different modes of selection (additive, dominant, multiplicative,
 24 recessive) occurring over 12 generations for the post-admixture simulations, over 50 generations for the
 25 selection in Native American simulations, and over 25 generations for the selection in European simulations.
 26 Each simulation for a given combination of parameters consisted of 10,000 advantageous SNPs with a starting
 27 allele frequency of the advantageous allele lower than 0.5. **(b)** The proportion of significant SNPs from (a) that
 28 were assigned to the correct simulated scenario of (left-to-right) post-admixture selection or selection in Native
 29 Americans or Europeans (using a likelihood ratio $> 1,000$ to make a call; otherwise 'Unclassified'). Rows give the
 30 true selection coefficient (legend at right), and the heatmap values give the classification rate. Rows with N.A.
 31 shows instances with less than 50 selected SNPs for which the classification rate is poorly estimated.

32 **Fig. 3. Performance of AdaptMix compared to existing methods. (a)** Power of AdaptMix and Ohana to detect
 33 selection occurring prior to admixture only in the Native American source of an admixed population. The grey
 34 line depicts Ohana's power with $K=4$ when testing for selection only in the ancestry component most
 35 representative of the Native American source, with the brown line testing under the general model. **(b)** Power
 36 of AdaptMix, Ohana, and two local ancestry deviation (LAD) approaches (RFMix, ELAI; Maples et al 2013, Guan
 37 2014) to detect selection occurring in an admixed population directly following the admixture event. The
 38 purple line depicts Ohana's power with $K=3$ when testing for selection only in the ancestry component most
 39 representative of the admixed population, with the green line testing under the general model. See Methods
 40 section for a detailed explanation of the simulation parameters employed for each scenario. Power for **(a)** and
 41 **(b)** is based on a P -value cutoff that resulted in a false-positive rate of 0.05 in neutral simulations.

42 **Fig. 4. Genome-wide selection scan in five Latin American cohorts.** Manhattan plot showing the genomic
 43 regions identified as selected via AdaptMix in each Latin American cohort. The dashed horizontal lines indicate

1 the P -values cutoffs corresponding to a false-positive rate of 5×10^{-5} based on neutral simulations. Different
2 shapes represent the most likely selection model. Names of genes associated with significant SNPs are shown.

3 **Fig 5. Regional selection plot at the HLA region in five Latin American cohorts.** The top plot shows the
4 $-\log_{10}(P\text{-values})$ of SNPs from AdaptMix, the middle plot shows Z -score values based on African local ancestry
5 deviations, and the bottom plot shows genes in the region shaded in grey. Genomic coordinates are in Mb
6 (build hg19 as reference) and genes shown include transcripts.

7 **Fig. 6. Genetic loci with signals of selection at immune-related genes. (a), (b) and (c)** Regional selection plot at
8 three candidate regions of selection encompassing two immune-related genes in the Chilean and one immune-
9 related gene in the Peruvian cohort. Each plot is composed of four panels (rows), consisting of $-\log_{10}(P\text{-values})$
10 of SNPs: (row 1) from AdaptMix; (row 2) associated with immune-related cell counts via GWAS (Chen et al
11 2020); (row 3) associated (as expression quantitative trait loci [eQTLs]) with expression of genes *CD101*, *PTPN2*
12 and *MIF* for (a)-(c), respectively (Schmiedel et al. 2018); with (row 4) depicting genes in the region (in Mb, build
13 hg19 as reference). Horizontal dashed lines give significance thresholds of (row 1) $P\text{-value} = 1 \times 10^{-5}$ based on
14 neutral simulations (row 2) $P\text{-value} = 1 \times 10^{-5}$ (blue line) and $P\text{-value} = 5 \times 10^{-8}$ (red line), and (row 3) $P\text{-value} = 1 \times 10^{-4}$.
15 **(d), (e) and (f)** Derived allele frequency (DAF) in admixed Latin Americans (white circles) stratified by proportion of
16 inferred Native American ancestry, for the SNPs highlighted (vertical dashed line) in
17 top row panels. The sizes of the circles are proportional to the number of individuals in that particular bin.
18 Lines give expected DAF under neutrality (grey), post-admixture selection (brown) or selection in the Native
19 source (black). Horizontal dashed red, blue, and green lines depict DAF for surrogates to Native American,
20 European, and African sources, respectively. AdaptMix's conclusions for these SNPs are selection that is (d)
21 post-admixture, (e) unclassified and (f) pre-admixture in the Native source.

22 **Fig. 7. Genetic loci with signals of selection at metabolic-related genes. (a) and (b)** Regional selection plot at
23 two candidate regions of selection encompassing metabolic-related genes in the Mexican and Peruvian
24 cohorts, respectively. Each plot is composed of four panels consisting of $-\log_{10}(P\text{-values})$ of SNPs: (row 1) from
25 AdaptMix; (row 2) from the UK Biobank GWAS; (row 3) associated (as eQTLs) with expression of *BRINP3* and
26 *HKDC1* for (a)-(b), respectively, (GTEx eQTL study); with (row 4) depicting genes in the region (in Mb, build hg19
27 as reference). Horizontal dashed lines give significance thresholds of (row 1) $P\text{-value} = 1 \times 10^{-5}$ based on
28 neutral simulations (row 2) $P\text{-value} = 1 \times 10^{-5}$ (blue line) and $P\text{-value} = 5 \times 10^{-8}$ (red line), and (row 3) $P\text{-value} = 1 \times 10^{-4}$. **(c)**
29 and **(d)** Derived allele frequency (DAF) in admixed Latin Americans (white circles) stratified by proportion of
30 inferred Native American ancestry, for the SNPs highlighted (vertical dashed line) in top row panels, both of
31 which were classified as reflecting selection in the Native American source. The sizes of the circles are
32 proportional to the number of individuals in that particular bin. Lines give expected DAF under neutrality
33 (grey), post-admixture selection (brown) or selection in the Native American source (black). Horizontal dashed
34 red, blue, and green lines depict DAF for surrogates to Native American, European, and African sources,
35 respectively.

