

## RESEARCH ARTICLE SUMMARY

## HUMAN GENETICS

## Reconstruction of the human amylase locus reveals ancient duplications seeding modern-day variation

Feyza Yilmaz<sup>†</sup>, Charikleia Karageorgiou<sup>†</sup>, Kwondo Kim<sup>†</sup>, Petar Pajic, Kendra Scheer, Human Genome Structural Variation Consortium, Christine R. Beck, Ann-Marie Torregrassa, Charles Lee<sup>\*</sup>, Omer Gokcumen<sup>\*</sup>

**INTRODUCTION:** Human adaptation to a wide range of diets is a hallmark of our species, sometimes even reflected in our genomic diversity. The amylase gene encodes an enzyme that digests starch, a complex carbohydrate found in many modern human diets. Genomic studies have found substantial variation in the number of amylase gene copies, which is believed to be an adaptive response to dietary changes among human populations, after the advent of agriculture. However, the sequence complexity of the amylase gene region has hindered our understanding of the evolution of its variation and functional implications over time.

**RATIONALE:** Recent technological advances have made it possible to resolve the sequence complexity of the amylase gene region with unprecedented accuracy and detail. Our study has reconstructed this region at nucleotide-level resolution in 98 individuals, using a combination of long-read sequencing and optical genome mapping technologies. We have now elucidated the mechanisms that have given rise to the

genetic and protein variation of amylase and provided insights into the evolutionary trajectory and potential functional effects of this genomic region throughout human history.

**RESULTS:** Our study has identified 30 distinct structural patterns of the amylase gene region across the genomes of 98 modern-day humans. We have found evidence for negative selection at the protein level to maintain the essential function of the amylase genes. Furthermore, we identified two distinct mechanisms, with different mutation rates, that produce the copy number variation and structural patterns seen for the salivary (*AMY1*) and the pancreatic (*AMY2A* and *AMY2B*) genes, respectively. Analysis of archaic hominin genomes showed that some Neanderthals harbored *AMY1* duplications. We also found that hunter-gatherers already had highly variable *AMY1* copy numbers as early as 45,000 years ago, followed by a significant increase in the *AMY1* copy number in the genomes of European farmers over the past 4000 years.

**CONCLUSION:** The molecular archaeology of the amylase region, one of the most structurally dynamic and fastest evolving regions of the human genome, has been comprehensively dissected in this study. Our findings suggest that the initial event leading to multiple *AMY1* copies occurred far before agricultural transitions, possibly even before the human-Neanderthal split. Our results are consistent with an evolutionary scenario where an initial duplication of the *AMY1* gene occurred ~800,000 years ago, leading to the generation of common structural patterns containing three *AMY1* genes. Moreover, our data further demonstrate that copy number variation in the *AMY1* and *AMY2* genes has emerged through two distinct mechanisms. Given the frequency increase in higher *AMY1* copy number patterns in the past 4000 years, selection may have acted on this existing *AMY1* copy number variation, consistent with an adaptive response to the increased role of starch in diets. Taken together, our study provides a robust framework that carefully contextualizes the impact of environmental factors and human lifestyles (such as specific dietary preferences) on the evolution of complex regions of the human genome. ■

The list of author affiliations is available in the full article online.

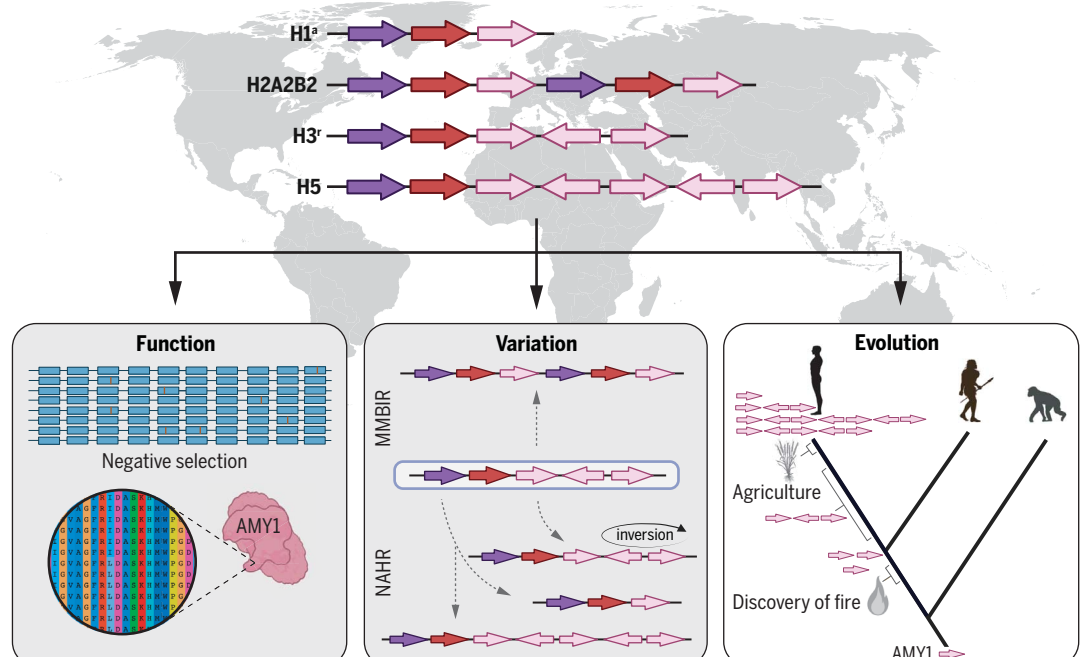
\*Corresponding author. Email: charles.lee@jax.org (C.L.); omergokc@buffalo.edu (O.G.)

<sup>†</sup>These authors contributed equally to this work.

Cite this article as F. Yilmaz et al., *Science* **386**, eadn0609 (2024). DOI: 10.1126/science.adn0609

**READ THE FULL ARTICLE AT**  
<https://doi.org/10.1126/science.adn0609>

**Implications of reconstructing complex genetic variation in the amylase locus.** The comprehensive map of the human amylase locus revealed structural variations such as duplications and inversions (top). Negative selection was observed on all amylase gene-coding sequences (bottom left). Two mechanisms were identified behind these variations: nonallelic homologous recombination (NAHR) and microhomology-mediated break-induced replication (MMBIR) (bottom-middle). Amylase gene duplications predate agriculture and possibly the human-Neanderthal split (bottom right). A putative adaptive increase in variation among European farmers was noted over the past 4000 years. [Figure created with BioRender]



## RESEARCH ARTICLE

## HUMAN GENETICS

# Reconstruction of the human amylase locus reveals ancient duplications seeding modern-day variation

Feyza Yilmaz<sup>1†</sup>, Charikleia Karageorgiou<sup>2†</sup>, Kwondo Kim<sup>1†</sup>, Petar Pajic<sup>2</sup>, Kendra Scheer<sup>2</sup>, Human Genome Structural Variation Consortium<sup>‡</sup>, Christine R. Beck<sup>1,3,4</sup>, Ann-Marie Torregrossa<sup>5,6</sup>, Charles Lee<sup>1\*</sup>, Omer Gokcumen<sup>2\*</sup>

Previous studies suggested that the copy number of the human salivary amylase gene, *AMY1*, correlates with starch-rich diets. However, evolutionary analyses are hampered by the absence of accurate, sequence-resolved haplotype variation maps. We identified 30 structurally distinct haplotypes at nucleotide resolution among 98 present-day humans, revealing that the coding sequences of *AMY1* copies are evolving under negative selection. Genomic analyses of these haplotypes in archaic hominins and ancient human genomes suggest that a common three-copy haplotype, dating as far back as 800,000 years ago, has seeded rapidly evolving rearrangements through recurrent nonallelic homologous recombination. Additionally, haplotypes with more than three *AMY1* copies have significantly increased in frequency among European farmers over the past 4000 years, potentially as an adaptive response to increased starch digestion.

Copy number variation at the amylase locus is frequently attributed to human health and adaptation (1). As such, this structurally variable locus is a prime target for research on the fundamental biology of gene duplications. There are two types of amylase genes, *AMY1* and *AMY2*, which are reported to be expressed in the salivary glands and pancreas, respectively (2). Both genes encode the amylase enzyme, which breaks down polymeric starch into simple sugar molecules, a crucial digestive process for starch-eating species (3). It has been shown that mammals that consume starch-rich diets underwent independent bursts of amylase gene duplications from the ancestral pancreatic *AMY2-like* gene (4). A great ape-specific duplication resulted in the formation of the salivary *AMY1* gene (5), which has since evolved to produce unusual copy number variations, ranging from 2 to 17 copies per diploid cell (1, 6). This variation is especially prominent in human populations with high starch consumption, particularly those with a history of agriculture (1, 6, 7). These evolutionary insights indicate that copy number variation at the amylase locus may play an adaptive role in shaping the metabolic response

to starchy diets, including the presence of microbes that break down amylase-resistant starch (8).

Given its adaptive and putative functional roles, duplications of the *AMY1* gene were linked to the advent of agriculture ~10,000 years ago (1). The lack of nucleotide-level resolution in evolutionary analysis has led to disagreements about the timing and functional importance of *AMY1* gene duplications in relation to starch-rich diets and human evolution (7, 9–12). To address this issue, we have resolved this locus at nucleotide level at a population scale across 98 individuals from different populations, using optical genome mapping and long-read sequencing techniques. The nucleotide-resolved haplotypes of this locus subsequently allowed us to conduct evolutionary genetic analyses on ancient human and archaic hominin genomes to investigate the timing of *AMY1* gene duplications within the context of agriculture.

## Results

### Structural haplotypes at the human amylase locus

The amylase locus in the human genome is a ~212.5-kilobase pair (kbp) region on chromosome 1 (GRCh38; chr1:103,554,220 to 103,766,732) which contains *AMY2B*, *AMY2A*, *AMY1A*, *AMY1B*, and *AMY1C* genes (Fig. 1A). This locus is largely composed of segmental duplications with >99% sequence similarity, which complicates its accurate assembly using short-read sequencing (fig. S1). Using the sequence similarity of segmental duplications and the labeling patterns from the optical genome mapping data from the GRCh38 reference assembly in silico map, we defined six distinct amylase segments overlapping

the amylase genes, depicted by colored arrows (Fig. 1A and table S1). Using optical genome mapping, which has been previously demonstrated to resolve similar complex regions (13–15), we constructed haplotype-resolved diploid assemblies for 98 individuals [ $n = 196$  sampled alleles (i.e., haploid sample size); table S2] and characterized this locus using the copy number and orientation of the amylase segments (Fig. 1A). This approach allowed us to identify 52 distinct amylase haplotypes (fig. S2, A and B, and table S3), of which 7 were previously reported (7) (fig. S3). These haplotypes were then classified on the basis of the number of amylase gene copies, adhering to the established nomenclature HXAYBZ, where HX represents *AMY1*, AY represents *AMY2A*, and BZ represents *AMY2B* copy numbers; superscripts “a” and “r” indicate ancestral and reference haplotypes, respectively (fig. S4) (7, 11). We subsequently defined 30 high-confidence haplotypes (from 117 observed alleles in 81 individuals) that were orthogonally supported by de novo assemblies on the basis of long-read sequencing (Fig. 1B) (16). This represents the first nucleotide-level reconstruction of the amylase locus at a population scale.

The length of the amylase haplotypes ranged from 111 kbp (H1<sup>a</sup>.1 and H1<sup>a</sup>.2) to 402 kbp (H7.1) (Fig. 1B), capturing those that are structurally identical to the GRCh38 (H3<sup>r</sup>.1) and the T2T-chm13 (H7.3) (17) reference assemblies (Fig. 1A). Four haplotypes, H1<sup>a</sup>.1 ( $n = 20$  out of 117), H3<sup>r</sup>.1 ( $n = 18$  out of 117), H3<sup>r</sup>.2 ( $n = 22$  out of 117), and H3<sup>r</sup>.4 ( $n = 21$  out of 117), were categorized as common, each with an allele frequency exceeding 5% across all studied populations. These four common haplotypes collectively constitute ~70% of all amylase haplotypes ( $n = 81$  out of 117) in this study. Despite our limited sample size, we found that the four common haplotypes exist in all continental regions (Fig. 1C and fig. S5). In addition, *AMY1* copy number variations do not exhibit a discernible geographic specificity ( $P$  value = 0.4312, Kruskal-Wallis rank sum test).

