

Supporting Information

Fu et al. 10.1073/pnas.1221359110

SI Text

1. SI Shotgun Sequencing

To determine the proportion of endogenous DNA, we carried out low-depth shotgun sequencing from the four libraries prepared from the Tianyuan skeleton. For this purpose, the barcoded libraries prepared from TY1301 (libraries B3071 and B3073) and TY1305 (libraries B3072 and B3074) were combined into two separate pools, each of which was sequenced on a lane of the Illumina Genome Analyzer IIx (FC-104-400x version 4 sequencing chemistry and PE-203-4001 cluster generation kit version 4) using a paired-end run with 76 + 7 cycles (1). An indexed control PhiX 174 library was spiked-in to yield 2–3% control reads (index 5'-TTGCCGC-3'). Base calling was performed with the machine-learning algorithm IBIS (2). Forward and reverse sequence reads overlapping by at least 11 bp were merged into single sequences to reconstruct full-length molecule sequences (3). These were used for further analysis. Merged reads were aligned against the human reference genome [National Center for Biotechnology Information (NCBI) accession no. 37/hg19] using BWA (4) with default parameters, and the output was converted to SAM/BAM format (5). The proportion of endogenous DNA in the four libraries ranges between 0.01% and 0.03% (Table S1), which makes the generation of sequences by whole-genome shotgun sequencing economically unfeasible. We therefore decided to proceed with hybridization enrichment.

2. SI Assembly of Mitochondrial DNA

Sequencing and processing of the raw data were performed as described for shotgun sequencing above. Mapping to the human mtDNA reference genome was done using an iterative mapping assembler (6) with a position-specific scoring matrix that takes into account the nucleotide misincorporation patterns found in ancient DNA sequences. To remove PCR duplicates, we built a consensus from sequences with identical start and end coordinates by retaining the base with the highest sum of quality scores at each position in the alignment. The average length of the mtDNA molecules is 59 bp (Fig. S1A). Mitochondrial coverage as determined from unique sequences ranges between 1.7- and 35.6-fold (Table S1). The consensus sequences obtained from the tibia and the femur are identical. The femur (TY1301) shows the better preservation of the two samples (0.48-fold mtDNA coverage per mg bone; Table S1). One of the libraries prepared from the femur (TY1301), which produced the highest coverage (35.6-fold), was used for further mtDNA analysis. To test whether the mtDNA fragments originated from one individual, the proportion of sequences that matched the consensus base at each position was calculated (Fig. S1B). The average support for the consensus base is 98.8%. Only 36 positions were covered with less than six sequences (the lowest threefold), but all sequences support the consensus base at these positions. The consensus support was below 80% for 8 out of 16,566 positions. However, manual inspection allowed a clear consensus call to be made for these eight cases: Five of the positions were incorrectly aligned, and another three positions showed more than one sequence with a C→T or G→A mismatch close to its end, suggesting that these substitutions represent nucleotide misincorporations due to cytosine deamination.

3. SI DNA Capture

3.1. DNA Library Amplification. To generate large quantities of amplified library for hybridization capture, the libraries prepared

from the femur TY1301 for shotgun sequencing (B3071 and B3073), and a further 39 libraries where deaminated cytosines (uracils) had been enzymatically removed, were reamplified in 24 100- μ L reactions using the primer pair “genomic R1” (5'-ACACTCTTTCCTACACGACGCTCTTCCGATCT-3') and “multiplex R2” (5'-GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT-3').

Amplification products from each library were pooled and purified using solid-phase reversible immobilization (SPRI) technology (7) as follows: An SPRI solution was prepared by combining 95 g PEG 8000 (Promega), 50 mL 5 M NaCl, 2.5 mL 1 M Tris-HCl (pH 8.0), 0.5 mL 0.5 M EDTA (pH 8.0), and 125 μ L Tween 20 and filling up to 250 mL with water. Five milliliters of carboxylated Sera-Mag Speedbeads (Distrilab BV) were washed twice in TE buffer, resuspended in 1 mL TE (10 mM Tris-HCl, 1 mM EDTA, pH 8.0), and added to the SPRI solution to obtain a ready-to-use 38% (wt/vol) PEG-SPRI suspension, which was stored in the refrigerator until used. Pooled PCR products from each library were mixed with an equal volume of SPRI suspension and incubated for 30 min at room temperature. Beads were pulled to the tube wall using a magnet rack and washed twice with 70% ethanol. After drying for 30 min, DNA was eluted in 60 μ L TE buffer. The concentrations of the amplified libraries were between 1,036 and 1,268 ng/ μ L as determined on a NanoDrop 1000 photospectrometer. The original two libraries (B3071 and B3073) and five of the new libraries were captured twice using probes for chromosome 21 and twice using probes for mtDNA. A set of 34 of the new libraries was captured twice using the probes for the admixture SNPs.

3.2. Target Enrichment. 3.2.1. Chromosome 21 capture probe design.

Using the human reference genome sequence (hg19), 11,701,943 probe sequences of 52 bases were extracted from chromosome 21 with 3-bp tiling. To eliminate repetitive sequences, probes containing 15-mer sequences that are overrepresented in the human genome were removed (8). This resulted in the removal of ~25% of the tiled probes, and 37,159 contiguous regions covered by probes remained (Fig. S4). Approximately 35 Mbp of sequence are present in the hg19 chromosome 21 reference assembly. Of these, 29.8 Mb (85%) were targeted by at least one probe, and 22.7 Mb (65%) were covered by 17 or 18 probes. A universal flanking sequence (5'-CACTGCGG-3') was attached to the 3' end of each of the 8,722,911 probe sequences. The resulting 60-base probes were printed on nine custom-designed 1 million-feature arrays (Agilent).

An additional array with 5,506 probes targeting the complete human mitochondrial genome (hg19 as provided by the University of California Santa Cruz Browser, <http://genome.ucsc.edu/cgi-bin/hgGateway>) using 3-bp tiling was designed as above. We removed 38 repetitive probes using the 15-mer filtering. The complete set of 5,468 probes that passed the filtering was printed 178 times in a 1 million-feature array generating a total of 973,304 probes.