36

1 References

- 2 Acuña-Alonzo V, Flores-Dorantes T, Kruit JK, Villarreal-Molina T, Arellano-Campos O, Hünemeier
3 T, Moreno-Estrada A, Ortiz-López MG, Villamil-Ramírez H, León-Mimila P. 2010. A functional
4 ABCA1 gene variant is associated with low HDL-cholesterol levels and shows evidence of positive
5 selection in Native Americans. *Human molecular genetics* 19:2877-2885.
- 6 Akaike H. 1974. A new look at the statistical model identification. *IEEE transactions on automatic
7 control* 19:716-723.
- 8 Alexander DH, Novembre J, Lange K. 2009. Fast model-based estimation of ancestry in unrelated
9 individuals. *Genome Res* 19:1655-1664.
- 10 Amorim CE, Nunes K, Meyer D, Comas D, Bortolini MC, Salzano FM, Hunemeier T. 2017. Genetic
11 signature of natural selection in first Americans. *Proc Natl Acad Sci U S A* 114:2195-2199.
- 12 Avila-Arcos MC, McManus KF, Sandoval K, Rodriguez-Rodriguez JE, Villa-Islas V, Martin AR, Luisi
13 P, Penaloza-Espinosa RI, Eng C, Huntsman S, et al. 2020. Population History and Gene Divergence in
14 Native Mexicans Inferred from 76 Human Exomes. *Mol Biol Evol* 37:994-1006.
- 15 Badillo Rivera KM, Nieves-Colon MA, Mendoza KS, Davalos VV, Lencinas LEE, Chen JW, Zhang
16 ET, Sockell A, Tello PO, Hurtado GM. 2021. Clotting factor genes are associated with preeclampsia in
17 high altitude pregnant women in the Peruvian Andes. medRxiv.
- 18 Balding DJ, Nichols RA. 1995. A method for quantifying differentiation between populations at multi-
19 allelic loci and its implications for investigating identity and paternity. *Genetica* 96:3-12.
- 20 Basu A, Tang H, Zhu X, Gu CC, Hanis C, Boerwinkle E, Risch N. 2008. Genome-wide distribution of
21 ancestry in Mexican Americans. *Hum Genet* 124:207-214.
- 22 Bersaglieri T, Sabeti PC, Patterson N, Vanderploeg T, Schaffner SF, Drake JA, Rhodes M, Reich DE,
23 Hirschhorn JN. 2004. Genetic signatures of strong recent positive selection at the lactase gene. *Am J
24 Hum Genet* 74:1111-1120.
- 25 Bhatia G, Patterson N, Sankararaman S, Price AL. 2013. Estimating and interpreting FST: the impact
26 of rare variants. *Genome research* 23:1514-1521.
- 27 Bhatia G, Tandon A, Patterson N, Aldrich MC, Ambrosone CB, Amos C, Bandera EV, Berndt SI,
28 Bernstein L, Blot WJ. 2014. Genome-wide scan of 29,141 African Americans finds no evidence of
29 directional selection since admixture. *The American Journal of Human Genetics* 95:437-444.
- 30 Borrego F. 2013. The CD300 molecules: an emerging family of regulators of the immune system.
31 *Blood, The Journal of the American Society of Hematology* 121:1951-1960.
- 32 Bouloc A, Bagot M, Delaire S, Bensussan A, Boumsell L. 2000. Triggering CD101 molecule on
33 human cutaneous dendritic cells inhibits T cell proliferation via IL- 10 production. *European journal of
34 immunology* 30:3132-3139.
- 35 Calandra T, Roger T. 2003. Macrophage migration inhibitory factor: a regulator of innate immunity.
36 *Nature Reviews Immunology* 3:791-800.
- 37 Cavalli-Sforza LL editor.; 1969.
- 38 Chacon-Duque JC, Adhikari K, Fuentes-Guajardo M, Mendoza-Revilla J, Acuna-Alonzo V, Barquera
39 R, Quinto-Sanchez M, Gomez-Valdes J, Everardo Martinez P, Villamil-Ramirez H, et al. 2018. Latin
40 Americans show wide-spread Converso ancestry and imprint of local Native ancestry on physical
41 appearance. *Nat Commun* 9:5388.
- 42 Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. 2015. Second-generation PLINK:
43 rising to the challenge of larger and richer datasets. *Gigascience* 4:s13742-13015-10047-13748.
- 44 Chen M-H, Raffield LM, Mousas A, Sakaue S, Huffman JE, Moscati A, Trivedi B, Jiang T, Akbari P,
45 Vuckovic D. 2020. Trans-ethnic and ancestry-specific blood-cell genetics in 746,667 individuals from
46 5 global populations. *Cell* 182:1198-1213. e1114.
- 47 Cheng JY, Stern AJ, Racimo F, Nielsen R. 2021. Detecting selection in multiple populations by
48 modelling ancestral admixture components. *Mol Biol Evol*.

- 1 Consortium TGP. 2015. A global reference for human genetic variation. *Nature* 526:68.
- 2 Cuadros-Espinoza S, Laval G, Quintana-Murci L, Patin E. 2021. The genomic signatures of natural
3 selection in admixed human populations. *bioRxiv*.
- 4 Delaneau O, Marchini J, Zagury J-F. 2012. A linear complexity phasing method for thousands of
5 genomes. *Nature methods* 9:179-181.
- 6 Deng L, Ruiz-Linares A, Xu S, Wang S. 2016. Ancestry variation and footprints of natural selection
7 along the genome in Latin American populations. *Sci Rep* 6:21766.
- 8 Ettinger NA, Duggal P, Braz RF, Nascimento ET, Beaty TH, Jeronimo SM, Pearson RD, Blackwell
9 JM, Moreno L, Wilson ME. 2009. Genetic admixture in Brazilians exposed to infection with
10 *Leishmania chagasi*. *Ann Hum Genet* 73:304-313.
- 11 Fan S, Hansen ME, Lo Y, Tishkoff SA. 2016. Going global by adapting local: A review of recent
12 human adaptation. *Science* 354:54-59.
- 13 Fumagalli M, Moltke I, Grarup N, Racimo F, Bjerregaard P, Jørgensen ME, Korneliussen TS, Gerbault
14 P, Skotte L, Linneberg A. 2015. Greenlandic Inuit show genetic signatures of diet and climate
15 adaptation. *Science* 349:1343-1347.
- 16 Galinsky KJ, Bhatia G, Loh PR, Georgiev S, Mukherjee S, Patterson NJ, Price AL. 2016. Fast
17 Principal-Component Analysis Reveals Convergent Evolution of ADH1B in Europe and East Asia. *Am
18 J Hum Genet* 98:456-472.
- 19 Ghoussaini M, Mountjoy E, Carmona M, Peat G, Schmidt EM, Hercules A, Fumis L, Miranda A,
20 Carvalho-Silva D, Buniello A. 2021. Open Targets Genetics: systematic identification of trait-
21 associated genes using large-scale genetics and functional genomics. *Nucleic acids research* 49:D1311-
22 D1320.
- 23 Gibbs RA, Belmont JW, Hardenbol P, Willis TD, Yu F, Yang H, Ch'ang L-Y, Huang W, Liu B, Shen
24 Y. 2003. The international HapMap project.
- 25 Giri A, Hellwege JN, Keaton JM, Park J, Qiu C, Warren HR, Torstenson ES, Kovesdy CP, Sun YV,
26 Wilson OD. 2019. Trans-ethnic association study of blood pressure determinants in over 750,000
27 individuals. *Nature genetics* 51:51-62.
- 28 Gravel S, Henn BM, Gutenkunst RN, Indap AR, Marth GT, Clark AG, Yu F, Gibbs RA, Bustamante
29 CD, Project G. 2011. Demographic history and rare allele sharing among human populations.
30 *Proceedings of the National Academy of Sciences* 108:11983-11988.
- 31 Gu S, Li H, Pakstis AJ, Speed WC, Gurwitz D, Kidd JR, Kidd KK. 2018. Recent Selection on a Class I
32 ADH Locus Distinguishes Southwest Asian Populations Including Ashkenazi Jews. *Genes (Basel)* 9.
- 33 Guan Y. 2014. Detecting structure of haplotypes and local ancestry. *Genetics* 196:625-642.
- 34 Guo C, Ludvik AE, Arlotto ME, Hayes MG, Armstrong LL, Scholtens DM, Brown CD, Newgard CB,
35 Becker TC, Layden BT. 2015. Coordinated regulatory variation associated with gestational
36 hyperglycaemia regulates expression of the novel hexokinase HKDC1. *Nature communications* 6:1-8.
- 37 Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. 2009. Inferring the joint
38 demographic history of multiple populations from multidimensional SNP frequency data. *PLoS
39 genetics* 5:e1000695.
- 40 Haak W, Lazaridis I, Patterson N, Rohland N, Mallick S, Llamas B, Brandt G, Nordenfelt S, Harney E,
41 Stewardson K. 2015. Massive migration from the steppe was a source for Indo-European languages in
42 Europe. *Nature* 522:207-211.
- 43 Haller BC, Messer PW. 2019. SLiM 3: forward genetic simulations beyond the Wright–Fisher model.
44 *Molecular biology and evolution* 36:632-637.
- 45 Hamid I, Korunes KL, Beleza S, Goldberg A. 2021. Rapid adaptation to malaria facilitated by
46 admixture in the human population of Cabo Verde. *Elife* 10.
- 47 Hancock AM, Witonsky DB, Ehler E, Alkorta-Aranburu G, Beall C, Gebremedhin A, Sukernik R,
48 Utermann G, Pritchard J, Coop G. 2010. Human adaptations to diet, subsistence, and ecoregion are due
49 to subtle shifts in allele frequency. *Proceedings of the National Academy of Sciences* 107:8924-8930.