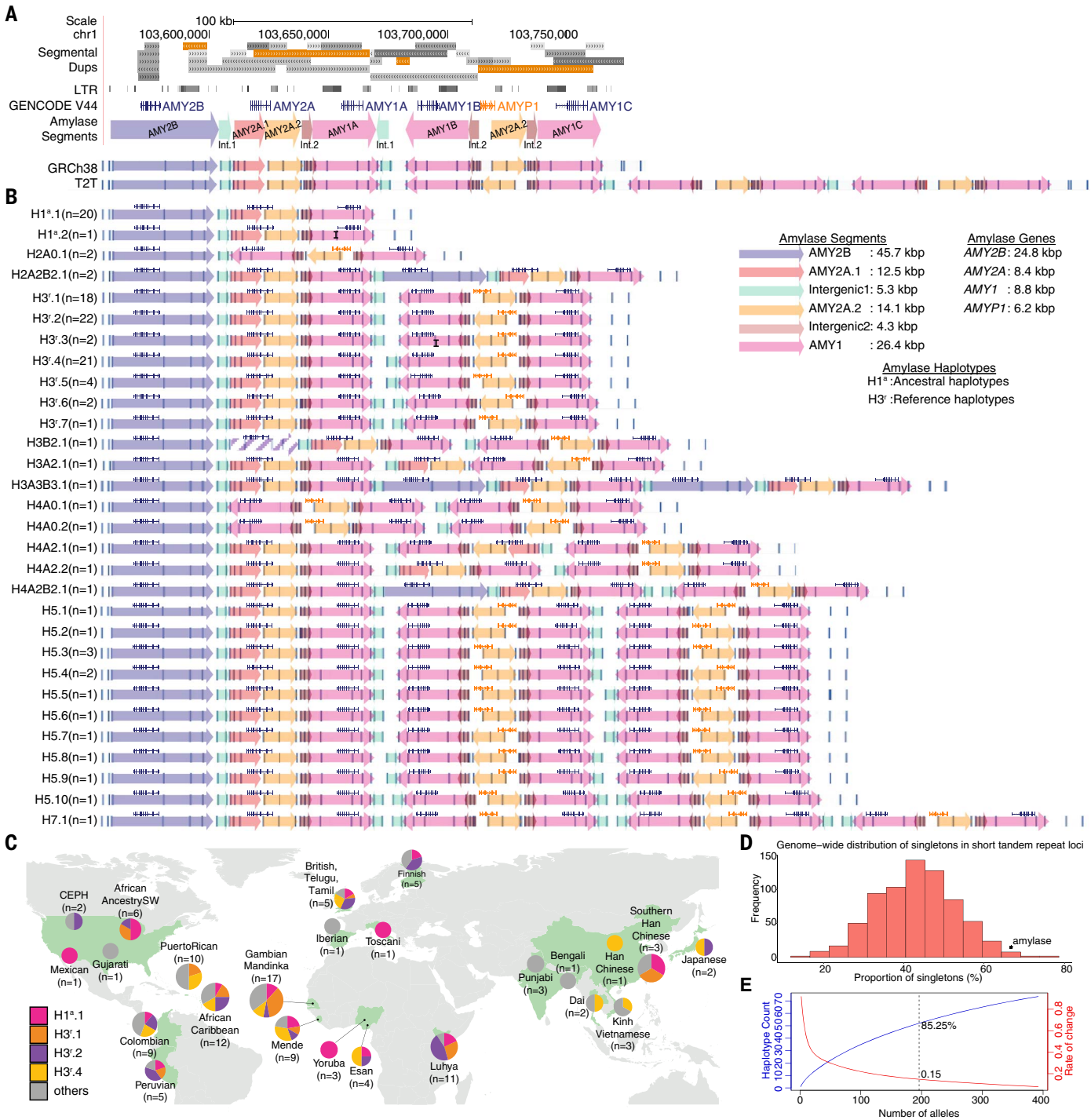
Out of 30 “high confidence” amylase haplotypes, we identified 19 (63%) as singletons, which are haplotypes that appear only once in our dataset. To infer the relative mutation rate of the amylase locus, we compared this number to that of tandem repeats. To avoid potential biases, we used the same 33 individuals that are present in both the tandem repeat database (EnsembleTR) and our dataset. Only 21 of the 30 distinct amylase haplotypes were present in these 33 individuals. To ensure the analysis was consistent across different genomic regions, we only considered matched tandem repeat loci that had exactly 21 detected alleles in the human population. This yielded 719 tandem repeat loci (fig. S6). We found that the proportion of singleton haplotypes was significantly higher for the amylase locus than

<sup>1</sup>The Jackson Laboratory for Genomic Medicine, Farmington, CT, USA. <sup>2</sup>Department of Biological Sciences, University at Buffalo, Buffalo, NY, USA. <sup>3</sup>University of Connecticut, Institute for Systems Genomics, Storrs, CT, USA. <sup>4</sup>The University of Connecticut Health Center, Farmington, CT, USA. <sup>5</sup>Department of Psychology, University at Buffalo, Buffalo, NY, USA. <sup>6</sup>University at Buffalo Center for Ingestive Behavior Research, University at Buffalo, Buffalo, NY, USA.

\*Corresponding author. Email: charles.lee@jax.org (C.L.); omegokc@buffalo.edu (O.G.)

†These authors contributed equally to this work.

‡Human Genome Structural Variation Consortium collaborators and affiliations are listed at the end of this paper.



**Fig. 1. Amylase structural haplotypes identified from present-day humans in this study.** (A) Segmental duplications (light to dark gray: 90 to 98% similarity; light to dark orange: >99% similarity), GENCODE V44 gene annotations, and long terminal repeats (LTRs) are represented as tracks. The lower panel shows amylase segments (colored arrows), and haplotype structures of GRCh38 and T2T-chm13 reference assemblies, represented as in silico maps with white backgrounds and vertical blue lines displaying optical mapping labels. The AMY2B segment overlaps the AMY2B gene, AMY2A.1 and AMY2A.2 segments overlap the AMY2A gene, and the AMY1 segment overlaps the AMY1 gene. (B) The high-confidence amylase structural haplotypes resolved in our dataset ( $n = 30$ ). The vertical black line in the second AMY1 segment of H3<sup>r</sup>.3

represents the polymorphic label present in three alleles. The diagonal stripes in the second AMY2B segment of H3B2.1 indicate that it is a partial copy of the first AMY2B segment. Haplotype IDs: HX: X denotes the number of AMY1 copies; AY: Y denotes the number of AMY2A copies; BZ: Z denotes the number of AMY2B copies. The superscript “a” denotes the ancestral amylase haplotype structure, and the superscript “r” denotes the reference amylase haplotype structure. The number in parentheses indicates the number of alleles. (C) The distribution of common amylase haplotypes across 26 population samples. (D) The proportion of singletons for tandem repeat loci (EnsembleTR) across the genome. For adequate comparison, we used the same individuals ( $n = 33$ ) for whom we were able to reconstruct amylase haplotypes in our dataset.



Additionally, we filtered the tandem repeat loci (719 loci, unit length 1 to 6 bp) that we analyzed to match the number of distinct alleles ( $n = 21$ ) observed in the amylase locus. The asterisk (\*) represents the proportion of singletons among all distinct haplotypes (~67%, 14 out of 21) detected at the amylase locus. (E) Rarefaction and extrapolation sampling curve based on 52 amylase haplotypes,

displaying how the number of distinct haplotypes (blue line) is projected to saturate with the increase in the number of alleles. The rate of change (red line, 0.15) indicates the number of previously unknown haplotypes discovered per unit increase in the number of analyzed alleles. The dashed line shows the proportion of estimated number of samples (85.25%) captured in our study.

**Fig. 2. The variants in amylase coding sequences and negative selection on three amylase gene types.** (A) The maximum likelihood phylogenetic tree of amylase coding sequences (left) and dN/dS estimate for each amylase gene type (right). The phylogenetic tree is rooted with a coding sequence from the sheep genome (Oar\_rambouillet\_v1.0).

The number of nucleotide and resultant amino acid changes that are paralog-specific are indicated. The numbers in parentheses indicate nucleotide and amino acid changes that are variable within the

AMY1 branch. The numbers to the right of each bar represent the FDR-adjusted  $P$  value for the likelihood ratio between  $H_{02}$  (dN/dS ratio is fixed to one on the foreground branches) and  $H_1$  (two dN/dS ratios are allowed on the foreground and background branches, respectively) for each gene type. (B) The positions of amino acid variants within and between AMY2B, AMY2A, and AMY1 protein sequences. The 211 and 366 positions are highlighted because they overlap with a conserved section of the amylase protein sequence and have a predicted functional impact [AlphaMissense (90)]. The coordinates are based on the residues of the amylase enzyme from UniProtKB (accession: PODUB6). Single-letter abbreviations for the amino acid residues are as follows: A, Ala; C, Cys; D, Asp; E, Glu; F, Phe; G, Gly; H, His; I, Ile; K, Lys; L, Leu; M, Met; N, Asn; P, Pro; Q, Gln; R, Arg; S, Ser; T, Thr; V, Val; W, Trp; and Y, Tyr.

the genome-wide average for the 719 tandem repeat loci (the observed empirical percentile = 0.017, Fig. 1D). This observation is informative for understanding the mutation rate at this locus because the allele frequency spectrum and proportion of singletons are determined by mutation rate and genetic drift (18). By using the same individuals and matching the number of distinct haplotypes in our comparison, we control for demographic biases and provide a relative estimation of the mutation rate for the amylase locus. Considering that short tandem repeats have mutation rates as low as  $10^{-8}$  (similar to single-nucleotide variant mutation rate) and, in some cases, can be as high as  $10^{-2}$  mutations per locus per generation (19), our analysis encompasses the entire range of mutation rates found in the genome. We acknowledge that an ideal future comparison would involve other amylase-like loci, which exhibit similar mutational mechanisms and levels of structural variation and that are resolved using comparable approaches once such databases become available. When we

repeated this analysis for the complex 3q29 locus, known for its segmental duplication-rich nature and high levels of structural variation and that has been resolved with similar approaches (15, 20), we found 11 (50%) singleton haplotypes (fig. S7). Thus, the amylase locus is mutating faster than a typical structural variation hotspot and 98.3% of analyzed tandem repeats.

To understand the extent of the amylase haplotypes that are captured in our study compared to what exists in the human population, we conducted rarefaction analysis on the 98 samples. We identified all common haplotypes with a frequency of  $\geq 5\%$  (fig. S8) and 85% of all haplotypes overall (Fig. 1E and fig. S9).

**Strong negative selection limits functional variation among amylase gene copies**

To systematically assess the selection pressure on the amylase genes coding sequences, we examined the degree of protein-coding sequence variation associated with our high-confidence

amylase haplotypes (30 haplotypes from 117 alleles). Gene annotation predicted 582 distinct intact protein-coding amylase gene copies in 117 alleles, and we experimentally validated these predictions using digital droplet polymerase chain reaction (ddPCR) on 18 randomly selected individuals (coefficient of determination  $R^2 = 0.94$ ; fig. S10A and tables S4 and S5) (21). The reconstruction of the coding sequence phylogeny revealed that all human amylase gene copies can be robustly clustered into three distinct types: AMY2B, AMY2A, and AMY1 (bootstrap value = 96%; Fig. 2A and fig. S11). We found that the AMY2B, AMY2A, and AMY1 genes had 23, 23, and 36 fixed coding sequence variations specific to each type, resulting in 6, 11, and 19 gene-specific amino acid differences, respectively (table S6). On the basis of the coding sequence alignment (22), we estimated the synonymous and nonsynonymous substitution (dN/dS) ratios using codeml (23) and found no evidence for lineage-specific selection pressure acting on any of the amylase gene types [false discovery rate (FDR) adjusted  $P$  value from  $\chi^2$  test > 0.05;

table S7). By contrast, all amylase gene types show significant signatures of negative selection (FDR adjusted  $P$  value from  $\chi^2$  test  $<0.05$ ; Fig. 2A and table S7). These observations suggest that negative selection (dN/dS ratio  $<1$ ) has acted to retain the amino acid sequence of amylase gene copies, both within and between the three amylase gene types. It is of note that we identified two amino acid variants at positions 211 and 366 (accession: P0DUB6) that could contribute to functional differences and may have biomedical relevance (Fig. 2B, fig. S12, and table S8) (24).

Our coding sequence data from gene annotations cover 582 intact amylase gene copies and confirmed the explicit difference between *AMY1*, *AMY2A*, and *AMY2B* genes, which is crucial to distinguish the expression of these genes across tissues. Indeed, according to existing data portals, Genotype-Tissue Expression and The Human Protein Atlas (25, 26), *AMY2A* and *AMY2B* are expressed in the pancreas and, to some level, in adipose and brain tissues, whereas the *AMY1* gene is expressed primarily in the parotid salivary gland (fig. S13). These observations are consistent with the idea that *AMY1* first emerged in the ancestor of great apes, resulting in a gain of expression in the parotid gland tissues (5). Furthermore, subsequent *AMY1* gene duplications in the human lineage appear to affect the dosage of amylase in the parotid salivary glands (1). Among the haplotypes that we identified in this study, we found 110 amylase pseudogenes and showed that they share a single phylogenetic origin from an ancestral incomplete gene duplication of the *AMY2A* gene (27) (fig. S14). Thus, the pseudogenization of *AMY2A* is more likely due to a single mutational event rather than a loss of constraint and the repeated occurrence of new loss-of-function variants.