3.2.2. “Admixture SNP” capture probe design. To further investigate what proportion of archaic admixture might be present in Tianyuan, we identified sites where seven individuals from seven different African populations (Bantu Kenya, Bantu South Africa, Biaka, Mandenka, Mbuti, San, and Yoruba) from the CEPH-Human Genome Diversity Panel (HGDP-CEPH) differ from both the Denisovan and Neandertal genomes (set A) and sites where all these African individuals match the Neandertal genome but are all different from the Denisovan (set B) (9). We designed a capture array for the 1,666 sites in set A and the 1,800 sites in set B.

Using the human reference genome sequence (hg19), we designed 52-bp-long probes tiled at 3-bp intervals across 105 bp

centered on each SNP. For each probe, alternatives carrying each of the two allelic variants of the SNP were included. To eliminate repetitive sequences, 49% of the probes containing 15-mer sequences that are overrepresented in the human genome were removed (8). The sequence 5'-CACTGCGG-3' was attached to the 3' end of each of the 124,163 probes. The resulting 60-base probes were printed on nine custom-designed 1 million-feature arrays (Agilent).

3.2.3. Generation of biotinylated capture probes. The array probes were cleaved and converted into probe libraries as follows. After adding 500 μ L elution solution (125 mM NaOH, 0.05% Tween 20), each array was assembled with a gasket slide in a hybridization chamber (8) and rotated for 6 h with 12 rpm at room temperature. The eluate was recovered using a syringe, neutralized by adding 19 μ L 20% acetic acid, and purified using the QIAquick Nucleotide Removal Kit (Qiagen). The purified probes were eluted in 20 μ L EB (10 mM Tris-Cl, pH 8.5). Using 2 μ L of the eluate, successful probe recovery was confirmed by denaturing PAGE. The first probe library adapter was added through a primer extension reaction with Bst polymerase. The 50- μ L reaction mixture contained 10 μ L of purified probes and 3 μ L of (24 U) Bst polymerase (large fragment; New England BioLabs), and in final concentrations 1 \times Thermopool buffer (New England BioLabs), 250 μ M each dNTP, and 1 μ M extension primer APL2 (5'-biotin-CGTGGATGAGGAGCCGAGTG-3'). After incubation for 1 min at 50 $^{\circ}$ C, 5 min at 15 $^{\circ}$ C, 5 min at 20 $^{\circ}$ C, 5 min at 25 $^{\circ}$ C, 5 min at 30 $^{\circ}$ C, and 10 min at 37 $^{\circ}$ C in a thermal cycler, the reaction was purified using the MinElute PCR Purification Kit (Qiagen). Subsequently, a blunt-end repair reaction was performed to remove 3' overhangs generated by Bst polymerase. The 40- μ L reaction mixture contained the complete eluate (20 μ L), 0.4 μ L (2 U) T4 DNA polymerase (Fermentas), and 0.4 μ L (4 U) Klenow fragment (Fermentas), and in final concentrations 1 \times Tango buffer (Fermentas) and 100 μ M each dNTP. After incubation for 15 min at 25 $^{\circ}$ C, the reaction was purified using the MinElute PCR Purification Kit, eluting in 20 μ L EB. The second probe library adapter was added by blunt-end ligation as follows. A double-stranded adapter was generated by combining 7 μ L water, 1 μ L T4 DNA ligase buffer (Fermentas), 1 μ L 100 μ M APL1 (5'-phosphate-ACACGCTGGTGCATCCCTAT-Pho-3'), and 1 μ L 100 μ M APL6 (5'-ATAGGGATCGC-ACCAGCGTGT-3'). The mixture was incubated for 10 s at 95 $^{\circ}$ C in a thermal cycler and slowly cooled to 14 $^{\circ}$ C at a rate of 0.1 $^{\circ}$ /s. Then, 3 μ L T4 DNA ligase buffer, 2 μ L water, 4 μ L 50% PEG 4000, and the eluate from the previous reaction were added. After mixing, 1 μ L (5 U) T4 DNA ligase (Fermentas) was added and the reaction was incubated for 30 min at room temperature and purified using the MinElute PCR Purification Kit. To remove adapter dimers, 20 μ L MyOne C1 streptavidin beads (Invitrogen) were washed twice with BWT+SDS buffer (1 M NaCl, 10 mM Tris-HCl, pH 8.0, 1 mM EDTA, 0.05% Tween 20, 0.5% SDS). The beads were resuspended in 180 μ L BWT+SDS buffer, the complete eluate from the ligation reaction was added (20 μ L), and the suspension was rotated at room temperature for 20 min. The beads were washed twice with 0.1 \times BWT buffer (0.1 M NaCl, 10 mM Tris-HCl, pH 8.0, 1 mM EDTA), resuspended in 25 μ L TT buffer (1 mM Tris-HCl, pH 8.0, 0.01% Tween 20), and incubated for 10 min at 95 $^{\circ}$ C to release the biotinylated strands, representing the final probe library. Using 1 mL of the probe library and a standard dilution series of known concentration, the number of probe library molecules was determined by quantitative (q)PCR using the primer pair APL5-APL6. Based on this assay, we estimated that on average 16,000 copies of each probe were recovered.

Probe libraries were amplified for nine cycles in 100- μ L PCR reactions (avoiding PCR plateau). Each reaction contained 10 μ L probe library (~6,000 copies per probe) and 1 μ L Herculase II Fusion DNA polymerase (Agilent), and in final concentrations