1 Harris DN, Ruczinski I, Yanek LR, Becker LC, Becker DM, Guio H, Cui T, Chilton FH, Mathias RA,
2 O'Connor TD. 2019. Evolution of hominin polyunsaturated fatty acid metabolism: from Africa to the
3 New World. *Genome biology and evolution* 11:1417-1430.

4 Hayes MG, Urbanek M, Hivert M-F, Armstrong LL, Morrison J, Guo C, Lowe LP, Scheftner DA,
5 Pluzhnikov A, Levine DM. 2013. Identification of HKDC1 and BACE2 as genes influencing glycemic
6 traits during pregnancy through genome-wide association studies. *Diabetes* 62:3282-3291.

7 Hellenthal G, Busby GBJ, Band G, Wilson JF, Capelli C, Falush D, Myers S. 2014. A genetic atlas of
8 human admixture history. *Science* 343:747-751.

9 Hodgson JA, Pickrell JK, Pearson LN, Quillen EE, Prista A, Rocha J, Soodyall H, Shriver MD, Perry
10 GH. 2014. Natural selection for the Duffy-null allele in the recently admixed people of Madagascar.
11 *Proc Biol Sci* 281:20140930.

12 Hoffmann TJ, Ehret GB, Nandakumar P, Ranatunga D, Schaefer C, Kwok P-Y, Iribarren C,
13 Chakravarti A, Risch N. 2017. Genome-wide association analyses using electronic health records
14 identify new loci influencing blood pressure variation. *Nature genetics* 49:54-64.

15 Homburger JR, Moreno-Estrada A, Gignoux CR, Nelson D, Sanchez E, Ortiz-Tello P, Pons-Estel BA,
16 Acevedo-Vasquez E, Miranda P, Langefeld CD, et al. 2015. Genomic Insights into the Ancestry and
17 Demographic History of South America. *PLoS Genet* 11:e1005602.

18 Joffe GM, Esterlitz JR, Levine RJ, Clemens JD, Ewell MG, Sibai BM, Catalano PM. 1998. The
19 relationship between abnormal glucose tolerance and hypertensive disorders of pregnancy in healthy
20 nulliparous women. *American journal of obstetrics and gynecology* 179:1032-1037.

21 Kanthimathi S, Liju S, Laasya D, Anjana RM, Mohan V, Radha V. 2016. Hexokinase domain
22 containing 1 (HKDC1) gene variants and their association with gestational diabetes mellitus in a south
23 indian population. *Annals of human genetics* 80:241-245.

24 Karlsson EK, Kwiatkowski DP, Sabeti PC. 2014. Natural selection and infectious disease in human
25 populations. *Nature Reviews Genetics* 15:379-393.

26 Kominsky DJ, Campbell EL, Colgan SP. 2010. Metabolic shifts in immunity and inflammation. *The*
27 *Journal of Immunology* 184:4062-4068.

28 Koscielny G, An P, Carvalho-Silva D, Cham JA, Fumis L, Gasparyan R, Hasan S, Karamanis N,
29 Maguire M, Papa E. 2017. Open Targets: a platform for therapeutic target identification and validation.
30 *Nucleic acids research* 45:D985-D994.

31 Lazaridis I, Patterson N, Mittnik A, Renaud G, Mallick S, Kirsanow K, Sudmant PH, Schraiber JG,
32 Castellano S, Lipson M, et al. 2014. Ancient human genomes suggest three ancestral populations for
33 present-day Europeans. *Nature* 513:409-413.

34 Lindo J, Haas R, Hofman C, Apatá M, Moraga M, Verdugo RA, Watson JT, Llave CV, Witonsky D,
35 Beall C. 2018. The genetic prehistory of the Andean highlands 7000 years BP though European
36 contact. *Science advances* 4:eaau4921.

37 Loh PR, Kichaev G, Gazal S, Schoech AP, Price AL. 2018. Mixed-model association for biobank-scale
38 datasets. *Nat Genet* 50:906-908.

39 Long JC. 1991. The genetic structure of admixed populations. *Genetics* 127:417-428.

40 Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, Hasz R, Walters G, Garcia F, Young N.
41 2013. The genotype-tissue expression (GTEx) project. *Nature genetics* 45:580-585.

42 Ludvik AE, Pusec CM, Priyadarshini M, Angueira AR, Guo C, Lo A, Hershenhouse KS, Yang G-Y,
43 Ding X, Reddy TE. 2016. HKDC1 is a novel hexokinase involved in whole-body glucose use.
44 *Endocrinology* 157:3452-3461.

45 Luisi P, García A, Berros JM, Motti JM, Demarchi DA, Alfaro E, Aquilano E, Argüelles C, Avena S,
46 Bailliet G. 2020. Fine-scale genomic analyses of admixed individuals reveal unrecognized genetic
47 ancestry components in Argentina. *PloS one* 15:e0233808.

48 Lumeng CN, Saltiel AR. 2011. Inflammatory links between obesity and metabolic disease. *The Journal*
49 *of clinical investigation* 121:2111-2117.