### Evolution of the copy number of *AMY1* gene

We show that all human amylase gene copies can be robustly clustered into three distinct types: *AMY2B*, *AMY2A*, and *AMY1*. However, to specifically study the evolution of the copy number of the salivary *AMY1* gene, we needed to identify the sequences within the amylase locus that were the most phylogenetically informative. To achieve this, we systematically evaluated the variation in 117 alleles by aligning sequences and identified a consensus sequence for each amylase segment (fig. S15, A and B). We identified an interval (from 22,850 to 26,730 bp) (22) within the *AMY1* segment (fig. S15C), where all observed *AMY1* segments ( $n = 337$ ) can be phylogenetically structured into three distinct clusters: *AMY1A* ( $n = 124$ ), *AMY1B* ( $n = 99$ ), and *AMY1C* ( $n = 114$ ) (fig. S15D). These clusters correspond to the segments represented in the GRCh38 reference assembly. Further, the chimpanzee (panTro6) and gorilla (gorGor6) reference genomes each contain

only one *AMY1C*-like segment. Despite the structural similarity of these haplotypes, we found that the nonhuman primate sequences cluster distinctly from those of the human H1<sup>a</sup>.1 haplotypes (figs. S16 and S17). These findings suggest that the common ancestor of humans and chimpanzees possessed a single *AMY1C*-like segment, and the *AMY1A* and *AMY1B* segments have evolved only in the human lineage (fig. S17). Note that the bonobo reference genome (panPan3) contains two *AMY1* segments: one ancestral and one resulting from an independent, bonobo-specific duplication, as determined through synteny and phylogenetic analysis (figs. S16 and S17). One of these duplications is likely nonfunctional owing to a previously reported disrupted coding sequence (1), a finding corroborated by the most recent annotation (RS\_2024\_02/NHGRI\_mPanPan1-v2.0\_pri).

To further understand the evolution of the *AMY1* gene copy number, we aligned all *AMY1* segments from the most common haplotypes, H3<sup>T</sup>.1 and H3<sup>T</sup>.2, that harbor all three *AMY1* segment types. Two independent Bayesian phylogenies based on these alignments indicate that the *AMY1B* segment arose from *AMY1C* ~140 to 270 thousand years ago (ka), followed by the duplication of *AMY1A* from the *AMY1B* segment ~120 to 240 ka (fig. S18 and table S9). Gene conversion between the GC-rich segmental duplications complicates time estimation based on a molecular clock and is a known phenomenon at the amylase locus (6, 28). Considering gene conversion between *AMY1* segments, the actual duplication dates are expected to be older than the estimates above. Some studies have used single-nucleotide variants in the flanking regions to infer the phylogenetic history of this locus, positing coalescence dates for human amylase locus at ~279 and ~450 ka (11, 12). However, as described previously (7), we found that the linkage disequilibrium between the flanking single-nucleotide variants and amylase haplotypes is low (e.g., H1<sup>a</sup>.1, average  $R^2 = \sim 0.26$  and median  $R^2 = \sim 0.03$ ) (fig. S19), complicating the time estimation of *AMY1* duplications using flanking regions. Therefore, our estimates avoid these complications and support the conclusion that the initial *AMY1* gene duplications substantially predated out-of-Africa migrations, by at least 30 thousand years (table S9).

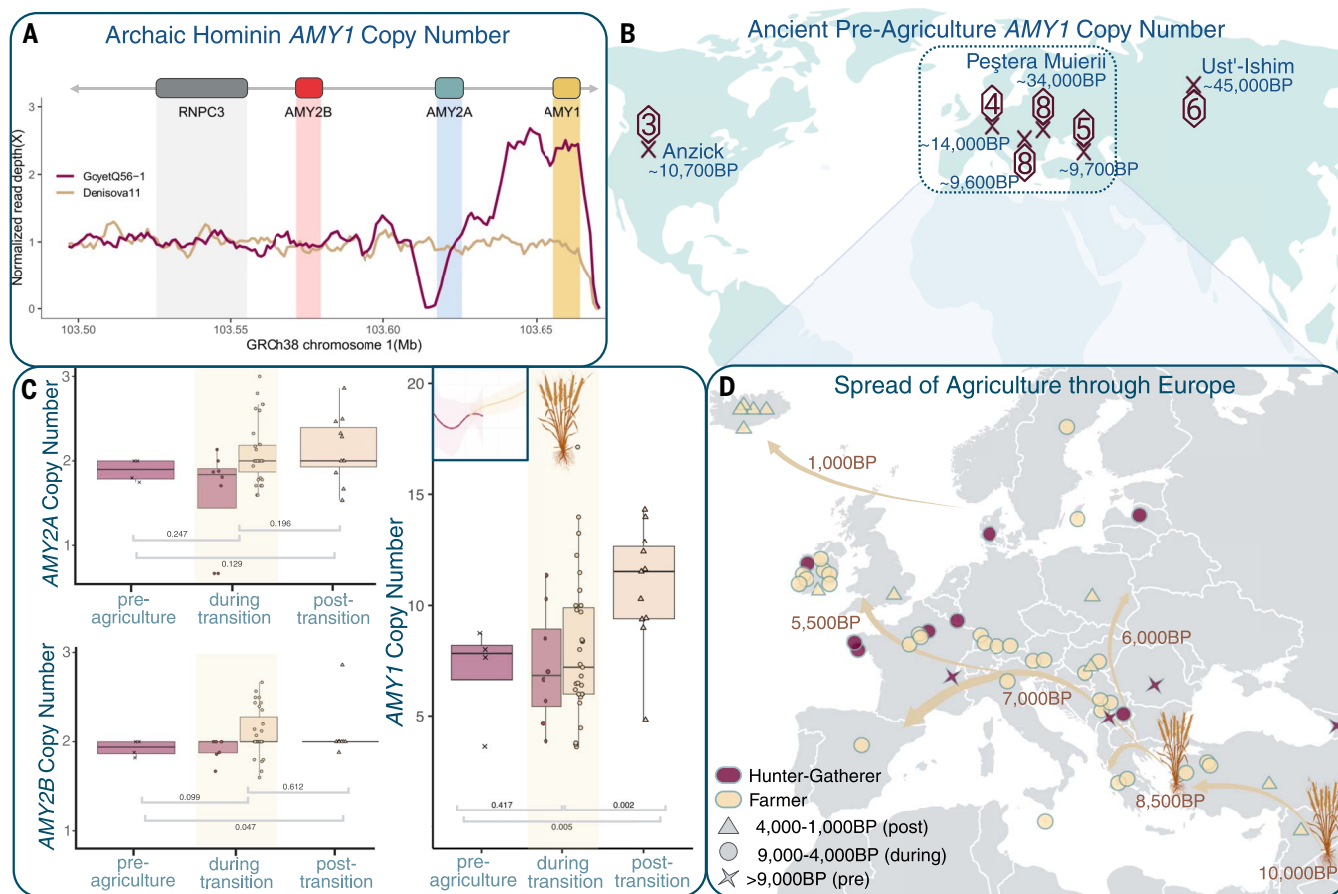
### *AMY1* copy number variation in archaic hominin genomes

A complementary approach for estimating the relative timing of gene duplications involves analyzing the read-depth of unique k-mers in ancient human and archaic hominin genomes. We first tested a k-mer approach using the GeneToCN algorithm (29) on short-read sequencing data to estimate the copy numbers of the *AMY1*, *AMY2A*, and *AMY2B* genes in 116 present-day human genomes (table S5).

Notably, the k-mer approach achieved an  $R^2 > 0.99$  correlation for 32 individuals, each with both haplotypes of the amylase locus reconstructed (fig. S10B). We tested our k-mer approach in 101 samples (a subset of the 116 individuals mentioned above) with digital droplet PCR validation data ( $R^2 = 0.95$ ; fig. S10A). These results suggest that the GeneToCN estimates are consistent with the digital droplet PCR estimations and are a viable option for estimating *AMY1* gene copy numbers in short-read whole-genome sequencing datasets.

We then aimed to estimate *AMY1* gene copy number in eight archaic hominin genomes using two approaches: (i) the validated k-mer method described above, as well as (ii) an independent read-depth analysis (table S10). Given the varying genome-wide coverage in most of these genomes, we needed to conduct a downsampling analysis to empirically determine that 1X and 5X genome-wide coverage provides  $>85\%$  and  $>95\%$  accuracy in estimating *AMY1* copy number, respectively (fig. S20 and table S11) (30). Controlling for GC bias and coverage within the amylase locus for each sample, we could reliably estimate the *AMY1* copy number for eight archaic hominin genomes (table S10). We found increased *AMY1* copy numbers in two Eastern and one Western Neanderthal, as well as in one Denisovan genome (Fig. 3A, fig. S21, and table S10). These include the Altai Neanderthal (2.6 copies), Denisova 2 (8 copies), GoyetQ56-1 (5.0 copies), and Mezmaiskaya 2 (4.7 copies). Previous read-depth analysis of Altai Neanderthal and Denisovan genomes found no evidence for an increase of *AMY1* gene copy number (10, 31). By incorporating additional archaic hominin genomes that have not been previously analyzed for *AMY1* copy number (total:  $n = 8$ ), we have now detected signatures of *AMY1* gene duplication in a total of four archaic hominins.

The *AMY1* duplication in these archaic hominins could be postulated by four scenarios. First, it is possible that because of the complexity of the amylase locus and the complications inherent in archaic hominin genomic sequencing, there may be a technical bias in our detection. However, the observation of duplications using two different approaches and in multiple genomes provides confidence in our results. Second, it is plausible that introgression into archaic hominins from humans may explain the presence of duplications in the former, especially in the light of recently revealed complex interactions and gene flow events between Neanderthals and present-day humans over time (32). The approach recently developed by Li and colleagues (32) will be valuable for more formally testing for introgression from and to Neanderthals, potentially underlying the origins of the observed *AMY1* duplications in Neanderthals. However,



**Fig. 3. Amylase gene duplication and the history of agriculture.** (A) The read depth of amylase locus spanning *RNPC3*, *AMY2B*, *AMY2A*, and *AMY1* genes (chr1:103494306 to 103668306) for GoyetQ56-1 (maroon line), a Neanderthal excavated in present-day Belgium, showing signatures of *AMY1* duplication, and for Denisova 11 (beige line), a hybrid hominin (Neanderthal and Denisovan) excavated from present-day Russia, showing signatures consistent with an ancestral single-copy *AMY1* haplotype. The maroon and beige lines indicate average read depths with 5-kbp window and 1-kbp step for each sample. The average read depth of each 5-kbp window was normalized by the average read depth of the *RNPC3* gene. Only uniquely mapped reads were used for this visualization. (B) A world map displaying the locations of ancient human samples. Sample locations are indicated with an "X", with corresponding hexagons showing the estimated *AMY1* copy number. Carbon dating (number of years before present; BP) estimated for each sample is indicated in blue. (C) Amylase copy number estimations from Europeans who were farmers (beige)

and hunter-gatherers (maroon). Samples are binned on the x axis according to three time periods; preagriculture (before the transition to agriculture, >9000 yr B.P.), during the transition (~9000 to 4000 yr B.P.), and posttransition (complete transition to agriculture, ~4000 to 1000 yr B.P.). The shape of data points corresponds to sample dating estimates (Xs for more than 9000 yr B.P., circles for 9000 to 4000 yr B.P., and triangles for 4000 to 1000 yr B.P.). The inset of the right panel shows nonparametric regression lines for *AMY1* copy number across time, for hunter-gatherers (maroon) and farmers (beige), with confidence intervals in lighter corresponding color, respectively. (D) Zoomed-in map showing the spread of agriculture into Europe from Asia. Major agricultural footholds are indicated by wheat pictograms. Tan arrows show general trends of human agricultural migration throughout Europe, with predicted time periods (58, 59). Ancient human samples that were analyzed for amylase copy number are annotated with shapes and colors for time period and lifestyle, respectively. [Figure created with BioRender]

the current conservative filtering approach ends up filtering ~89% of the bases in the amylase locus. Thus, the specific signals of introgression, even if they exist, remain hidden. Third, it is plausible that the duplication evolved independently in the archaic hominin lineage. However, we argue that two independent duplications (one in humans and another in archaic hominins) in less than a million years are unlikely, given that initial duplications from single *AMY1* copy haplotypes are rare in non-human primates (4) (fig. S16). The fourth and, in our opinion, the most plausible scenario is that the *AMY1* gene might already have been

copy number variable before the human-Neanderthal/Denisovan divergence [~800 ka (33)], albeit to a limited extent as compared to what is observed in present-day humans. Overall, our results suggest a complex history of *AMY1* duplications, which will be further scrutinized as more high-coverage archaic hominin genomes become available.