1 \times Herculase II reaction buffer, 250 μ M each dNTP, and 400 nM the primers APL5 (5'-CGTGGATGAGGAGCCGAGTG-3') and APL6. An initial denaturation step of 2 min at 95 $^{\circ}$ C was followed by nine cycles of denaturation at 95 $^{\circ}$ C for 20 s, annealing at 60 $^{\circ}$ C for 30 s and elongation at 72 $^{\circ}$ C for 30 s, and a final extension step at 72 $^{\circ}$ C for 5 min. PCR products were purified using the MinElute PCR Purification Kit and eluted in 20 μ L EB. Seven microliters of each product was then loaded on a 3% low-melting/1% high-melting agarose gel with SYBR Safe (Invitrogen). Narrow bands around 94 bp were excised from the gel to remove a faint smear of below-full-length probes. DNA was isolated from the gel slices using the MinElute Gel Extraction Kit, eluting in 30 μ L EB. One microliter of each eluate was used for qPCR to verify the success of gel extraction and determine an optimal cycle number for the subsequent amplification (avoiding PCR plateau). Between 8 and 19 μ L of gel-excised probes was used as template for amplification reactions in 100- μ L volumes, containing 1 μ L Herculase II Fusion DNA polymerase and in final concentrations 1 \times Herculase II reaction buffer, 250 μ M each dNTP, and 400 nM the primers APL5 and APL4 (5'-GGATTCTAATACGACTCACTATAGGGATCGCACCAGC-GTGT-3'). An initial denaturation step of 2 min at 95 $^{\circ}$ C was followed by eight cycles of denaturation at 95 $^{\circ}$ C for 20 s, annealing at 60 $^{\circ}$ C for 30 s and elongation at 72 $^{\circ}$ C for 30 s, and a final extension step at 72 $^{\circ}$ C for 5 min. The PCR products were purified using the MinElute PCR Purification Kit, eluted in 40 μ L EB, and quantified using a NanoDrop photospectrometer. Concentrations of the amplified probe libraries varied between 69 and 97 ng/ μ L. The nine probe libraries for chromosome 21 were pooled in equimolar ratio. At this stage, the probe library contains a T7 promoter sequence (introduced by APL4), in principle allowing for the generation of RNA capture probes following Gnirke et al. (10). However, because we did not observe substantial differences in the performance of DNA and RNA probe capture in earlier experiments, we chose to perform all captures of the Tianyuan libraries with DNA probes.

Single-stranded biotinylated DNA probes were generated in single-primed linear amplification reactions using a biotinylated primer. Because high amounts of template were required for these reactions, the chromosome 21 and mtDNA probe libraries were further amplified using the primer pair APL2 and APL6 under the conditions described above. For each probe set, 96 100- μ L reactions were prepared, containing 200 μ L template and 2 μ L Herculase II Fusion DNA polymerase, and in final concentrations 1 \times Herculase II reaction buffer, 250 μ M each dNTP, and 400 nM APL2. An initial denaturation step of 2 min at 95 $^{\circ}$ C was followed by 20 cycles of denaturation at 95 $^{\circ}$ C for 20 s, and annealing at 60 $^{\circ}$ C for 20 s and elongation at 72 $^{\circ}$ C for 20 s. All reactions were pooled in a Falcon tube and mixed with a double volume of 38% PEG-SPRI suspension (~20 mL). All other steps of the SPRI purification were performed as described above. Probes were eluted in 150 μ L TE buffer and quantified using a NanoDrop photospectrometer (~250 ng/ μ L). DNA probes were stored at -20 $^{\circ}$ C until used.

3.2.4. Hybridization capture. For each hybridization reaction, a sample library pool (15 μ L total volume) was created by combining 6 μ L sample library (~2 μ g), 5.25 μ L water, 2.5 μ L 1 mg/mL human Cot-1 DNA (Invitrogen), 0.25 μ L 10 mg/mL salmon sperm DNA (Invitrogen), 0.5 μ L 500 μ M BO4 (5'-GTGACTGGAGTTCA-GACGTGTGCTCTCCGATCT-phosphate-3'), and 0.5 μ L 500 μ M BO10 (5'-AGATCGGAAGAGCGTCTGTAGGGAAAGAG-TGT-phosphate-3'). The sample library pool was incubated for 5 min at 95 $^{\circ}$ C and then 5 min at 65 $^{\circ}$ C, and held at room temperature afterward. A probe pool was created by diluting 300 ng single-stranded DNA probes with water to obtain a total volume of 4 μ L. Hybridization buffer was prepared by combining 1.7 mL Hi-RPM buffer (aCGH Hybridization Kit; Agilent) and 300 μ L 50 \times Denhardt's solution (Sigma-Aldrich). Hybridization re-

actions were assembled in 96-well plates by adding 20 μL hybridization buffer and the complete sample library pool to the probe pool. The reactions were mixed and incubated at 62 °C for 2 d. For each reaction, 30 μL MyOne T1 streptavidin beads (Invitrogen) were washed once with 150 μL wash buffer 1 (1 \times SSC, 0.01% SDS) for 15 min at room temperature, three times with 120 μL HWT buffer [1 \times AmpliTaq Gold buffer without MgCl_2 (Applied Biosystems), 0.02% Tween 20] for 10 min at 60 °C, and once with 150 μL wash buffer 3 (0.1 \times SSC, 0.05% Tween-20) for 5 min at room temperature. Before each wash step, liquid was collected in the bottom of the wells by briefly spinning the plate at 2,000 $\times g$ in a centrifuge. The plate was then placed on a 96-well ring magnet plate and the supernatant was removed. Wash buffer was added and the plate was sealed with strip caps. Beads were fully resuspended by vortexing for 5–10 s. During room-temperature wash steps, the plate was taped to a rotator. For high-temperature wash steps, the plate was placed in a thermal cycler (with lid heating turned off) and inverted several times during incubation. After the last wash step, to eluate the capture library molecules, the beads were resuspended in 19 μL melt solution (125 mM NaOH, 0.05% Tween 20). After incubation for 15 min at room temperature, the supernatant was transferred to a fresh plate and mixed with 0.7 μL 20% acetic acid and 190 μL PN buffer (Qiagen). Each capture eluate was then purified in a separate MinElute spin column following Qiagen's instructions for using the Nucleotide Removal Kit. DNA was eluted in 30 μL TT buffer.

To verify the successful retrieval of library molecules, molecule numbers were estimated from 1 μL of capture eluate by qPCR (~1E8 total molecules for most libraries). The remaining 29 μL was amplified in 100- μL reactions, containing 1 μL Herculanase II Fusion DNA polymerase, and in final concentrations 1 \times Herculanase II reaction buffer, 250 μM each dNTP, and 400 nM the primers genomic R1 and multiplex R2. An initial denaturation step of 2 min at 95 °C was followed by 28 cycles of denaturation at 95 °C for 30 s, annealing at 60 °C for 30 s and elongation at 72 °C for 30 s, and a final extension step at 72 °C for 5 min. Amplified libraries were purified using a 1:1 ratio of 38% PEG-SPRI suspension as described above. The amplified capture eluates were eluted in 15 μL TE and their concentrations were determined using a NanoDrop photospectrometer (~170 ng/ μL on average).