1 Maples BK, Gravel S, Kenny EE, Bustamante CD. 2013. RFMix: a discriminative modeling approach
2 for rapid and robust local-ancestry inference. *The American Journal of Human Genetics* 93:278-288.
3 Mathieson I. 2020. Limited evidence for selection at the FADS locus in Native American populations.
4 *Molecular biology and evolution* 37:2029-2033.
5 Mathieson I, Lazaridis I, Rohland N, Mallick S, Patterson N, Roodenberg SA, Harney E, Stewardson
6 K, Fernandes D, Novak M, et al. 2015. Genome-wide patterns of selection in 230 ancient Eurasians.
7 *Nature* 528:499-503.
8 Moreno-Estrada A, Gignoux CR, Fernandez-Lopez JC, Zakharia F, Sikora M, Contreras AV, Acuna-
9 Alonzo V, Sandoval K, Eng C, Romero-Hidalgo S, et al. 2014. Human genetics. The genetics of
10 Mexico recapitulates Native American substructure and affects biomedical traits. *Science* 344:1280-
11 1285.
12 Moreno-Estrada A, Gravel S, Zakharia F, McCauley JL, Byrnes JK, Gignoux CR, Ortiz-Tello PA,
13 Martinez RJ, Hedges DJ, Morris RW, et al. 2013. Reconstructing the population genetic history of the
14 Caribbean. *PLoS Genet* 9:e1003925.
15 Norris ET, Rishishwar L, Chande AT, Conley AB, Ye K, Valderrama-Aguirre A, Jordan IK. 2020.
16 Admixture-enabled selection for rapid adaptive evolution in the Americas. *Genome Biol* 21:29.
17 Ochoa D, Hercules A, Carmona M, Suveges D, Gonzalez-Uriarte A, Malangone C, Miranda A, Fumis
18 L, Carvalho-Silva D, Spitzer M. 2021. Open Targets Platform: supporting systematic drug-target
19 identification and prioritisation. *Nucleic acids research* 49:D1302-D1310.
20 Osuna-Ramos JF, Reyes-Ruiz JM, Del Ángel RM. 2018. The role of host cholesterol during flavivirus
21 infection. *Frontiers in cellular and infection microbiology* 8:388.
22 Pasaniuc B, Sankararaman S, Torgerson DG, Gignoux C, Zaitlen N, Eng C, Rodriguez-Cintron W,
23 Chapela R, Ford JG, Avila PC. 2013. Analysis of Latino populations from GALA and MEC studies
24 reveals genomic loci with biased local ancestry estimation. *Bioinformatics* 29:1407-1415.
25 Patterson N, Moorjani P, Luo Y, Mallick S, Rohland N, Zhan Y, Genschoreck T, Webster T, Reich D.
26 2012. Ancient admixture in human history. *Genetics* 192:1065-1093.
27 Pickup J, Crook M. 1998. Is type II diabetes mellitus a disease of the innate immune system?
28 *Diabetologia* 41:1241-1248.
29 Pierron D, Heiske M, Razafindrazaka H, Pereda-Loth V, Sanchez J, Alva O, Arachiche A, Boland A,
30 Olaso R, Deleuze JF, et al. 2018. Strong selection during the last millennium for African ancestry in the
31 admixed population of Madagascar. *Nat Commun* 9:932.
32 Poulter M, Hollox E, Harvey CB, Mulcare C, Peuhkuri K, Kajander K, Sarner M, Korpela R, Swallow
33 DM. 2003. The causal element for the lactase persistence/non-persistence polymorphism is located in a
34 1 Mb region of linkage disequilibrium in Europeans. *Ann Hum Genet* 67:298-311.
35 Pulit SL, Stoneman C, Morris AP, Wood AR, Glastonbury CA, Tyrrell J, Yengo L, Ferreira T, Marouli
36 E, Ji Y. 2019. Meta-analysis of genome-wide association studies for body fat distribution in 694 649
37 individuals of European ancestry. *Human molecular genetics* 28:166-174.
38 Racimo F, Gokhman D, Fumagalli M, Ko A, Hansen T, Moltke I, Albrechtsen A, Carmel L, Huerta-
39 Sánchez E, Nielsen R. 2017. Archaic adaptive introgression in TBX15/WARS2. *Molecular biology and
40 evolution* 34:509-524.
41 Racimo F, Marnetto D, Huerta-Sánchez E. 2017. Signatures of archaic adaptive introgression in
42 present-day human populations. *Molecular biology and evolution* 34:296-317.
43 Refoyo-Martínez A, da Fonseca RR, Halldórsdóttir K, Árnason E, Mailund T, Racimo F. 2019.
44 Identifying loci under positive selection in complex population histories. *Genome research* 29:1506-
45 1520.
46 Reynolds AW, Mata-Miguez J, Miro-Herrans A, Briggs-Cloud M, Sylestine A, Barajas-Olmos F,
47 Garcia-Ortiz H, Rzhetskaya M, Orozco L, Raff JA, et al. 2019. Comparing signals of natural selection
48 between three Indigenous North American populations. *Proc Natl Acad Sci U S A* 116:9312-9317.

- 1 Rishishwar L, Conley AB, Wigington CH, Wang L, Valderrama-Aguirre A, Jordan IK. 2015. Ancestry,
2 admixture and fitness in Colombian genomes. *Sci Rep* 5:12376.
- 3 Robbins GR, Wen H, Ting JP-Y. 2014. Inflammasomes and metabolic disorders: old genes in modern
4 diseases. *Molecular cell* 54:297-308.
- 5 Ruiz-Linares A, Adhikari K, Acuna-Alonzo V, Quinto-Sanchez M, Jaramillo C, Arias W, Fuentes M,
6 Pizarro M, Everardo P, de Avila F, et al. 2014. Admixture in Latin America: geographic structure,
7 phenotypic diversity and self-perception of ancestry based on 7,342 individuals. *PLoS Genet*
8 10:e1004572.
- 9 Rumold CU, Aldenderfer MS. 2016. Late Archaic–Early Formative period microbotanical evidence for
10 potato at Jiskairumoko in the Titicaca Basin of southern Peru. *Proceedings of the National Academy of*
11 *Sciences* 113:13672-13677.
- 12 Santoscoy- Ascencio G, Baños- Hernández CJ, Navarro- Zarza JE, Hernández- Bello J, Bucala R,
13 López- Quintero A, Valdés- Alvarado E, Parra- Rojas I, Illades- Aguiar B, Muñoz- Valle JF. 2020.
14 Macrophage migration inhibitory factor promoter polymorphisms are associated with disease activity
15 in rheumatoid arthritis patients from Southern Mexico. *Molecular genetics & genomic medicine*
16 8:e1037.
- 17 Schmiedel BJ, Singh D, Madrigal A, Valdovino-Gonzalez AG, White BM, Zapardiel-Gonzalo J, Ha B,
18 Altay G, Greenbaum JA, McVicker G. 2018. Impact of genetic polymorphisms on human immune cell
19 gene expression. *Cell* 175:1701-1715. e1716.
- 20 Shuai K, Liu B. 2003. Regulation of JAK–STAT signalling in the immune system. *Nature Reviews*
21 *Immunology* 3:900-911.
- 22 Sibai BM. 2003. Diagnosis and management of gestational hypertension and preeclampsia. *Obstetrics*
23 *& Gynecology* 102:181-192.
- 24 Sirugo G, Williams SM, Tishkoff SA. 2019. The Missing Diversity in Human Genetic Studies. *Cell*
25 177:1080.
- 26 Soares LR, Tsavalier L, Rivas A, Engleman EG. 1998. V7 (CD101) ligation inhibits TCR/CD3-induced
27 IL-2 production by blocking Ca²⁺ flux and nuclear factor of activated T cell nuclear translocation. *The*
28 *Journal of Immunology* 161:209-217.
- 29 Tan Y-X, Hu S-M, You Y-P, Yang G-L, Wang W. 2019. Replication of previous genome-wide
30 association studies of HKDC1, BACE2, SLC16A11 and TMEM163 SNPs in a gestational diabetes
31 mellitus case–control sample from Han Chinese population. *Diabetes, metabolic syndrome and obesity:*
32 *targets and therapy* 12:983.
- 33 Tang H, Choudhry S, Mei R, Morgan M, Rodriguez-Cintrón W, Burchard EG, Risch NJ. 2007. Recent
34 genetic selection in the ancestral admixture of Puerto Ricans. *Am J Hum Genet* 81:626-633.
- 35 Van Dijk M, Mulders J, Poutsma A, Könt AA, Lachmeijer AM, Dekker GA, Blankenstein MA,
36 Oudejans CB. 2005. Maternal segregation of the Dutch preeclampsia locus at 10q22 with a new
37 member of the winged helix gene family. *Nature genetics* 37:514-519.
- 38 van Dijk M, Oudejans C. 2011. STOX1: key player in trophoblast dysfunction underlying early onset
39 preeclampsia with growth retardation. *Journal of pregnancy* 2011.
- 40 Vicente M, Priehodova E, Diallo I, Podgorna E, Poloni ES, Cerny V, Schlebusch CM. 2019.
41 Population history and genetic adaptation of the Fulani nomads: inferences from genome-wide data and
42 the lactase persistence trait. *BMC Genomics* 20:915.
- 43 Vicuna L, Klimenkova O, Norambuena T, Martinez FI, Fernandez MI, Shchur V, Eyheramendy S.
44 2020. Postadmixture Selection on Chileans Targets Haplotype Involved in Pigmentation,
45 Thermogenesis and Immune Defense against Pathogens. *Genome Biol Evol* 12:1459-1470.
- 46 Villarreal-Molina MT, Flores-Dorantes MT, Arellano-Campos O, Villalobos-Comparan M, Rodríguez-
47 Cruz M, Miliar-García A, Huertas-Vazquez A, Menjivar M, Romero-Hidalgo S, Wachter NH. 2008.
48 Association of the ATP-binding cassette transporter A1 R230C variant with early-onset type 2 diabetes
49 in a Mexican population. *Diabetes* 57:509-513.