#### The *AMY1* copy number has increased in the past 4000 years among European farmers

To explore the changes in frequency of *AMY1* copy number since the migration out of Africa ~60 ka (34), we analyzed the genomes of 68

ancient human genomes (table S12). The oldest genome analyzed was the Ust'-Ishim sample from Siberia (~45,000 years before the present (yr B.P.)), which has six *AMY1* gene copies per diploid cell. Similarly, the oldest modern human from Europe, the Peștera Muierii sample from Romania (~34,000 yr B.P.) has eight *AMY1* gene copies per diploid cell. These copy numbers indicate that high *AMY1* gene copies (defined here as ≥6 copies per diploid cell) had already spread across Eurasia as far back as ~45,000 yr B.P. (35) (Fig. 3B).

We next analyzed the ancient human genomes in the context of agricultural development



and found a general trend where the *AMY1* gene copy number is significantly higher among samples excavated from archaeologically agricultural contexts compared to those from hunter-gatherer contexts ( $P = 0.023$ ; fig. S22). To further investigate this trend, we examined ancient human genomes from Europe, where we have a clear timeline of the Neolithic transition. Specifically, the Neolithic transition of Europe started with the cultural and genetic influx from Anatolian farmers around ~9000 yr B.P. and progressed into Northwestern Europe by 5000 yr B.P. (36), with small, isolated groups of hunter-gatherers persisting until at least 4000 B.P. (37). Our analysis encompasses the geographic and temporal span of this transition. The oldest sample in our dataset from an agricultural European archaeological context is AKT16 from Anatolia, dated to 8547 yr B.P. (38), and the youngest sample from a hunter-gatherer archaeological context is SRA62 from Ireland, dated to 5215 yr B.P. (39). On the basis of the dates and archaeological context of the samples described in their respective studies (table S12), we parsed the ancient European human genomes into the following periods for visualization purposes: (i) the preagricultural period (> 9000 yr B.P.), where all our samples are from hunter-gatherers; (ii) the agricultural transition period (9000 to 4000 yr B.P.), representing the long period of transition to agriculture in Europe where both hunter-gatherer and agriculturalist groups coexisted; and (iii) the postagricultural period, during which all of Europe has completely transitioned into an agriculturalist life style (<4000 yr B.P.) (38–59). We found that preagricultural genomes already harbored four to eight *AMY1* copies per diploid cell (table S12). We also observed a consistent and significant increase in *AMY1* copy numbers across these periods ( $P = 0.005$ ; Fig. 3, C and D) and found similar trends for the *AMY2A* genes (Fig. 3C) and non-European samples (figs. S23 and S24). These findings support the notion that amylase haplotypes with higher numbers of amylase gene copies have increased in frequency over the past 4000 years. We found no significant differences in *AMY1* copy number between samples excavated from agricultural versus hunter-gatherer archaeological contexts during the “agricultural transition” period when farmers and hunter-gatherers shared the same habitat.

Taken together, these findings are consistent with the idea that either neutral evolution or weak-adaptive forces, perhaps due to preagricultural experimentation with food processing techniques (60), such as flour production from wild cereals (61), have retained the wide range of standing *AMY1* copy number variation in preagricultural Europe. The gradual increase in starch availability as Europe transitioned to an agricultural lifestyle (62) may

underlie selective forces acting on high-copy number haplotypes, explaining the increase of *AMY1* copy number in late postagriculture European populations. Given that the exact mechanisms through which *AMY1* copy number may confer adaptive advantage are unknown and the dietary intake even within agricultural and hunter-gatherer groups is highly diverse (63), it is challenging to reach definitive conclusions that the high *AMY1* copy number had an adaptive role during the agricultural transition. Additional samples with comprehensive archaeological context will help further elucidate the potential adaptive role of amylase in relation to the specific composition of ancient diets, as well as demographic considerations such as population replacements and gene flow from different groups that shaped the European Neolithic gene pool.

#### Multiple mutational mechanisms underlie the amylase copy number variation

We next investigated the likely mutational origins of the present-day amylase haplotypes. To do this, we first examined the relationship among the four most common haplotypes, H3<sup>T</sup>.1, H3<sup>T</sup>.2, H3<sup>T</sup>.4, and H1<sup>A</sup>.1. Building upon the previously reported link between one-copy and three-copy haplotypes (64) and the fact that we observe only three types of *AMY1* segments, we propose an evolutionary model linking the ancestral chimpanzee-like haplotype (H1<sup>A</sup>-like) to the common three-copy haplotypes (H3<sup>T</sup>.1 and H3<sup>T</sup>.2) (fig. S25). According to this model, the initial duplications of *AMY1* starting from H1<sup>A</sup>-like ancestral haplotype led to the emergence of H3<sup>T</sup> haplotypes. Given that multiple mutational steps are required to move from one-copy common haplotypes to other common haplotypes and that we do not observe any intermediate haplotypes in the present-day human genomes, the duplication from one-copy to three-copy haplotypes is likely to have occurred only once in the human lineage. This is also supported by the absence of duplications in the amylase segments within the H1<sup>A</sup>.1 haplotype, which would impede nonallelic homologous recombination (NAHR) events (65). By contrast, H3<sup>T</sup>.1 haplotypes harbor copies of identical and unidirectional sequences (e.g., *AMY2A.2* segments), providing an ideal template for recurrent NAHR events leading to the diverse haplotypes that we see today.

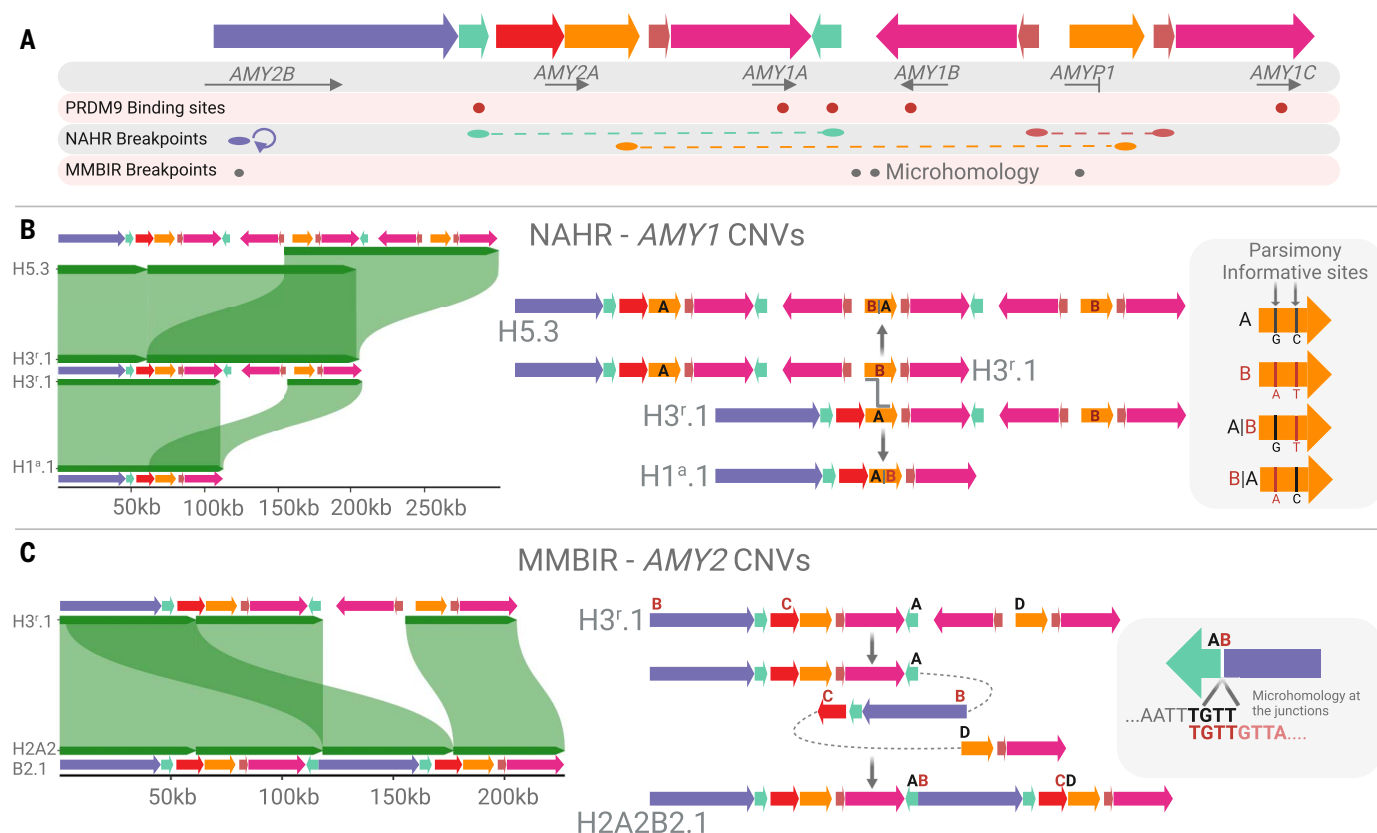
To further investigate the mutational relationships between the H3<sup>T</sup> haplotypes and other extant haplotypes, we used dotplots and sequence alignments to delineate the breakpoints of structural differences in amylase haplotypes (66, 67). In parallel, we conducted a scan for PRDM9-binding motifs within the amylase locus to pinpoint possible recombination sites (Fig. 4A, fig. S26, and table S13). By integrating all these observations, we were

able to construct a putative evolutionary path with the fewest plausible mutational steps that can explain the origin of present-day amylase haplotypes starting from NAHR-prone H3<sup>T</sup> haplotypes (figs. S27 and S28) (66). Our proposed evolutionary model of mutational events is consistent with our hypothesis regarding the central role of the H3<sup>T</sup> haplotypes in seeding extant haplotype variation and offers three major insights into the evolution of the amylase locus in humans.

First, we found evidence for recurrent NAHR events among common haplotypes (H3<sup>T</sup>.1, H3<sup>T</sup>.2, and H3<sup>T</sup>.4) harboring the *AMY1A* and *AMY1B* segments with breakpoints in the *AMY2A.2* segment. These NAHR events, which may occur among different haplotype combinations, could concurrently result in the duplication and deletion of two *AMY1* gene copies (e.g., Fig. 4B and fig. S28) (66). Therefore, although other less likely scenarios are possible, NAHR-based deletion and duplications underlie copy number variation of the *AMY1* segments and thus *AMY1* genes. Specifically, as this proposed mechanism always adds or deletes two copies of the *AMY1* genes, our finding explains how most human diploid genomes harbor even-numbered *AMY1* gene copies (7) (fig. S29). Thus, the majority of H1<sup>A</sup> haplotypes among present-day humans may have predominantly originated from H3<sup>T</sup> haplotypes. If true, this hypothesis explains the homoplastic occurrence of H1<sup>A</sup> haplotypes across the amylase phylogenetic tree (fig. S19C) and the lack of an out-of-Africa signal in H1<sup>A</sup> nucleotide diversity, which would be expected if the H1<sup>A</sup> haplotypes arose before out-of-Africa migrations (table S14). We argue that although the H1<sup>A</sup> haplotype is structurally nearly identical to the ancestral (chimpanzee) human amylase haplotype, the H1<sup>A</sup> haplotypes observed in extant humans have arisen recurrently from H3<sup>T</sup> haplotypes.

Second, we characterized three microhomology-mediated break-induced replication events and identified the accompanying microhomologies at the breakpoint junctions (Fig. 4C) (66). Even though these three rearrangements constitute only five alleles (H2A2B2.1, H3<sup>T</sup>.6, and H3B2.1) (~4%), they hold substantial biological relevance because H2A2B2.1 and H3B2.1 harbor duplications of the *AMY2* genes. Taken together, different mechanisms drive the copy number variation of the salivary *AMY1* genes and pancreatic *AMY2* genes, with the lower copy number variation of *AMY2* genes explained by the slower rate of nonrecurrent microhomology-mediated break-induced replication events (68).

Third, recurrent NAHR-mediated inversion events at the amylase locus, similar to those described previously (69), underpin the mutational connections between the common H3<sup>T</sup>.1, H3<sup>T</sup>.2, and H3<sup>T</sup>.4 haplotypes (fig. S30), as well as several other inversions among extant haplotypes



**Fig. 4. The evolutionary and mutational connections among common haplotypes.** (A) The structural variation breakpoints and recombination hotspots in the amylase locus. The colored arrows represent amylase segments. The PRDM9 binding sites are represented with red dots. The nonallelic homologous recombination (NAHR) breakpoints are represented with purple, green, and orange dots, and dashed lines. The microhomology-mediated break-induced replication breakpoints are represented with gray dots. (B) NAHR-mediated duplication and deletion of the *AMY1A-AMY1B* cluster. The *AMY2.2* (orange) segment serves as the

recombination substrate for the crossover, resulting in the duplication or deletion of the *AMY1A-AMY1B* cluster as illustrated in the middle panel. Chimeric *AMY2.2* segments have been identified, using parsimony informative sites within the *AMY2.2* segment (right panel). (C) Microhomology-mediated break-induced replication-based copy number gain resulting in the formation of H2A2B2.1. The middle panel shows the mutational mechanism. Four nucleotides of microhomology internal to the breakends were identified at the breakpoint junction (right panel). [Figure created using BioRender]

(66). Given that inversions underlie the structural differences between most of the amylase haplotypes, their functional and adaptive relevance presents an intriguing avenue for future research.