Three microliters (~500 ng) of capture eluate from the first round of hybridization was used as template for a second round of hybridization, which was performed under the same conditions except for a reduced incubation time (1 d). Capture eluates were again quantified by qPCR (values were one to two orders of magnitude higher), amplified (23 cycles), purified with SPRI beads, and quantified on a NanoDrop photospectrometer.

3.2.5. Library pooling and multiplex sequencing. The amplified capture eluates were diluted 25-fold using TE buffer (to obtain concentrations of ~5 ng/ μL). To enable highly accurate multiplex sequencing (1), indexes were added to both library adapters using 5'-tailed PCR primers in 50- μL reactions containing 1 μL template and 0.5 μL Herculanase II Fusion DNA polymerase, and in final concentrations 1 \times Herculanase II reaction buffer, 250 μM each dNTP, and 400 nM each indexing primer. Amplification was performed with a small cycle number to avoid the formation of heteroduplicates in PCR plateau. An initial denaturation step of 2 min at 95 °C was followed by six cycles of denaturation at 95 °C for 30 s, annealing at 60 °C for 30 s and elongation at 72 °C for 30 s, and a final extension step at 72 °C for 5 min. Fifty microliters of PB buffer (Qiagen) and 1 μL 3 M sodium acetate were added to each reaction. By pooling the sample/PB mixes in equimolar ratios, several library pools were generated. Libraries from chromosome 21 and mtDNA capture were kept in separate pools. After adding the twofold volume of PB buffer to each pool, the libraries were purified using the MinElute PCR Purification Kit

and eluted in 25 μL TE. DNA concentration was determined using a DNA 1000 chip on a Bioanalyzer 2100 (Agilent).

Sequencing and raw sequence processing were performed as described in SI Text, section 1. We generated one lane of sequence data from the chromosome 21 captures. mtDNA captures were sequenced on one-quarter and the SNP capture was sequenced on one-fifth of an Illumina Genome Analyzer Ix lane. Alignments were generated by mapping the merged reads against the human reference genome (NCBI accession no. 37/hg19) using BWA (4) with default parameters.

3.3. Evaluating the Efficiency of the Target Enrichment Method. To evaluate whether the enriched libraries had been sequenced to exhaustion, we performed a subsampling analysis by randomly drawing subsets of aligned sequences and counting the number of unique sequences in each subset based on identical alignment start and end coordinates. Because capture had been performed in replicates for each library, we performed this analysis separately for the sequences from each replicate as well as the combined sequences from both replicates (Fig. S5). The subsampling curves provide two insights. First, sequencing depth was sufficient to almost completely exhaust the complexity of the enriched libraries, that is, deeper sequencing would not considerably increase the number of sequences of unique DNA fragments. Second, the trajectories of the subsampling plots are virtually identical, irrespective of whether sequences from capture replicates are analyzed separately or in combination. This indicates that the number of sequences obtained from unique DNA fragments is merely dependent on sequencing depth and that replicate captures from the same libraries are not required.

We next wanted to assess what proportion of the target molecules existing in the libraries we successfully captured and sequenced via hybridization enrichment. For this purpose, we made use of the whole-genome shotgun sequences that had been generated from some of the libraries (B3071 and B3073) that were also used for capture. We first searched for sequences that aligned to the capture target regions of chromosome 21 in the whole-genome shotgun data and identified 52 such sequences. Based on identical start and end alignment coordinates, we next checked how many of these sequences were also present in the chromosome 21 capture data. For B3071, we found that 18 out of 23 (78%) DNA fragments with a length of ≥ 35 bp and a mapping quality of ≥ 30 are also represented in the capture dataset. For B3073, we determined a similar number (74%; 14 out of 19 sequences). These results demonstrate that the capture strategy presented in this study efficiently captures most of the target molecules present in a library.

3.4. Chromosome 21 Sequence Coverage and mtDNA Contamination Estimate. In total, we obtained 9,373,365 sequences from five uracil-DNA-glycosylase (UDG)-treated libraries with a length of at least 35 bp, 4,406,261 of which (46.8%) aligned to chromosome 21 with a mapping quality of at least 30. For each library, sequences that map to the same outer reference coordinates were replaced by the sequence with the highest sum of base qualities (11). Through this approach, 789,925 sequences were found to be unique, and each unique molecule was sequenced on average 5.6 times. When combining sequences from all five libraries, 19.9 Mbp of 29.8 Mbp of the Tianyuan chromosome 21 target regions were captured, and average coverage is 1.75-fold.

To estimate human mitochondrial contamination in each of the libraries used for nuclear DNA capture, we focused on three positions where the Tianyuan mitochondrial consensus sequence differs from at least 99% of 311 present-day human mitochondrial genomes (12) (position 5,348: C \rightarrow T; position 5,836: A \rightarrow G; position 11,257: C \rightarrow T) (Fig. 3B). The positions 5,348: C \rightarrow T and 11,257: C \rightarrow T may appear to be due to deamination; however, this seems unlikely, because (i) the libraries are UDG-treated, (ii) the

majority of reads agree on the nucleotide, and (iii) these positions are not near the ends of the majority of reads, as is typical for deaminations. We then counted the unique fragments that cover these positions to determine how many differ from the consensus sequence, that is, are putative contaminants. We estimate human mitochondrial contamination across all libraries to be on average 1.1% (highest, 3.1%; lowest, 0.1%).

3.5. Chromosome 21 Sequence Determination. We disregarded sites where all 13 individuals have the same genotype; 171,853 sites are variable in at least one of among the 13 individuals. To ensure that genotyping errors do not dominate the differences between individuals, we required the difference in Phred-scaled likelihoods between the two most likely genotypes to be at least 50 (corresponding to an error rate of no more than 10^{-5}). At sites where we could not call a diploid genotype in this way, we considered whether the likelihoods of the two most likely homozygous genotypes differ by at least 50, and if so called the most likely haploid genotype instead. A total of 87,243 sites pass this filter.