- 1 Wang S, Lewis CM, Jakobsson M, Ramachandran S, Ray N, Bedoya G, Rojas W, Parra MV, Molina
2 JA, Gallo C, et al. 2007. Genetic variation and population structure in native Americans. *PLoS Genet*
3 3:e185.
- 4 Warren HR, Evangelou E, Cabrera CP, Gao H, Ren M, Mifsud B, Ntalla I, Surendran P, Liu C, Cook
5 JP. 2017. Genome-wide association analysis identifies novel blood pressure loci and offers biological
6 insights into cardiovascular risk. *Nature genetics* 49:403-415.
- 7 Warrington NM, Beaumont RN, Horikoshi M, Day FR, Helgeland Ø, Laurin C, Bacelis J, Peng S, Hao
8 K, Feenstra B. 2019. Maternal and fetal genetic effects on birth weight and their relevance to cardio-
9 metabolic risk factors. *Nature genetics* 51:804-814.
- 10 Weissgerber TL, Mudd LM. 2015. Preeclampsia and diabetes. *Current diabetes reports* 15:1-10.
- 11 Yi X, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZXP, Pool JE, Xu X, Jiang H, Vinckenbosch N,
12 Korneliussen TS. 2010. Sequencing of 50 human exomes reveals adaptation to high altitude. *Science*
13 329:75-78.
- 14 Zamudio S. 2007. High-altitude hypoxia and preeclampsia. *Frontiers in bioscience: a journal and*
15 *virtual library* 12:2967.
- 16 Zhou Q, Zhao L, Guan Y. 2016. Strong selection at MHC in Mexicans since admixture. *PLoS genetics*
17 12:e1005847.
- 18 Zhu Z, Guo Y, Shi H, Liu C-L, Panganiban RA, Chung W, O'Connor LJ, Himes BE, Gazal S,
19 Hasegawa K. 2020. Shared genetic and experimental links between obesity-related traits and asthma
20 subtypes in UK Biobank. *Journal of Allergy and Clinical Immunology* 145:537-549.
- 21
22

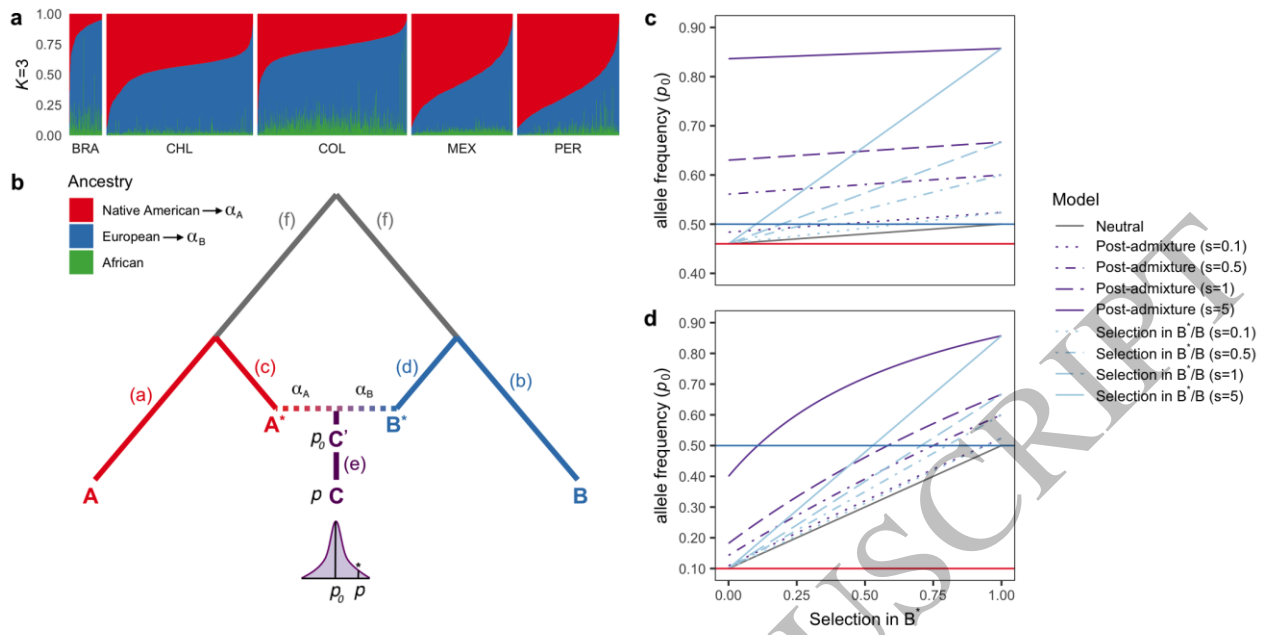


Figure 1
170x88 mm (8.3 x DPI)