## Discussion

In this study, we have dissected the evolution of the amylase locus. First, we hypothesize that the initial duplications of the *AMY1* genes occurred once through multiple duplications, evolving from a one-copy ancestral haplotype to the three-copy present-day haplotypes frequently observed in our dataset. Analysis of archaic hominin genomes suggests that these initial duplications may have occurred well before the split of the human-Neanderthal/Denisovan. This observation is concordant with the recent evidence of Neanderthal starch consumption (70), and perhaps the availability of cooked starch in archaic hominins made possible through the domestication of fire (71) (Fig. 5).

Second, we hypothesize that selection acted on abundant standing *AMY1* copy number

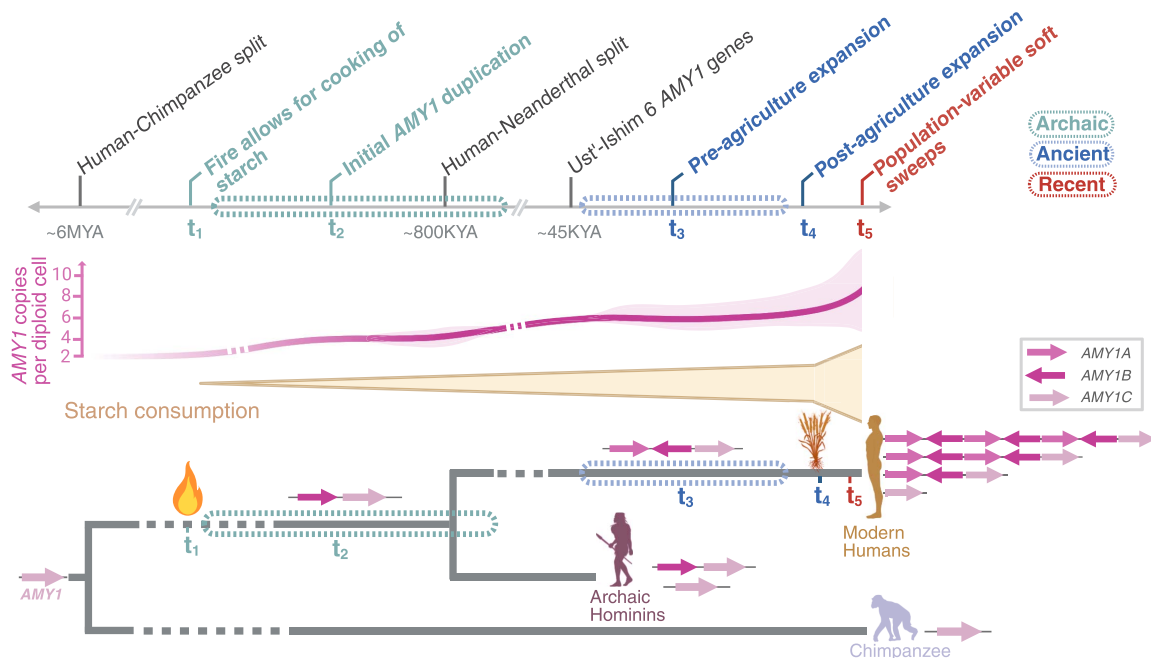
variation at this locus rather than on de novo variants. We observed a wide range of *AMY1* gene copy number variation (~three to nine copies) in samples that predate agriculture. We further found that late agricultural populations consistently harbored amylase haplotypes with higher *AMY1* copy numbers (Fig. 5). However, the lack of linkage disequilibrium between the flanking single-nucleotide polymorphisms (SNPs) and the *AMY1* copy number variation prevented us from conducting haplotype-based tests of positive selection acting on *AMY1* copy number. We hypothesize that partial soft sweeps involving preexisting amylase haplotypes may have influenced *AMY1* copy number variation in relation to historical starch consumption patterns in different populations. The effects of amylase gene duplications on taste preferences and starch metabolism may have also predisposed humans to prefer and tolerate the consumption of wild grains, as reported for the Mesolithic groups in the Balkans (72), facilitating the adoption of

starch-rich diets and the eventual domestication of these plants. Overall, our results support the complex narrative of the transition to agriculture, which includes the replacement of Western hunter-gatherers by Anatolian farmers in Europe (73), potentially bringing with them amylase haplotypes that harbor higher *AMY1* gene copies. Similarly, transient interactions between hunter-gatherer and agricultural groups (38) could explain the similar copy numbers observed between these groups during the transition period. Given the unknowns concerning how *AMY1* copy number affects metabolic function and its adaptive value under different life histories and starch consumption levels, it is challenging to draw definitive conclusions. Further, the agricultural transition varied across different regions and periods, involving diverse types of starches, which likely influenced the putative role of *AMY1* copy number in local adaptation. Anthropologically contextualized studies in Indigenous populations, such as the Andeans



**Fig. 5. An evolutionary model of the human amylase genes and resulting hypotheses.**

Top: A timeline of human amylase locus evolution based on the results of this study, with relevant events indicated on top. Middle: A schematic view showing the increase in *AMY1* copy number variation and the mean number of *AMY1* copies in present-day human populations throughout history as starch consumption increases. Bottom: A phylogenetic representation of the hypothesized amylase duplication timeline. [Figure created with BioRender]



(74, 75), could provide further insights into the relationships between dietary practices, metabolic outcomes, and amylase genetic variation.

Third, we found that NAHR and microhomology-mediated break-induced replication underlie the copy number variation of *AMY1* and *AMY2* genes, respectively, explaining the different rates of their evolution. The extremely high rate of structural variation due to NAHR led to substantial *AMY1* copy number variation with distinct mutational propensities. For example, our evolutionary model involving common H3<sup>r</sup> haplotypes suggests recurrent NAHR events mediated by highly similar sequences (> 99%), resulting in duplications or deletions of both *AMY1A* and *AMY1B* genes. By contrast, the H1<sup>a</sup> haplotypes harboring a single *AMY1* copy and divergent *AMY2* genes are less susceptible to NAHR events. Therefore, the mutation types and rates at the amylase locus may differ depending on extant haplotype variation in a population, especially in bottlenecked populations such as Indigenous Americans (76). It is a distinct possibility that a bottlenecked population ends up with a very high frequency of H1<sup>a</sup> due to drift. In this case, the absence of segments with highly similar sequences within H1<sup>a</sup> haplotype would mitigate recurrent NAHR events, resulting in a slower accumulation of variation in this population. By contrast, if one of the larger amylase haplotypes were to become prevalent as a result of drift, the rate of variation would increase exponentially. Within this general context, one interesting question for future work is whether

larger amylase haplotypes experience negative selection due to increased genomic instability.

Taken together, our study underscores how gene duplications in early human history provided the genetic foundation for dietary flexibility during agricultural innovations, contributing to modern human evolution.

### Materials and methods summary

In this study, samples ( $n = 98$ ) from the 1000 Genomes Project (77), as part of the datasets from the Human Genome Structural Variation Consortium (HGSVC) (20), and Human Pan-genome Reference Consortium (HPRC) (78), and the Genome In a Bottle (79) were analyzed. For each sample, datasets from Bionano Genomics optical genome mapping and PacBio HiFi sequencing were used. A ~212.5-kbp region on chromosome 1 containing amylase genes was characterized using optical genome mapping and blastn (BLASTN 2.9.0+) alignments (80). The amylase segments were identified and validated across human and nonhuman primate samples. De novo assembly of optical genome maps was performed using Bionano Solve v3.5, followed by alignment and haplotype reconstruction using Bionano Access v1.7 software (81). PacBio HiFi long-read sequencing data were also de novo assembled using hifiasm 0.16.1-r375 (82). Population-specific patterns in the *AMY1* gene copy numbers were evaluated using the Kruskal-Wallis test. The singleton proportions at the amylase locus were compared to those of genome-wide tandem repeats. Rarefaction analysis was conducted

to assess haplotype diversity at the amylase locus.

To predict amylase gene coordinates and copy number, we used Exonerate v2.4.0 “protein2genome” tool (83) with default settings and a maximum intron length of 20 kbp, employing amylase protein sequences from UniProt (84). Overlapping hits were clustered, and the best hit was selected. Hits translating into full-length polypeptides (511 amino acids) were kept. Any conflicts with predicted amylase gene copy numbers were manually curated. De novo genome annotations were obtained using the T2T-chm13 gff annotation by liftoff v1.6.3 (85) and compared with the homology-based annotations. ddPCR with custom primers was used to validate *AMY1* copy numbers in human DNA samples, with Hind III digestion prior to ddPCR. Alignments of amylase segments and coding sequences were constructed using MAFFT v7.310 (86) and PAGAN v1.53 (87), incorporating chimpanzee sequences as an out-group. Phylogenetic trees were reconstructed using IQ-TREE v2.2.0 (88) and visualized with FigTree v1.4.4. dN/dS ratios were estimated using PAML v4.10.6 codeml (23), testing for gene type-specific selection pressures. Functional analysis of amylase protein sequences was performed with CLUSTALO v1.2.3 (89) and AlphaMissense (90). The expression data of amylase genes were analyzed using transcripts per million (TPM) data from the Genotype-Tissue Expression (GTEx) portal (26), focusing on three tissues with the highest expression values. Because major salivary glands were not included in GTEx,

supplementary TPM data from Saitou *et al.* (91) were used, considering AMY1A, AMY1B, and AMY1C transcripts collectively as representative of AMY1 expression. Sequence origins and breakpoints in AMY1 segments were identified using BLAT v37x1 (92), nucmer v3.1, and mummerplot v3.5 (93, 94). AMY1 duplication events were dated using BEASTv 2.7.5 (95), without tree topology optimization for H3<sup>F.1</sup> and H3<sup>F.2</sup> haplotypes independently, and results were visualized with FigTree v1.4.4 (<http://tree.bio.ed.ac.uk/software/figtree/>).

We used GeneToCN (29) to determine the amylase copy number in 68 ancient human and 8 archaic hominin genomes. Raw sequencing data were sourced from the European Nucleotide Archive. GeneToCN, using 25-mers, calculated the copy numbers of AMY1, AMY2A, and AMY2B genes. This analysis was validated with ddPCR for present-day humans. AMY1, AMY2A, and AMY2B copy numbers were compared between hunter-gatherer and agricultural lifestyles using the Wilcoxon rank sum test. The copy number trends over time were also analyzed. To assess coverage bias, eight ancient genomes were downsampled, and their AMY1 copy numbers were recalculated and compared to full genome estimates. This analysis considered potential batch effects and sample diversity.

We obtained raw sequencing reads of archaic hominin genomes ( $n = 38$ ) from public datasets, removed adapters, and merged overlapping paired-end reads using leeHom v1.2.17 (96). For the Chagyrskaya Neanderthal genome, preprocessed reads were used. Clean reads were mapped to the GRCh38 reference genome using bwa v0.7.17 aln (97). GC bias was assessed and corrected with deepTools v3.5.1 (98), resulting in the exclusion of 29 genomes. Average read-depth at both genome-wide level and the amylase locus was calculated with mosdepth v0.3.5 (99), excluding one more sample for low coverage ( $<1\times$ ). GC-corrected reads were realigned to T2T-chm13 and GRCh38 assemblies with bwa v0.7.17 mem (97) and further aligned to a reference containing a single AMY1 gene copy. Amylase gene copy numbers were calculated from the average read-depth of each amylase gene normalized by the genome-wide average.

To identify potential PRDM9 binding sites in the H3<sup>F.1</sup> haplotype, we used FIMO v5.5.4 from the MEME package (100) with the degenerate 13-bp motif “CCNCCNTNCCNC” and a nonredundant DNA database background frequency matrix. Hits with a FIMO score above 10 and a  $P$  value below 0.00011 were considered. To investigate structural variation mechanisms, we aligned haplotypes using nucmer v3.1 (93, 94) and visualized alignments with mummerplot and SVbyEye (<https://github.com/daewoooo/SVbyEye/tree/master>). Potential breakpoints were identified and 20-kbp

sequences upstream and downstream were aligned using MAFFT v7.522 (86). We evaluated replication and DNA recombination-based processes and inferred NAHR when breakpoints were within paralogous segments and lacked replication-based sequence motifs.