To also use the information from the sites called as haploid genotypes in subsequent analyses, haploid calls were counted as zero or one reference allele out of one observation, whereas the diploid calls were counted as zero, one, or two reference alleles out of two observations.

To compute pairwise distances, we ignored sites where two or more alternative alleles were seen among the individuals analyzed. A total of 86,525 sites passed this filter. The distance between two genotypes was defined as the difference in the number of reference alleles (Table S2). To calculate the distance between two individuals, the distances between genotypes were summed over all sites.

3.6. Phylogenetic Reconstruction of Chromosome 21. TreeMix estimates a tree (Fig. S64) where there is 100% bootstrap support for Tianyuan clustering together with Karitiana, Han, and Dai (100 bootstrap replicates), irrespective of whether an admixture event between the Denisovan and Papuan populations is taken into account (Fig. 2 and Fig. S64).

3.7. Archaic Admixture. We note that the distance between the Denisovan and the Tianyuan chromosome 21 sequences (43,893) is similar to the distance between the Denisovan and the Papuan sequences (43,935) and smaller than the distances between the Denisovan and the other Asian sequences (45,160–47,535) (Table 1). We therefore explored whether a population related to the Denisovan individual may have contributed genes to the ancestors of the Tianyuan individual, as is the case with present-day Melanesians (13, 14) and has been suggested for some

mainland Asian populations (15), although this is more likely to represent Neandertal gene flow (16). Although the residuals (Fig. S6B) from the TreeMix analysis (which can be interpreted as signals of admixture) detect the previously described admixture signal between Denisovans and Papuans (Fig. 2), there is no indication of admixture between Denisovans and the Tianyuan individual that exceeds that seen between Denisovans and any other population analyzed.

To further investigate what proportion of archaic admixture might be present in Tianyuan, we identified sites where seven individuals from seven different HGDP-CEPH Africans (Bantu Kenya, Bantu South Africa, Biaka, Mandenka, Mbuti, San, and Yoruba) differ from both the Denisovan and Neandertal genomes (set A), and sites where all of these African individuals match the Neandertal genome but differ from the Denisovan genome (set B). The sharing of alleles seen in both archaic genomes (set A) tends to detect archaic human admixture both from Neandertals and Denisovans or groups related to them, whereas sharing of alleles seen in the Denisovan but not the Neandertal (set B) tends to show admixture with Denisovans or groups related to them.

We randomly selected one individual from each of the 44 populations in mainland Eurasia and the Americas represented in the HGDP-CEPH and counted the number of alleles each individual shares with the archaic individuals in the two sets of SNPs (Fig. 3). We similarly analyzed the 11 sequenced present-day humans and the Tianyuan. Here 1,499 sites in set A and 1,618 sites in set B were scored in all individuals analyzed (Fig. S7). We find no evidence for Denisovan gene flow into the population from which the Tianyuan individual is derived, above what might be present in all mainland Asian populations.

We note that there are at least two potential explanations for the relatively small distance between the Tianyuan individual and the Denisovan individual other than Denisovan gene flow. First, Denisovan and Tianyuan might share some artifacts present in ancient DNA sequences. However, whereas errors are present in the Tianyuan sequences due to its low coverage, the Denisovan genome sequence is of ~ 30 -fold coverage and the DNA fragments are sequenced to such a high redundancy that its error rate is lower than the present-day sequences analyzed here (16). Second, the smaller distance between the Tianyuan and Denisovan individuals may be due to the fact that the 40,000-year-old Tianyuan individual as well as the perhaps equally old (or older) Denisovan individual are closer to their common ancestral population than present-day populations. This explanation is compatible with the observation that the numbers of substitutions inferred to have occurred between the Denisovan genome sequence and the common ancestor of humans and chimpanzees is smaller than for present-day humans (15, 16).

- Kircher M, Sawyer S, Meyer M (2012) Double indexing overcomes inaccuracies in multiplex sequencing on the Illumina platform. *Nucleic Acids Res* 40(1):e3.
- Kircher M, Stenzel U, Kelso J (2009) Improved base calling for the Illumina Genome Analyzer using machine learning strategies. *Genome Biol* 10(8):R83.
- Kircher M, Heyn P, Kelso J (2011) Addressing challenges in the production and analysis of Illumina sequencing data. *BMC Genomics* 12:382.
- Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25(14):1754–1760.
- Li H, et al.; 1000 Genome Project Data Processing Subgroup (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25(16):2078–2079.
- Briggs AW, et al. (2009) Targeted retrieval and analysis of five Neandertal mtDNA genomes. *Science* 325(5938):318–321.
- DeAngelis MM, Wang DG, Hawkins TL (1995) Solid-phase reversible immobilization for the isolation of PCR products. *Nucleic Acids Res* 23(22):4742–4743.
- Hodges E, et al. (2007) Genome-wide in situ exon capture for selective resequencing. *Nat Genet* 39(12):1522–1527.
- Patterson N, et al. (2012) Ancient admixture in human history. *Genetics* 192(3):1065–1093.
- Gnirke A, et al. (2009) Solution hybrid selection with ultra-long oligonucleotides for massively parallel targeted sequencing. *Nat Biotechnol* 27(2):182–189.
- Kircher M (2012) Analysis of high-throughput ancient DNA sequencing data. *Methods Mol Biol* 840:197–228.
- Green RE, et al. (2008) A complete Neandertal mitochondrial genome sequence determined by high-throughput sequencing. *Cell* 134(3):416–426.
- Reich D, et al. (2010) Genetic history of an archaic hominin group from Denisova Cave in Siberia. *Nature* 468(7327):1053–1060.
- Reich D, et al. (2011) Denisova admixture and the first modern human dispersals into Southeast Asia and Oceania. *Am J Hum Genet* 89(4):516–528.
- Skoglund P, Jakobsson M (2011) Archaic human ancestry in East Asia. *Proc Natl Acad Sci USA* 108(45):18301–18306.
- Meyer M, et al. (2012) A high-coverage genome sequence from an archaic Denisovan individual. *Science* 338(6104):222–226.