1
2
3
4

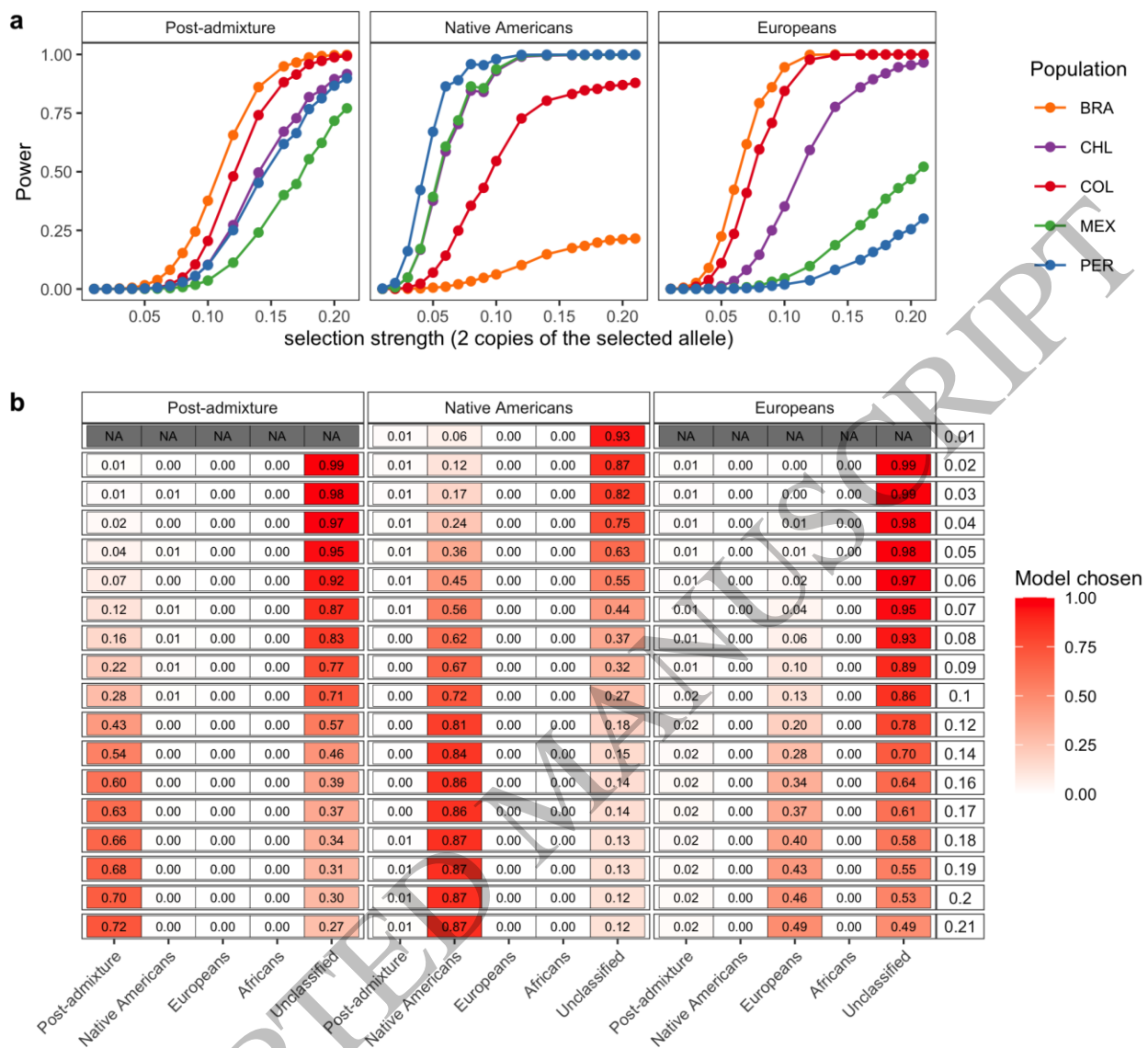


Figure 2
170x158 mm (8.3 x DPI)

1
2
3
4

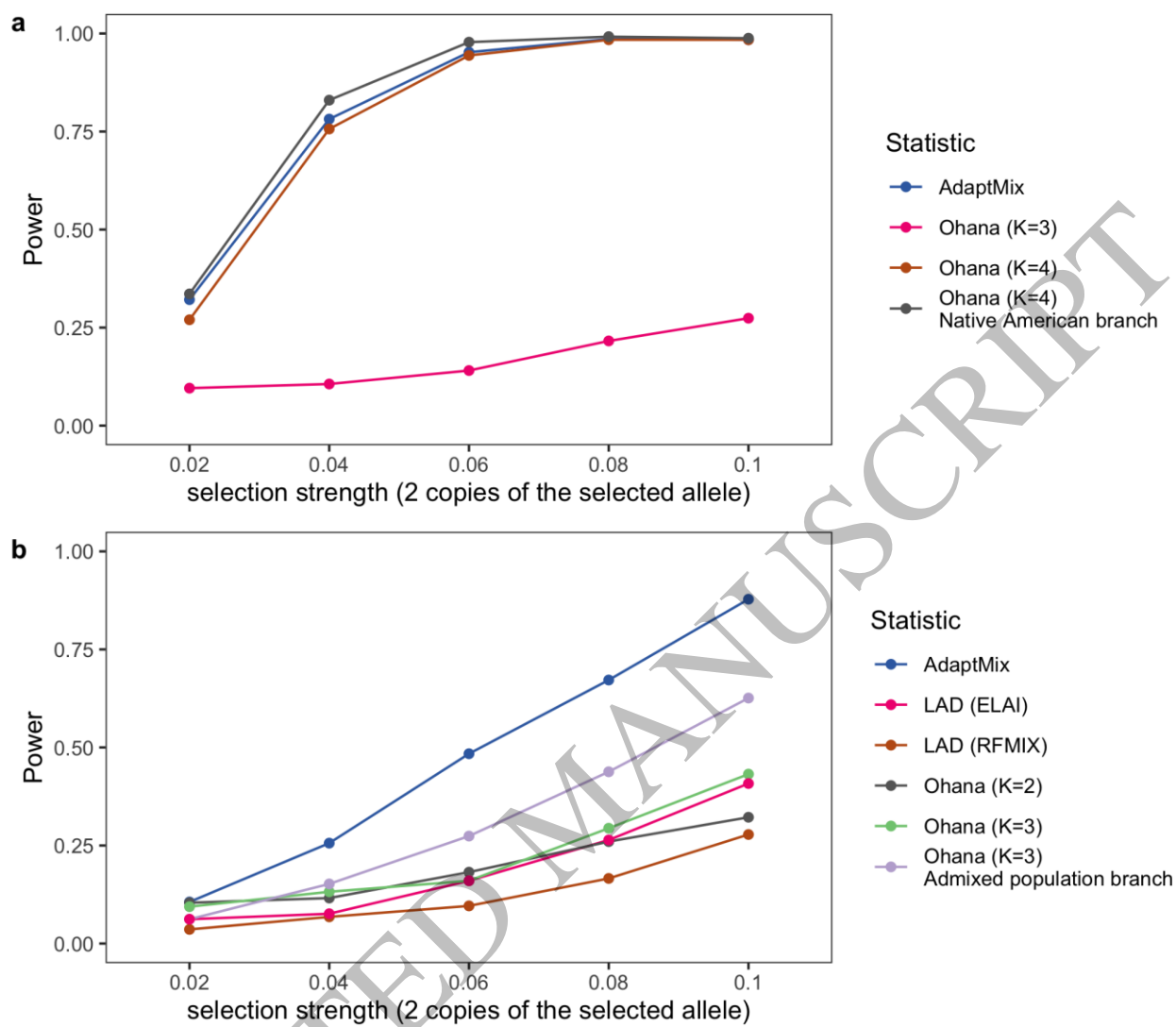


Figure 3
170x149 mm (8.3 x DPI)