## REFERENCES AND NOTES

- G. H. Perry *et al.*, Diet and the evolution of human amylase gene copy number variation. *Nat. Genet.* **39**, 1256–1260 (2007). doi: [10.1038/ng2123](https://doi.org/10.1038/ng2123); pmid: [17828263](https://pubmed.ncbi.nlm.nih.gov/17828263/)
- T. Nishide *et al.*, Sequences of cDNAs for human salivary and pancreatic  $\alpha$ -amylases. *Gene* **28**, 263–270 (1984). doi: [10.1016/0378-1119\(84\)90265-8](https://doi.org/10.1016/0378-1119(84)90265-8); pmid: [6610603](https://pubmed.ncbi.nlm.nih.gov/6610603/)
- C. Peyrot des Gachons, P. A. S. Breslin, Salivary Amylase: Digestion and Metabolic Syndrome. *Curr. Rev. Oral Biol. Med.* **16**, 102 (2016). doi: [10.1007/s11892-016-0794-7](https://doi.org/10.1007/s11892-016-0794-7); pmid: [27640169](https://pubmed.ncbi.nlm.nih.gov/27640169/)
- P. Pajic *et al.*, Independent amylase gene copy number bursts correlate with dietary preferences in mammals. *eLife* **8**, e44628 (2019). doi: [10.7554/eLife.44628](https://doi.org/10.7554/eLife.44628); pmid: [31084707](https://pubmed.ncbi.nlm.nih.gov/31084707/)
- M. H. Meisler, C. N. Ting, The remarkable evolutionary history of the human amylase genes. *Crit. Rev. Oral Biol. Med.* **4**, 503–509 (1993). doi: [10.1177/10454411930040033501](https://doi.org/10.1177/10454411930040033501); pmid: [7690604](https://pubmed.ncbi.nlm.nih.gov/7690604/)
- P. C. Groot, W. H. Mager, R. R. Frants, Interpretation of polymorphic DNA patterns in the human alpha-amylase multigene family. *Genomics* **10**, 779–785 (1991). doi: [10.1016/0888-7543\(91\)90463-0](https://doi.org/10.1016/0888-7543(91)90463-0); pmid: [1679752](https://pubmed.ncbi.nlm.nih.gov/1679752/)
- C. L. Usher *et al.*, Structural forms of the human amylase locus and their relationships to SNPs, haplotypes and obesity. *Nat. Genet.* **47**, 921–925 (2015). doi: [10.1038/ng.3340](https://doi.org/10.1038/ng.3340); pmid: [26098870](https://pubmed.ncbi.nlm.nih.gov/26098870/)
- A. C. Poole *et al.*, Human Salivary Amylase Gene Copy Number Impacts Oral and Gut Microbiomes. *Cell Host Microbe* **25**, 553–564.e7 (2019). doi: [10.1016/j.chom.2019.03.001](https://doi.org/10.1016/j.chom.2019.03.001); pmid: [30974084](https://pubmed.ncbi.nlm.nih.gov/30974084/)
- S. Mathieson, I. Mathieson, FADS1 and the Timing of Human Adaptation to Agriculture. *Mol. Biol. Evol.* **35**, 2957–2970 (2018). doi: [10.1093/molbev/msy180](https://doi.org/10.1093/molbev/msy180); pmid: [30272210](https://pubmed.ncbi.nlm.nih.gov/30272210/)
- G. H. Perry, L. Kistler, M. A. Kelaita, A. J. Sams, Insights into hominin phenotypic and dietary evolution from ancient DNA sequence data. *J. Hum. Evol.* **79**, 55–63 (2015). doi: [10.1016/j.jhevol.2014.10.018](https://doi.org/10.1016/j.jhevol.2014.10.018); pmid: [25563409](https://pubmed.ncbi.nlm.nih.gov/25563409/)
- D. Bolognini *et al.*, Global diversity, recurrent evolution, and recent selection on amylase structural haplotypes in humans. *bioRxiv*, doi: [10.1101/2024.02.07.579378](https://doi.org/10.1101/2024.02.07.579378) (2024).
- C. E. Inchley *et al.*, Selective sweep on human amylase genes postdates the split with Neanderthals. *Sci. Rep.* **6**, 37198 (2016). doi: [10.1038/srep37198](https://doi.org/10.1038/srep37198); pmid: [27853181](https://pubmed.ncbi.nlm.nih.gov/27853181/)
- W. Demareel *et al.*, The 22q11 low copy repeats are characterized by unprecedented size and structural variability. *Genome Res.* **29**, 1389–1401 (2019). doi: [10.1101/gr.248682.119](https://doi.org/10.1101/gr.248682.119); pmid: [31481461](https://pubmed.ncbi.nlm.nih.gov/31481461/)
- Y. Mostovoy *et al.*, Genomic regions associated with microdeletion/microduplication syndromes exhibit extreme diversity of structural variation. *Genetics* **217**, iyaa038 (2021). doi: [10.1093/genetics/iyaa038](https://doi.org/10.1093/genetics/iyaa038); pmid: [33724415](https://pubmed.ncbi.nlm.nih.gov/33724415/)
- F. Yilmaz *et al.*, High level of complexity and global diversity of the 3q29 locus revealed by optical mapping and long-read sequencing. *Genome Med.* **15**, 35 (2023). doi: [10.1186/s13073-023-01184-5](https://doi.org/10.1186/s13073-023-01184-5); pmid: [37165454](https://pubmed.ncbi.nlm.nih.gov/37165454/)
- F. Yilmaz, Amylase Segments, Haplotypes, Liftoff Gene Annotations and OrthogonalValidation-MoleculeSupport, Zenodo (2024); <https://zenodo.org/doi/10.5281/zenodo.13247145>
- S. Nurk *et al.*, The complete sequence of a human genome. *Science* **376**, 44–53 (2022). doi: [10.1126/science.abj6987](https://doi.org/10.1126/science.abj6987); pmid: [35357919](https://pubmed.ncbi.nlm.nih.gov/35357919/)
- A. Harpak, A. Bhaskar, J. K. Pritchard, Mutation Rate Variation is a Primary Determinant of the Distribution of Allele Frequencies in Humans. *PLoS Genet.* **12**, e1006489 (2016). doi: [10.1371/journal.pgen.1006489](https://doi.org/10.1371/journal.pgen.1006489); pmid: [27977673](https://pubmed.ncbi.nlm.nih.gov/27977673/)
- M. Verbiest *et al.*, Mutation and selection processes regulating short tandem repeats give rise to genetic and phenotypic diversity across species. *J. Evol. Biol.* **36**, 321–336 (2023). doi: [10.1111/jeb.14106](https://doi.org/10.1111/jeb.14106); pmid: [36289560](https://pubmed.ncbi.nlm.nih.gov/36289560/)
- P. Ebert *et al.*, Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* **372**, eabf7117 (2021). doi: [10.1126/science.abf7117](https://doi.org/10.1126/science.abf7117); pmid: [33632895](https://pubmed.ncbi.nlm.nih.gov/33632895/)

- P. Pajic, Digital Droplet PCR Protocol for Copy Number Variation Estimation, Zenodo (2024). <https://doi.org/10.5281/ZENODO.13155899>
- K. Kim, AmylaseSegmentCDSAlignments, DatingAMY1, DNdSAnalysis, OriginBreakPointAMY1, Protein2genome, GeneAnnotation, TandemRepeatAnalysis, Zenodo (2024); <https://doi.org/10.5281/zenodo.13169856>
- Z. Yang, PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**, 1586–1591 (2007). doi: [10.1093/molbev/msm088](https://doi.org/10.1093/molbev/msm088); pmid: [17483113](https://pubmed.ncbi.nlm.nih.gov/17483113/)
- P. Pajic, Functional analysis data (AlphaMissense) for Amylase proteins, Zenodo (2024); <https://doi.org/10.5281/ZENODO.13156521>
- M. Uhlen *et al.*, Tissue-based map of the human proteome. *Science* **347**, 1260419 (2015). doi: [10.1126/science.1260419](https://doi.org/10.1126/science.1260419); pmid: [25613900](https://pubmed.ncbi.nlm.nih.gov/25613900/)
- F. Aguet *et al.*, The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science* **369**, 1318–1330 (2020). doi: [10.1126/science.aaz1776](https://doi.org/10.1126/science.aaz1776); pmid: [32913098](https://pubmed.ncbi.nlm.nih.gov/32913098/)
- D. Carpenter *et al.*, Obesity, starch digestion and amylase: Association between copy number variants at human salivary (AMY1) and pancreatic (AMY2) amylase genes. *Hum. Mol. Genet.* **24**, 3472–3480 (2015). doi: [10.1093/hmg/ddv098](https://doi.org/10.1093/hmg/ddv098); pmid: [25788522](https://pubmed.ncbi.nlm.nih.gov/25788522/)
- M. R. Vollger *et al.*, Increased mutation and gene conversion within human segmental duplications. *Nature* **617**, 325–334 (2023). doi: [10.1038/s41586-023-05895-y](https://doi.org/10.1038/s41586-023-05895-y); pmid: [37165237](https://pubmed.ncbi.nlm.nih.gov/37165237/)
- F.-D. Pajuste, M. Remm, GeneToCN: An alignment-free method for gene copy number estimation directly from next-generation sequencing reads. *Sci. Rep.* **13**, 17765 (2023). doi: [10.1038/s41598-023-44636-z](https://doi.org/10.1038/s41598-023-44636-z); pmid: [37853040](https://pubmed.ncbi.nlm.nih.gov/37853040/)
- K. Scheer, Zenodo Downsampling Information, Zenodo (2024); <https://doi.org/10.5281/ZENODO.13251208>
- K. Prüfer *et al.*, The complete genome sequence of a Neanderthal from the Altai Mountains. *Nature* **505**, 43–49 (2014). doi: [10.1038/nature12886](https://doi.org/10.1038/nature12886); pmid: [24352235](https://pubmed.ncbi.nlm.nih.gov/24352235/)
- L. Li, T. J. Comi, R. F. Bierman, J. M. Akey, Recurrent gene flow between Neanderthals and modern humans over the past 200,000 years. *Science* **385**, eadi1768 (2024). doi: [10.1126/science.adil768](https://doi.org/10.1126/science.adil768); pmid: [38991054](https://pubmed.ncbi.nlm.nih.gov/38991054/)
- A. Gómez-Robles, Dental evolutionary rates and its implications for the Neanderthal-modern human divergence. *Sci. Adv.* **5**, eaaw1268 (2019). doi: [10.1126/sciadv.aaw1268](https://doi.org/10.1126/sciadv.aaw1268); pmid: [31106274](https://pubmed.ncbi.nlm.nih.gov/31106274/)
- S. Mallick *et al.*, The Simons Genome Diversity Project: 300 genomes from 142 diverse populations. *Nature* **538**, 201–206 (2016). doi: [10.1038/nature18964](https://doi.org/10.1038/nature18964); pmid: [27654912](https://pubmed.ncbi.nlm.nih.gov/27654912/)
- L. Vallini *et al.*, The Persian plateau served as hub for Homo sapiens after the main out of Africa dispersal. *Nat. Commun.* **15**, 1882 (2024). doi: [10.1038/s41467-024-46161-7](https://doi.org/10.1038/s41467-024-46161-7); pmid: [38528002](https://pubmed.ncbi.nlm.nih.gov/38528002/)
- M. Furholt, Mobility and Social Change: Understanding the European Neolithic Period after the Archaeogenetic Revolution. *J. Archaeol. Res.* **29**, 481–535 (2021). doi: [10.1007/s10814-020-09153-x](https://doi.org/10.1007/s10814-020-09153-x)
- L. Saag *et al.*, Extensive Farming in Estonia Started through a Sex-Biased Migration from the Steppe. *Curr. Biol.* **27**, 2185–2193.e6 (2017). doi: [10.1016/j.cub.2017.06.022](https://doi.org/10.1016/j.cub.2017.06.022); pmid: [28712569](https://pubmed.ncbi.nlm.nih.gov/28712569/)
- N. Marchi *et al.*, The genomic origins of the world's first farmers. *Cell* **185**, 1842–1859.e18 (2022). doi: [10.1016/j.cell.2022.04.008](https://doi.org/10.1016/j.cell.2022.04.008); pmid: [35561686](https://pubmed.ncbi.nlm.nih.gov/35561686/)
- L. M. Cassidy *et al.*, A dynastic elite in monumental Neolithic society. *Nature* **582**, 384–388 (2020). doi: [10.1038/s41586-020-2378-6](https://doi.org/10.1038/s41586-020-2378-6); pmid: [32555485](https://pubmed.ncbi.nlm.nih.gov/32555485/)
- J.-P. Bocquet-Appel, S. Naji, M. Vander Linden, J. Kozłowski, Understanding the rates of expansion of the farming system in Europe. *J. Archaeol. Sci.* **39**, 531–546 (2012). doi: [10.1016/j.jas.2011.10.010](https://doi.org/10.1016/j.jas.2011.10.010)
- K. Wang *et al.*, High-coverage genome of the Tyrolean Iceman reveals unusually high Anatolian farmer ancestry. *Cell Genomics* **3**, 100377 (2023). doi: [10.1016/j.xgen.2023.100377](https://doi.org/10.1016/j.xgen.2023.100377); pmid: [37719142](https://pubmed.ncbi.nlm.nih.gov/37719142/)
- A. Seguin-Orlando *et al.*, Heterogeneous Hunter-Gatherer and Steppe-Related Ancestries in Late Neolithic and Bell Beaker Genomes from Present-Day France. *Curr. Biol.* **31**, 1072–1083.e10 (2021). doi: [10.1016/j.cub.2020.12.015](https://doi.org/10.1016/j.cub.2020.12.015); pmid: [33434506](https://pubmed.ncbi.nlm.nih.gov/33434506/)
- E. Svensson *et al.*, Genome of Peștera Muierii skull shows high diversity and low mutational load in pre-glacial Europe. *Curr. Biol.* **31**, 2973–2983.e9 (2021). doi: [10.1016/j.cub.2021.04.045](https://doi.org/10.1016/j.cub.2021.04.045); pmid: [34010592](https://pubmed.ncbi.nlm.nih.gov/34010592/)