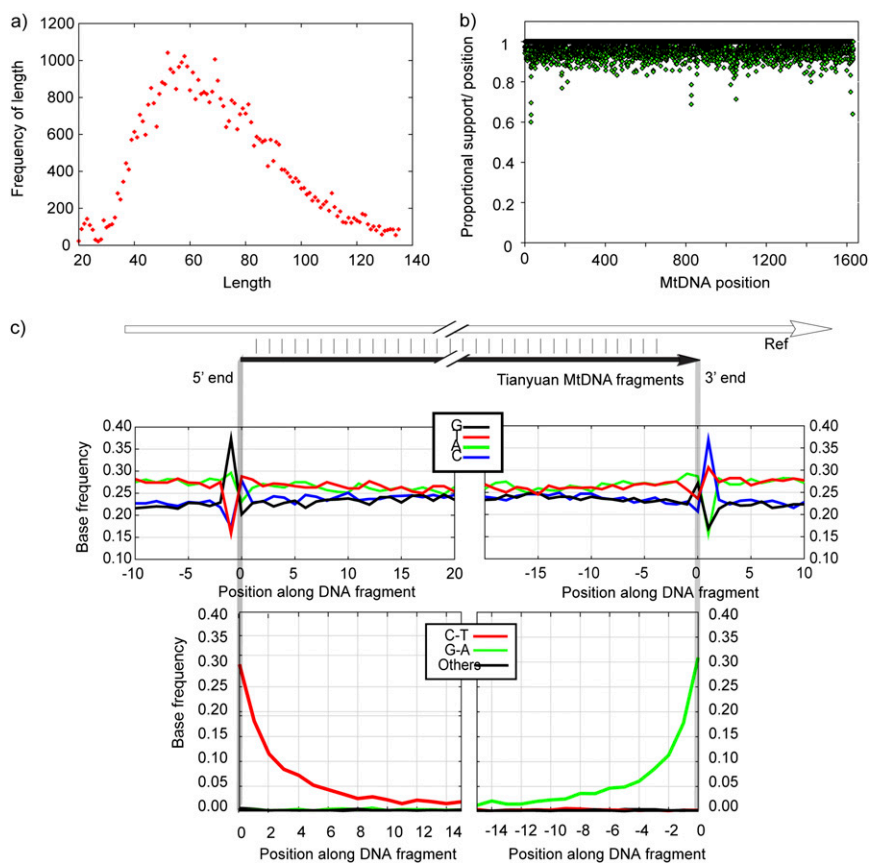


Fig. S1. (A) Length distribution of Tianyuan mtDNA fragments. (B) Fraction of bases agreeing with the consensus sequence along the mitochondrial genome. (C) DNA fragmentation and nucleotide misincorporation patterns inferred from Tianyuan mtDNA fragments. The base composition of the reference genome around fragment ends (Upper) indicates preferential strand breakage 3' of the purines. Increased frequency of C→T and G→A substitutions (Lower) close to fragment ends indicates cytosine deamination. The frequency X→Y is given as the fraction of bases at the indicated positions along the DNA fragments and inferred to be X according to the consensus sequence, which is observed as Y in the individual DNA fragments.