1
2
3
4

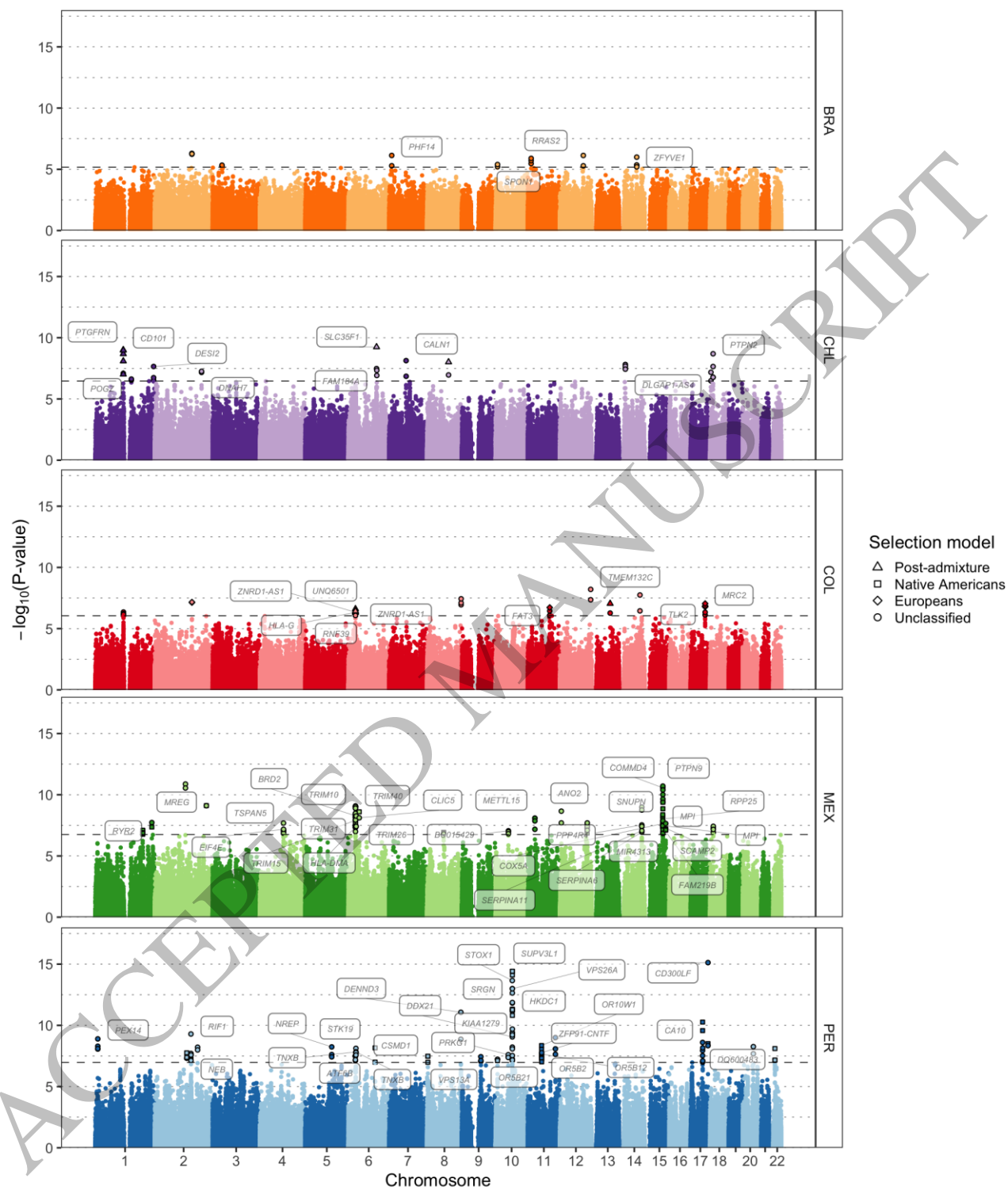
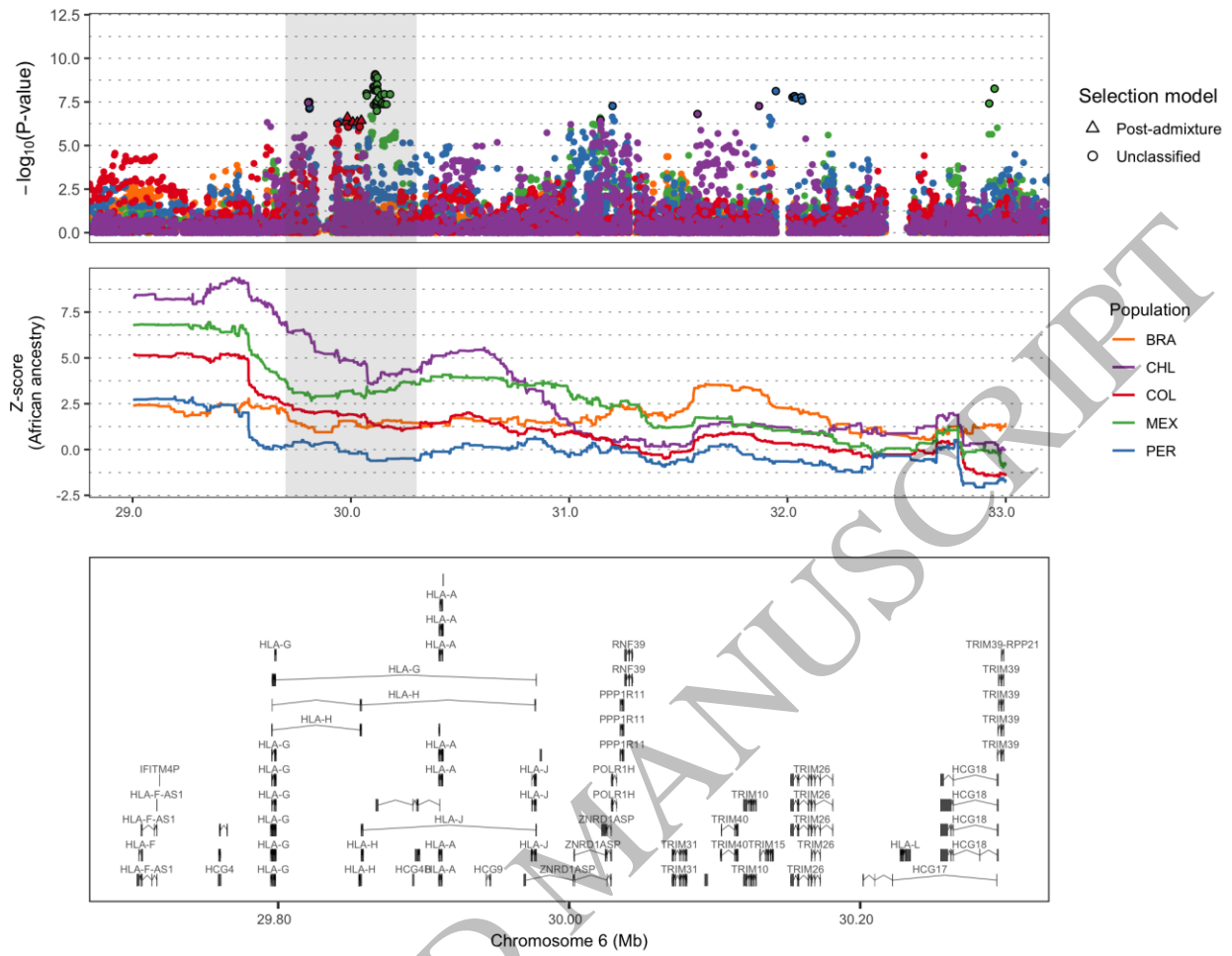