44. B. Ariano *et al.*, Ancient Maltese genomes and the genetic geography of Neolithic Europe. *Curr. Biol.* **32**, 2668–2680.e6 (2022). doi: [10.1016/j.cub.2022.04.069](https://doi.org/10.1016/j.cub.2022.04.069); pmid: [35588742](https://pubmed.ncbi.nlm.nih.gov/35588742/)
45. M. Rasmussen *et al.*, Ancient human genome sequence of an extinct Palaeo-Eskimo. *Nature* **463**, 757–762 (2010). doi: [10.1038/nature08835](https://doi.org/10.1038/nature08835); pmid: [20148029](https://pubmed.ncbi.nlm.nih.gov/20148029/)
46. I. Lazaridis *et al.*, Ancient human genomes suggest three ancestral populations for present-day Europeans. *Nature* **513**, 409–413 (2014). doi: [10.1038/nature13673](https://doi.org/10.1038/nature13673); pmid: [25230663](https://pubmed.ncbi.nlm.nih.gov/25230663/)
47. S. Schiffels *et al.*, Iron Age and Anglo-Saxon genomes from East England reveal British migration history. *Nat. Commun.* **7**, 10408 (2016). doi: [10.1038/ncomms10408](https://doi.org/10.1038/ncomms10408); pmid: [26783965](https://pubmed.ncbi.nlm.nih.gov/26783965/)
48. C. Gamba *et al.*, Genome flux and stasis in a five millennium transect of European prehistory. *Nat. Commun.* **5**, 5257 (2014). doi: [10.1038/ncomms6257](https://doi.org/10.1038/ncomms6257); pmid: [25334030](https://pubmed.ncbi.nlm.nih.gov/25334030/)
49. E. R. Jones *et al.*, Upper Palaeolithic genomes reveal deep roots of modern Eurasians. *Nat. Commun.* **6**, 8912 (2015). doi: [10.1038/ncomms9912](https://doi.org/10.1038/ncomms9912); pmid: [26567969](https://pubmed.ncbi.nlm.nih.gov/26567969/)
50. M. E. Allentoft *et al.*, Population genomics of post-glacial western Eurasia. *Nature* **625**, 301–311 (2024). doi: [10.1038/s41586-023-06865-0](https://doi.org/10.1038/s41586-023-06865-0); pmid: [38200295](https://pubmed.ncbi.nlm.nih.gov/38200295/)
51. P. Maisano Delser *et al.*, A curated dataset of modern and ancient high-coverage shotgun human genomes. *Sci. Data* **8**, 202 (2021). doi: [10.1038/s41597-021-00980-1](https://doi.org/10.1038/s41597-021-00980-1); pmid: [34349118](https://pubmed.ncbi.nlm.nih.gov/34349118/)
52. M. Srigyan *et al.*, Bioarchaeological evidence of one of the earliest Islamic burials in the Levant. *Commun. Biol.* **5**, 554 (2022). doi: [10.1038/s42003-022-03508-4](https://doi.org/10.1038/s42003-022-03508-4); pmid: [35672445](https://pubmed.ncbi.nlm.nih.gov/35672445/)
53. L. M. Cassidy *et al.*, Neolithic and Bronze Age migration to Ireland and establishment of the insular Atlantic genome. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 368–373 (2016). doi: [10.1073/pnas.1518445113](https://doi.org/10.1073/pnas.1518445113); pmid: [26712024](https://pubmed.ncbi.nlm.nih.gov/26712024/)
54. F. Sánchez-Quinto *et al.*, Megalithic tombs in western and northern Neolithic Europe were linked to a kindred society. *Proc. Natl. Acad. Sci. U.S.A.* **116**, 9469–9474 (2019). doi: [10.1073/pnas.1818037116](https://doi.org/10.1073/pnas.1818037116); pmid: [30988179](https://pubmed.ncbi.nlm.nih.gov/30988179/)
55. L. G. Simões *et al.*, Genomic ancestry and social dynamics of the last hunter-gatherers of Atlantic France. *Proc. Natl. Acad. Sci. U.S.A.* **121**, e2310545121 (2024). doi: [10.1073/pnas.2310545121](https://doi.org/10.1073/pnas.2310545121); pmid: [38408241](https://pubmed.ncbi.nlm.nih.gov/38408241/)
56. C. Valdiosera *et al.*, Four millennia of Iberian biomolecular prehistory illustrate the impact of prehistoric migrations at the far end of Eurasia. *Proc. Natl. Acad. Sci. U.S.A.* **115**, 3428–3433 (2018). doi: [10.1073/pnas.1717762115](https://doi.org/10.1073/pnas.1717762115); pmid: [29531053](https://pubmed.ncbi.nlm.nih.gov/29531053/)
57. S. S. Ebenesersdóttir *et al.*, Ancient genomes from Iceland reveal the making of a human population. *Science* **360**, 1028–1032 (2018). doi: [10.1126/science.aar2625](https://doi.org/10.1126/science.aar2625); pmid: [29853688](https://pubmed.ncbi.nlm.nih.gov/29853688/)
58. M. Lipson *et al.*, Parallel palaeogenomic transects reveal complex genetic history of early European farmers. *Nature* **551**, 368–372 (2017). doi: [10.1038/nature24476](https://doi.org/10.1038/nature24476); pmid: [29144465](https://pubmed.ncbi.nlm.nih.gov/29144465/)
59. N. E. Altınışık *et al.*, A genomic snapshot of demographic and cultural dynamism in Upper Mesopotamia during the Neolithic Transition. *Sci. Adv.* **8**, eabo3609 (2022). doi: [10.1126/sciadv.abo3609](https://doi.org/10.1126/sciadv.abo3609); pmid: [36332018](https://pubmed.ncbi.nlm.nih.gov/36332018/)
60. A. Revedin *et al.*, Thirty thousand-year-old evidence of plant food processing. *Proc. Natl. Acad. Sci. U.S.A.* **107**, 18815–18819 (2010). doi: [10.1073/pnas.1006993107](https://doi.org/10.1073/pnas.1006993107); pmid: [20956317](https://pubmed.ncbi.nlm.nih.gov/20956317/)
61. D. R. Piperno, E. Weiss, I. Holst, D. Nadel, Processing of wild cereal grains in the Upper Palaeolithic revealed by starch grain analysis. *Nature* **430**, 670–673 (2004). doi: [10.1038/nature02734](https://doi.org/10.1038/nature02734); pmid: [15295598](https://pubmed.ncbi.nlm.nih.gov/15295598/)
62. J. Jovanović, R. C. Power, C. de Beccelièvre, G. Goude, S. Stefanović, Microbotanical evidence for the spread of cereal use during the Mesolithic-Neolithic transition in the Southeastern Europe (Danube Gorges): Data from dental calculus analysis. *J. Archaeol. Sci.* **125**, 105288 (2021). doi: [10.1016/j.jas.2020.105288](https://doi.org/10.1016/j.jas.2020.105288)
63. L. Betti *et al.*, Climate shaped how Neolithic farmers and European hunter-gatherers interacted after a major slowdown from 6,100 BCE to 4,500 BCE. *Nat. Hum. Behav.* **4**, 1004–1010 (2020). doi: [10.1038/s41562-020-0897-7](https://doi.org/10.1038/s41562-020-0897-7); pmid: [32632332](https://pubmed.ncbi.nlm.nih.gov/32632332/)
64. P. C. Groot *et al.*, Evolution of the human alpha-amylase multigene family through unequal, homologous, and inter- and intrachromosomal crossovers. *Genomics* **8**, 97–105 (1990). doi: [10.1016/0888-7543\(90\)90230-R](https://doi.org/10.1016/0888-7543(90)90230-R); pmid: [2081604](https://pubmed.ncbi.nlm.nih.gov/2081604/)
65. C. Karageorgiou, Amylase locus genomic analysis, Zenodo (2024); <https://doi.org/10.5281/ZENODO.13170898>.
66. See the supplementary text.
67. C. Karageorgiou, Amylase locus genomic analysis, Zenodo (2024); <https://doi.org/10.5281/ZENODO.13138560>.
68. F. Zhang *et al.*, The DNA replication FoSTeS/MMBIR mechanism can generate genomic, genic and exonic complex rearrangements in humans. *Nat. Genet.* **41**, 849–853 (2009). doi: [10.1038/ng.399](https://doi.org/10.1038/ng.399); pmid: [19543269](https://pubmed.ncbi.nlm.nih.gov/19543269/)
69. D. Porubsky *et al.*, Recurrent inversion polymorphisms in humans associate with genetic instability and genomic disorders. *Cell* **185**, 1986–2005.e26 (2022). doi: [10.1016/j.cell.2022.04.017](https://doi.org/10.1016/j.cell.2022.04.017); pmid: [35525246](https://pubmed.ncbi.nlm.nih.gov/35525246/)
70. J. A. Fellows Yates *et al.*, The evolution and changing ecology of the African hominid oral microbiome. *Proc. Natl. Acad. Sci. U.S.A.* **118**, e2021655118 (2021). doi: [10.1073/pnas.2021655118](https://doi.org/10.1073/pnas.2021655118); pmid: [33972424](https://pubmed.ncbi.nlm.nih.gov/33972424/)
71. C. Larbey, S. M. Mentzer, B. Liguoz, S. Wur, M. K. Jones, Cooked starchy food in hearths ca. 120 kya and 65 kya (MIS 5e and MIS 4) from Klasies River Cave, South Africa. *J. Hum. Evol.* **131**, 210–227 (2019). doi: [10.1016/j.jhevol.2019.03.015](https://doi.org/10.1016/j.jhevol.2019.03.015); pmid: [31182202](https://pubmed.ncbi.nlm.nih.gov/31182202/)
72. E. Cristiani *et al.*, Wild cereal grain consumption among Early Holocene foragers of the Balkans predates the arrival of agriculture. *eLife* **10**, e72976 (2021). doi: [10.7554/eLife.72976](https://doi.org/10.7554/eLife.72976); pmid: [34850680](https://pubmed.ncbi.nlm.nih.gov/34850680/)
73. G. M. Kilg *et al.*, The Demographic Development of the First Farmers in Anatolia. *Curr. Biol.* **26**, 2659–2666 (2016). doi: [10.1016/j.cub.2016.07.057](https://doi.org/10.1016/j.cub.2016.07.057); pmid: [27498567](https://pubmed.ncbi.nlm.nih.gov/27498567/)
74. J. Lindo *et al.*, The genetic prehistory of the Andean highlands 7000 years BP through European contact. *Sci. Adv.* **4**, eaau4921 (2018). doi: [10.1126/sciadv.aau4921](https://doi.org/10.1126/sciadv.aau4921); pmid: [30417096](https://pubmed.ncbi.nlm.nih.gov/30417096/)
75. K. Jorgensen, O. A. Garcia, M. Kiyamu, T. D. Brutsaert, A. W. Bigham, Genetic adaptations to potato starch digestion in the Peruvian Andes. *Am. J. Biol. Anthropol.* **180**, 162–172 (2023). doi: [10.1002/ajpa.24656](https://doi.org/10.1002/ajpa.24656)
76. J. V. Moreno-Mayer *et al.*, Terminal Pleistocene Alaskan genome reveals first founding population of Native Americans. *Nature* **553**, 203–207 (2018). doi: [10.1038/nature25173](https://doi.org/10.1038/nature25173); pmid: [29323294](https://pubmed.ncbi.nlm.nih.gov/29323294/)
77. G. R. Abecasis *et al.*, A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061–1073 (2010). doi: [10.1038/nature09534](https://doi.org/10.1038/nature09534); pmid: [20981092](https://pubmed.ncbi.nlm.nih.gov/20981092/)
78. W.-W. Liao *et al.*, A draft human pangome reference. *Nature* **617**, 312–324 (2023). doi: [10.1038/s41586-023-05896-x](https://doi.org/10.1038/s41586-023-05896-x); pmid: [37165242](https://pubmed.ncbi.nlm.nih.gov/37165242/)
79. J. M. Zook *et al.*, Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci. Data* **3**, 160025 (2016). doi: [10.1038/sdata.2016.25](https://doi.org/10.1038/sdata.2016.25); pmid: [27271295](https://pubmed.ncbi.nlm.nih.gov/27271295/)
80. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, D. J. Lipman, Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990). doi: [10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2); pmid: [2231712](https://pubmed.ncbi.nlm.nih.gov/2231712/)
81. H. Cao *et al.*, Rapid detection of structural variation in a human genome using nanochannel-based genome mapping technology. *Gigascience* **3**, 34 (2014). doi: [10.1186/2047-217X-3-34](https://doi.org/10.1186/2047-217X-3-34); pmid: [25671094](https://pubmed.ncbi.nlm.nih.gov/25671094/)
82. H. Cheng, G. T. Concepcion, X. Feng, H. Zhang, H. Li, Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat. Methods* **18**, 170–175 (2021). doi: [10.1038/s41592-020-01056-5](https://doi.org/10.1038/s41592-020-01056-5); pmid: [33526886](https://pubmed.ncbi.nlm.nih.gov/33526886/)
83. G. S. C. Slater, E. Birney, Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics* **6**, 31 (2005). doi: [10.1186/1471-2105-6-31](https://doi.org/10.1186/1471-2105-6-31); pmid: [15713233](https://pubmed.ncbi.nlm.nih.gov/15713233/)
84. UniProt Consortium, UniProt: A worldwide hub of protein knowledge. *Nucleic Acids Res.* **47**, D506–D515 (2019); pmid: [30395287](https://pubmed.ncbi.nlm.nih.gov/30395287/)
85. A. Shumate, S. L. Salzberg, Liftoft: Accurate mapping of gene annotations. *Bioinformatics* **37**, 1639–1643 (2020). doi: [10.1093/bioinformatics/btaa1016](https://doi.org/10.1093/bioinformatics/btaa1016); pmid: [3320174](https://pubmed.ncbi.nlm.nih.gov/3320174/)
86. K. Katoh, D. M. Standley, MAFFT multiple sequence alignment software version 7: Improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013). doi: [10.1093/molbev/mst010](https://doi.org/10.1093/molbev/mst010); pmid: [23296960](https://pubmed.ncbi.nlm.nih.gov/23296960/)
87. A. Löytynoja, A. J. Vilella, N. Goldman, Accurate extension of multiple sequence alignments using a phylogeny-aware graph algorithm. *Bioinformatics* **28**, 1684–1691 (2012). doi: [10.1093/bioinformatics/bts198](https://doi.org/10.1093/bioinformatics/bts198); pmid: [22531217](https://pubmed.ncbi.nlm.nih.gov/22531217/)
88. B. Q. Minh *et al.*, IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020). doi: [10.1093/molbev/msaa015](https://doi.org/10.1093/molbev/msaa015); pmid: [32011700](https://pubmed.ncbi.nlm.nih.gov/32011700/)
89. F. Sievers, D. G. Higgins, Clustal Omega for making accurate alignments of many protein sequences. *Protein Sci.* **27**, 135–145 (2018). doi: [10.1002/pro.3290](https://doi.org/10.1002/pro.3290); pmid: [28884485](https://pubmed.ncbi.nlm.nih.gov/28884485/)
90. J. Cheng *et al.*, Accurate proteome-wide missense variant effect prediction with AlphaMissense. *Science* **381**, eadg7492 (2023). doi: [10.1126/science.adg7492](https://doi.org/10.1126/science.adg7492); pmid: [37733863](https://pubmed.ncbi.nlm.nih.gov/37733863/)
91. M. Saitou *et al.*, Functional Specialization of Human Salivary Glands and Origins of Proteins Intrinsic to Human Saliva. *Cell Rep.* **33**, 108402 (2020). doi: [10.1016/j.celrep.2020.108402](https://doi.org/10.1016/j.celrep.2020.108402); pmid: [33207190](https://pubmed.ncbi.nlm.nih.gov/33207190/)
92. W. J. Kent, BLAT—The BLAST-like alignment tool. *Genome Res.* **12**, 656–664 (2002). pmid: [11932250](https://pubmed.ncbi.nlm.nih.gov/11932250/)
93. G. Marçais *et al.*, MUMmer4: A fast and versatile genome alignment system. *PLOS Comput. Biol.* **14**, e1005944 (2018). doi: [10.1371/journal.pcbi.1005944](https://doi.org/10.1371/journal.pcbi.1005944); pmid: [29373581](https://pubmed.ncbi.nlm.nih.gov/29373581/)
94. S. Kurtz *et al.*, Versatile and open software for comparing large genomes. *Genome Biol.* **5**, R12 (2004). doi: [10.1186/gb-2004-5-2-r12](https://doi.org/10.1186/gb-2004-5-2-r12); pmid: [14759262](https://pubmed.ncbi.nlm.nih.gov/14759262/)
95. R. Bouckaert *et al.*, BEAST 2: A software platform for Bayesian evolutionary analysis. *PLOS Comput. Biol.* **10**, e1003537 (2014). doi: [10.1371/journal.pcbi.1003537](https://doi.org/10.1371/journal.pcbi.1003537); pmid: [24272319](https://pubmed.ncbi.nlm.nih.gov/24272319/)
96. G. Renaud, U. Stenzel, J. Kelso, leeHom: Adaptor trimming and merging for Illumina sequencing reads. *Nucleic Acids Res.* **42**, e141 (2014). doi: [10.1093/nar/gku699](https://doi.org/10.1093/nar/gku699); pmid: [25100869](https://pubmed.ncbi.nlm.nih.gov/25100869/)
97. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009). doi: [10.1093/bioinformatics/btp324](https://doi.org/10.1093/bioinformatics/btp324); pmid: [19451168](https://pubmed.ncbi.nlm.nih.gov/19451168/)
98. F. Ramírez *et al.*, deepTools2: A next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.* **44**, W160–W1605 (2016). doi: [10.1093/nar/gkw257](https://doi.org/10.1093/nar/gkw257); pmid: [27079975](https://pubmed.ncbi.nlm.nih.gov/27079975/)
99. B. S. Pedersen, A. R. Quinlan, Mosdepth: Quick coverage calculation for genomes and exomes. *Bioinformatics* **34**, 867–868 (2018). doi: [10.1093/bioinformatics/btx699](https://doi.org/10.1093/bioinformatics/btx699); pmid: [29096012](https://pubmed.ncbi.nlm.nih.gov/29096012/)
100. T. L. Bailey, J. Johnson, C. E. Grant, W. S. Noble, The MEME Suite. *Nucleic Acids Res.* **43**, W39–W49 (2015). doi: [10.1093/nar/gkv416](https://doi.org/10.1093/nar/gkv416); pmid: [25953851](https://pubmed.ncbi.nlm.nih.gov/25953851/)