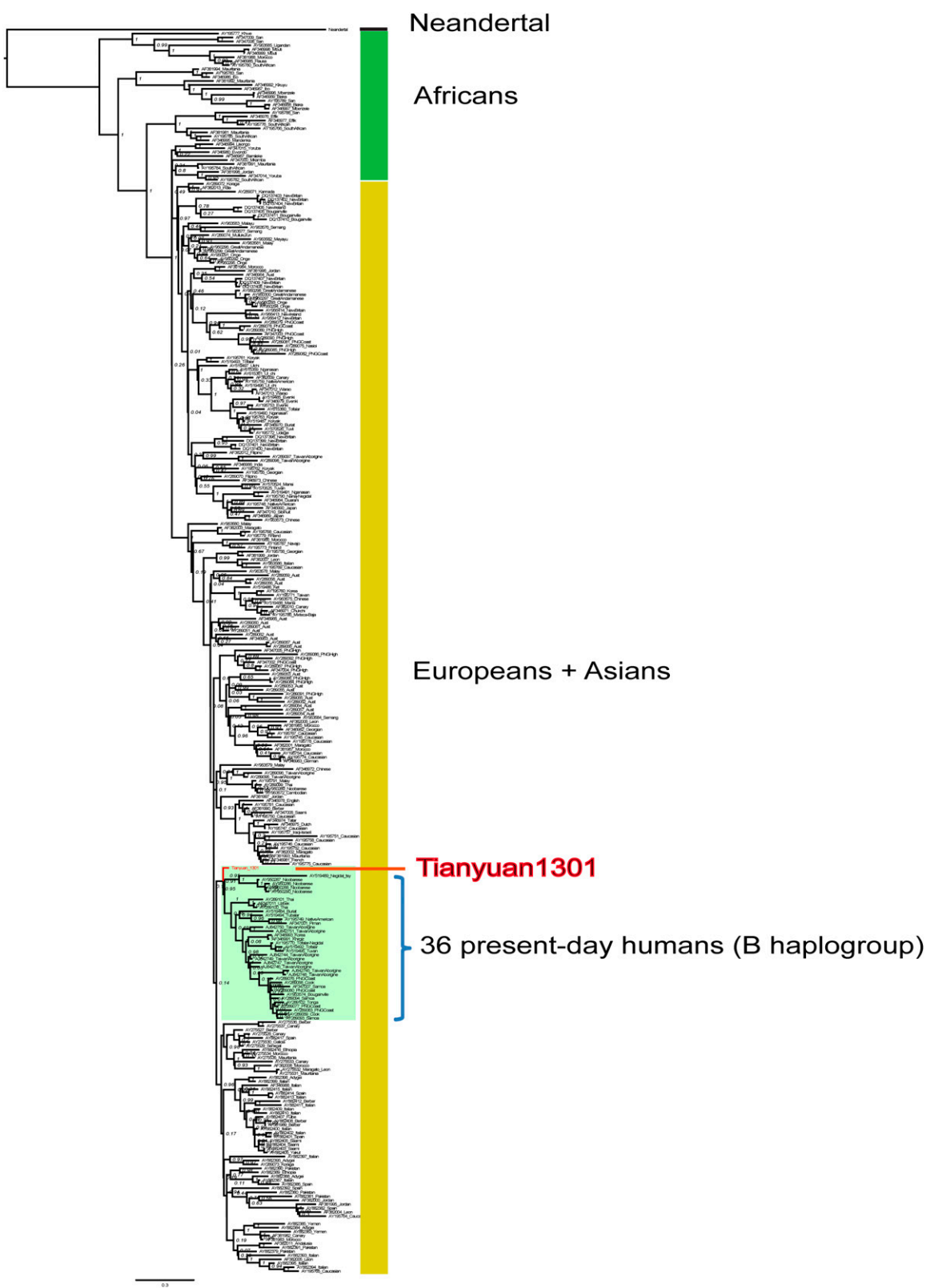


Fig. S2. Phylogenetic tree of mtDNAs from Tianyuan and 311 present-day humans. The phylogeny was estimated with a Bayesian approach under a GTR+I model of sequence evolution, using a Neandertal mtDNA (Vindija 33.25) (1) as an outgroup. Bar colors indicate different populations. The clade of Tianyuan and its closest 36 present-day humans is in green.

1. Briggs AW, et al. (2009) Targeted retrieval and analysis of five Neandertal mtDNA genomes. *Science* 325(5938):318–321.

a)				b)			
	8272	8282	16189	5348	5836	11257	
rCRS	CCTATAGCA	CCCCCTCTA	CCCCCTC	rCRS	C	rCRS	C
Tianyuan	CCTATAGCA	CCCCCTCTA	T	Tianyuan	T	Tianyuan	T
AY950286	CCTATAGCA	CCCCCTCTA	C	AY950286	C	AY950286	C
AY950287	CCTATAGCA	CCCCCTCTA	C	AY950287	C	AY950287	C
AY950288	CCTATAGCA	CCCCCTCTA	C	AY950288	C	AY950288	C
AY950290	CCTATAGCA	CCCCCTCTA	C	AY950290	C	AY950290	C
AY519489	CCTATAGCA	CCCCCTCTA	C	AY519489	C	AY519489	C
AJ842744	CCTATAGCA	CCCCCTCTA	C	AJ842744	C	AJ842744	C
AJ842746	CCTATAGCA	CCCCCTCTA	C	AJ842746	C	AJ842746	C
AJ842747	CCTATAGCA	CCCCCTCTA	C	AJ842747	C	AJ842747	C
AJ842749	CCTATAGCA	CCCCCTCTA	C	AJ842749	C	AJ842749	C
AY289076	CCTATAGCA	CCCCCTCTA	C	AY289076	C	AY289076	C
AY289083	CCTATAGCA	CCCCCTCTA	C	AY289083	C	AY289083	C
AF347007	CCTATAGCA	CCCCCTCTA	C	AF347007	C	AF347007	C
AY289068	CCTATAGCA	CCCCCTCTA	C	AY289068	C	AY289068	C
AY289077	CCTATAGCA	CCCCCTCTA	C	AY289077	C	AY289077	C
AY289080	CCTATAGCA	CCCCCTCTA	C	AY289080	C	AY289080	C
AY289094	CCTATAGCA	CCCCCTCTA	C	AY289094	C	AY289094	C
AY289102	CCTATAGCA	CCCCCTCTA	C	AY289102	C	AY289102	C
AY963574	CCTATAGCA	CCCCCTCTA	C	AY963574	C	AY963574	C
AY289069	CCTATAGCA	CCCCCTCTA	C	AY289069	C	AY289069	C
AY289093	CCTATAGCA	CCCCCTCTA	C	AY289093	C	AY289093	C
AJ842745	CCTATAGCA	CCCCCTCTA	C	AJ842745	C	AJ842745	C
AJ842748	CCTATAGCA	CCCCCTCTA	C	AJ842748	C	AJ842748	C
AF346991	CCTATAGCA	CCCCCTCTA	C	AF346991	C	AF346991	C
AF346993	CCTATAGCA	CCCCCTCTA	C	AF346993	C	AF346993	C
AY195770	CCTATAGCA	CCCCCTCTA	C	AY195770	C	AY195770	C
AY519492	CCTATAGCA	CCCCCTCTA	C	AY519492	C	AY519492	C
AY519495	CCTATAGCA	CCCCCTCTA	C	AY519495	C	AY519495	C
AJ842750	CCTATAGCA	CCCCCTCTA	C	AJ842750	C	AJ842750	C
AJ842751	CCTATAGCA	CCCCCTCTA	C	AJ842751	C	AJ842751	C
AF347001	CCTATAGCA	CCCCCTCTA	C	AF347001	C	AF347001	C
AY195749	CCTATAGCA	CCCCCTCTA	C	AY195749	C	AY195749	C
AY519494	CCTATAGCA	CCCCCTCTA	C	AY519494	C	AY519494	C
AY519484	CCTATAGCA	CCCCCTCTA	C	AY519484	C	AY519484	C
AY289101	CCTATAGCA	CCCCCTCTA	C	AY289101	C	AY289101	C
AF347011	CCTATAGCA	CCCCCTCTA	C	AF347011	C	AF347011	C
AY289100	CCTATAGCA	CCCCCTCTA	C	AY289100	C	AY289100	C

Fig. S3. (A) A 9-bp deletion and a substitution (at position 16,189) present in Tianyuan mtDNA and present-day human mtDNAs belonging to haplogroup B. (B) Three positions where at least 308 of 311 present-day human mtDNAs agree and the Tianyuan mtDNA is different.