Figure 4
170x204 mm (8.3 x DPI)

1
2
3
4



1
2
3
4

Figure 5
170x130 mm (8.3 x DPI)

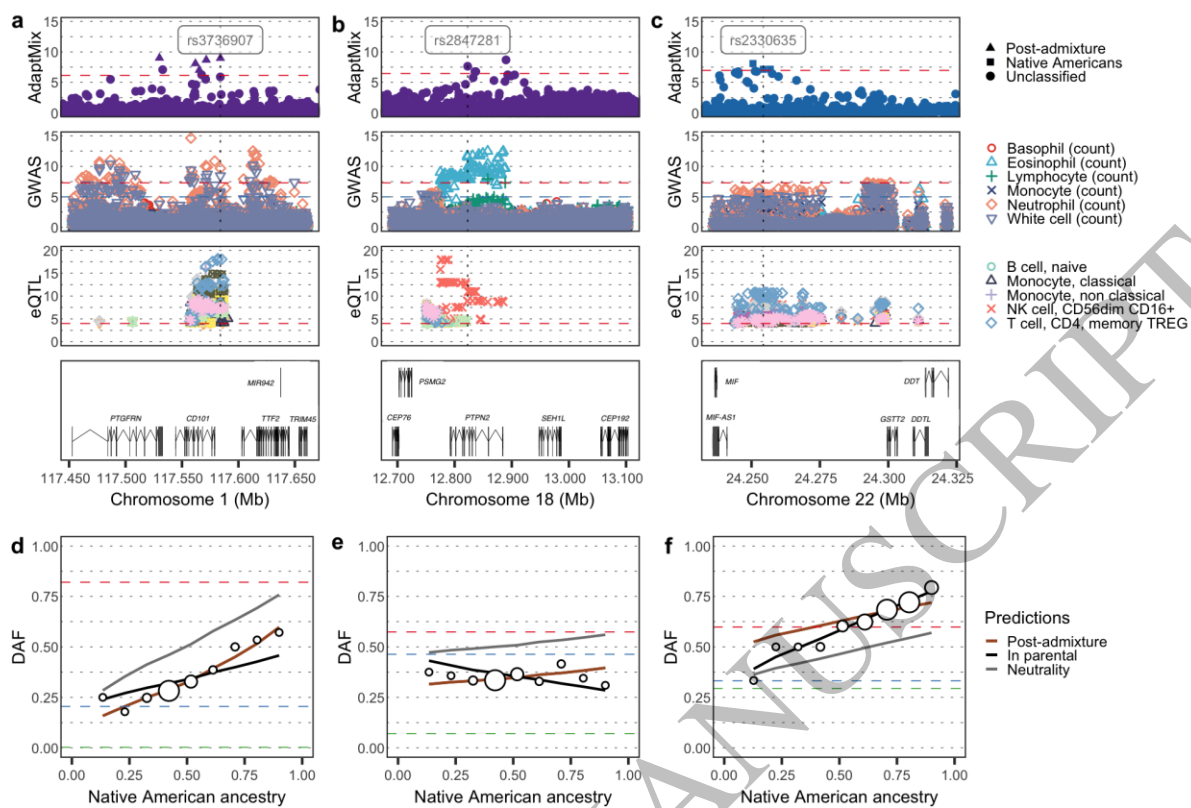


Figure 6
 170x110 mm (8.3 x DPI)

1
 2
 3
 4

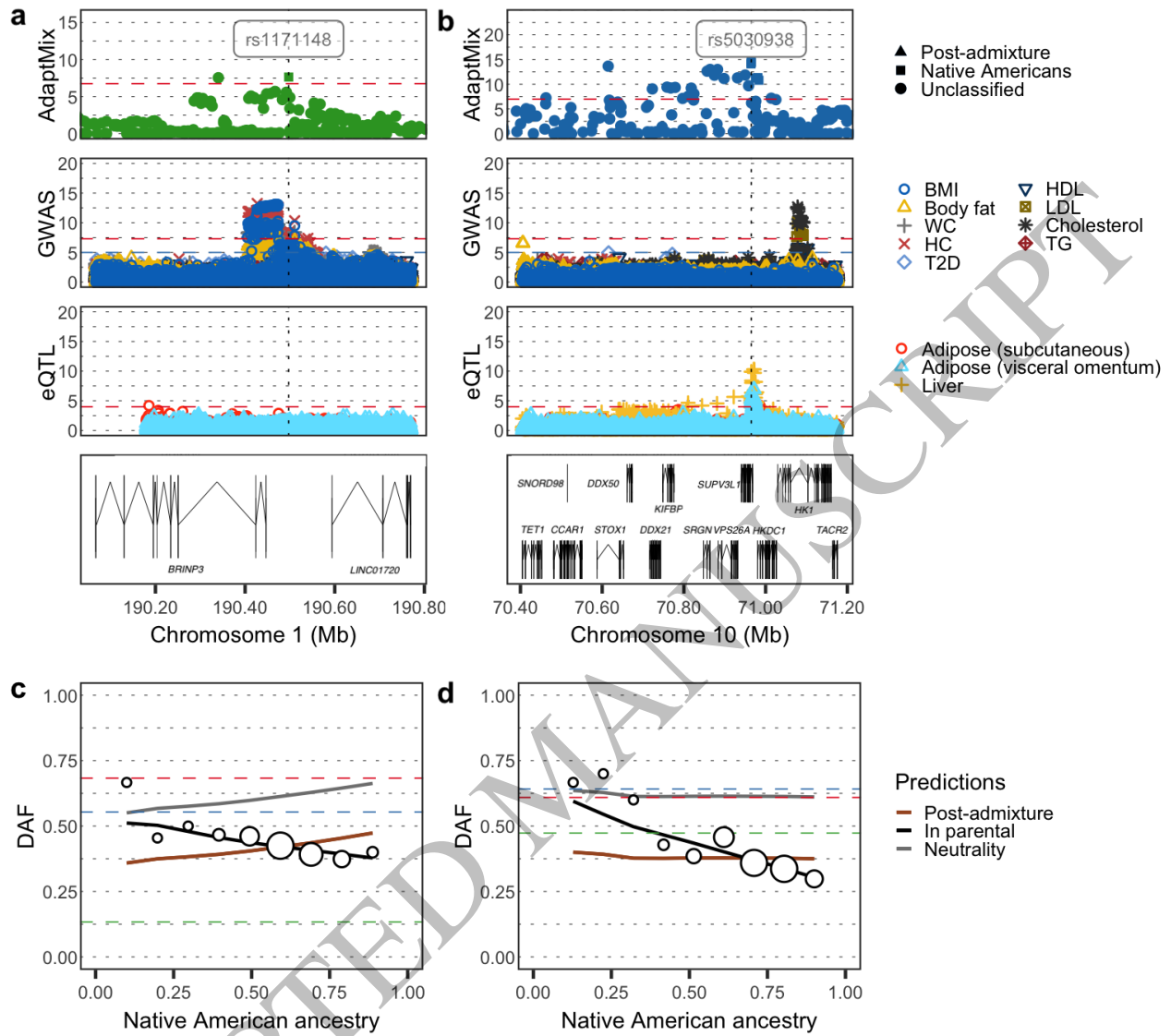


Figure 7
170x157 mm (8.3 x DPI)

1
2
3