# ACKNOWLEDGMENTS

We thank S. Syed, V. Reyes-Ortiz, N. Moskwa, P. Hallast, C. Robinett, S. Ruhl, V. Albert, V. Lynch, D. Taylor, L. Speidel, and J. Novembre for technical help, discussions, suggestions, and feedback on the manuscript; J. Akey and L. Li for help analyzing Neanderthal introgression data; the Human Genome Structural Variation Consortium and the Human Pangenome Reference Consortium for making their data publicly available; The Jackson Laboratory Scientific Services, including the Genome Technologies Service, for expert assistance with the work described herein; and Research IT for computational infrastructure and support. The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the Office of the Director of the National Institutes of Health and by NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS. The data used for the analyses described in this manuscript were obtained from the GTEx Portal on 1 December 2023. We are grateful to the people who generously contributed their samples to the 1000 Genomes Project. **Funding:** C.L. and F.Y. were supported by the National Institutes of Health (NIH) U24HG007497; K.K. was supported by The Jackson Laboratory Postdoctoral Scholar Award; O.G. was supported by National Science Foundation (NSF) awards 2049947 and 2123284; C.R.B. was supported by NIH GM133600. **Author contributions:** O.G. and C.L. conceived the study. F.Y. performed the analysis and interpretation of the Human Genome Structural Variation Consortium (HGSVC) and Human Pangenome Reference Consortium (HPRC) Bionano Genomics optical mapping, PacBio HiFi sequencing, and phased assemblies; and performed haplotype-resolved assemblies of HGSVC PacBio HiFi samples using hifiasm and amylase haplotype detection using HGSVC and HPRC datasets. C.K. performed the mutational mechanisms analyses, breakpoint characterization, linkage disequilibrium, principal component analysis, and phylogenetic analyses. C.K. contributed to evolutionary and functional analysis. K.K. and F.Y. performed gene annotation, K.K. performed singleton analysis, selection and phylogenetic analyses of amylase coding sequences and segments, and interpretation and time estimation of initial *AMY1* gene duplication events. F.Y. and K.K. performed archaic



hominin genome processing. P.P. performed ddPCR validation experiments, analysis of functional site differences in amylase amino acid sequences, and data visualization. K.S. calculated amylase gene copy numbers and performed downsampling in ancient human genomes. F.Y., C.K., K.K., P.P., C.R.B., A-M.T., C.L., and O.G. drafted and critically revised the article. Final approval of the version to be published was given by F.Y., C.K., K.K., P.P., K.S., C.R.B., A-M.T., C.L., and O.G. All authors read and approved the final manuscript. **Competing interests:** C.L. is a scientific advisory board member of Nabsys and Genome Insight. **Data and materials availability:** Data files used by this project are available by the following project IDs: Bionano Genomics optical genome mapping: PRJEB41077, PRJEB58376, PRJEB77842, PRJNA315896, PRJNA701308; PacBio HiFi: PRJEB58376, PRJEB75190, PRJNA701308; hifiasm assemblies: PRJEB78558; HPRC phased assemblies: PRJNA701308; ancient human genomes: PRJEB11364, PRJEB11995, PRJEB22660, PRJEB24629, PRJEB26760, PRJEB31045, PRJEB33172, PRJEB36297, PRJEB36854, PRJEB38008, PRJEB41240, PRJEB4604, PRJEB50857, PRJEB56570, PRJEB6272, PRJEB64656, PRJEB6622, PRJEB71770, PRJEB9783, PRJNA229448, PRJNA240906, PRJNA295861, PRJNA299403, PRJNA433631, PRJNA46213, PRJNA670050, PRJNA778930; archaic hominins: PRJEB10597, PRJEB10828, PRJEB1265, PRJEB2065, PRJEB20653, PRJEB21157, PRJEB21195, PRJEB21870, PRJEB21875,

PRJEB21881, PRJEB21882, PRJEB21883, PRJEB24663, PRJEB29475, PRJEB3092, PRJEB31410, PRJEB55327. Amylase segments, haplotypes, orthogonal validation of haplotypes and molecule support, and liftoff gene annotation files (16); amylase segments CDS alignments, dating of *AMY1*, dN/dS analysis, origin of *AMY1* breakpoint, protein2genome gene annotations, tandem repeat analysis files (22); ddPCR protocol (21) and functional analysis files (24); principal component analysis, breakpoint junction characterization, genomic rearrangement analysis, alignments, nucleotide diversity estimation, linkage disequilibrium estimation, phylogenetic trees, PRDM9 binding sites (67); structural variation mechanisms files (65); and ancient genomes downsampling analysis file information (30) are available at Zenodo. **License information:** Copyright © 2024 the authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original US government works. <https://www.sciencemag.org/about/science-licenses-journal-article-reuse>

**Human Genome Structural Variation Consortium**

Peter A. Audano<sup>7</sup>, Olanrewaju Austine-Orimoloye<sup>8</sup>, Christine R. Beck<sup>7,9</sup>, Evan E. Eichler<sup>10</sup>, Pille Hallast<sup>7</sup>, William T. Harvey<sup>10</sup>, Alex R. Hastie<sup>11</sup>, Kendra Hoekzema<sup>10</sup>, Sarah Hunt<sup>8</sup>, Jan O. Korbel<sup>12</sup>, Jennifer Kordosky<sup>10</sup>, Charles Lee<sup>7</sup>, Alexandra P. Lewis<sup>10</sup>, Tobias Marschall<sup>13</sup>, Katherine M. Munson<sup>10</sup>, Andy Pang<sup>11</sup>, Feyza Yilmaz<sup>7</sup>

<sup>7</sup>The Jackson Laboratory for Genomic Medicine, Farmington, CT, USA. <sup>8</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK. <sup>9</sup>The University of Connecticut Health Center, Farmington, CT, USA. <sup>10</sup>Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA. <sup>11</sup>Bionano Genomics, San Diego, CA, USA. <sup>12</sup>European Molecular Biology Laboratory, Genome Biology Unit, Heidelberg, Germany. <sup>13</sup>Institute for Medical Biometry and Bioinformatics, Medical Faculty, Heinrich Heine University, Düsseldorf, Germany.

**SUPPLEMENTARY MATERIALS**

[science.org/doi/10.1126/science.adn0609](https://science.org/doi/10.1126/science.adn0609)

Materials and Methods

Supplementary Text

Figs. S1 to S40

Tables S1 to S15

References (101–131)

MDAR Reproducibility Checklist

Submitted 27 November 2023; resubmitted 27 May 2024

Accepted 24 September 2024

Published online 17 October 2024

10.1126/science.adn0609