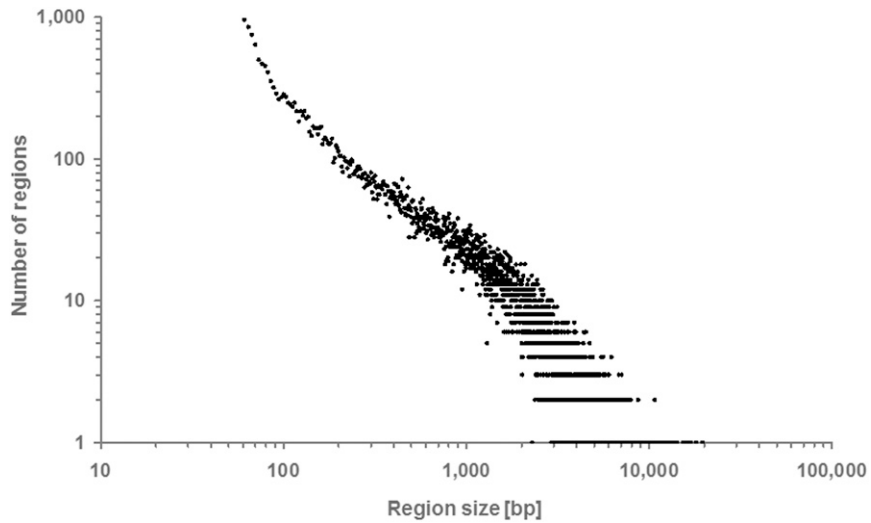


Fig. S4. Size distribution of chromosome 21 target regions.

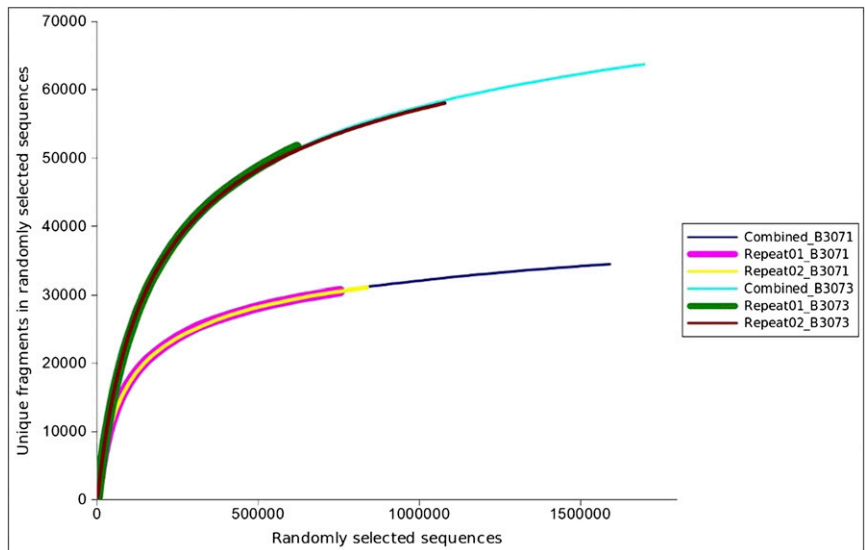


Fig. S5. Subsampling plots obtained from sequences from replicate captures of chromosome 21. Plots are shown for two of the libraries only (B3071 and B3073), but the results are consistent with those seen for all libraries.

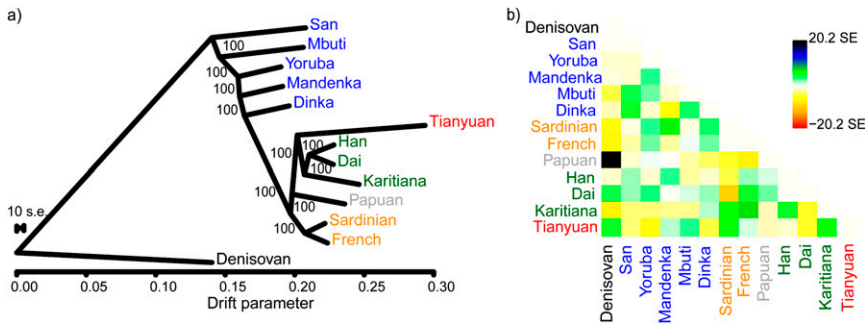


Fig. S6. (A) Maximum-likelihood tree of chromosome 21 sequences of the Tianyuan individual, 11 present-day humans, and the Denisovan individual. (B) Residual matrix of chromosome 21 sequences of the Tianyuan individual, 11 present-day humans, and the Denisovan individual. Residuals above zero represent populations that are closely related to each other according to the allele covariance matrix. SE indicates the average SE across each pair of populations. Darker colors represent stronger signal.

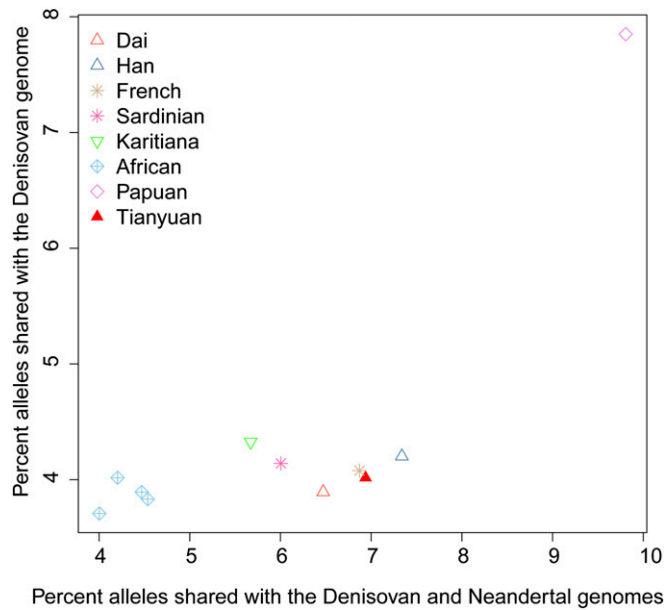


Fig. S7. Proportions of alleles shared with Neandertal and Denisovan genomes in the Tianyuan and 11 present-day individuals. The x axis shows the percent of alleles that match the Neandertal and Denisovan genomes at sites where both of these differ from the seven Africans, and the y axis indicates the percent of alleles that match the Denisovan genome where this differs from the Neandertal as well as the seven Africans.

Table S1. Shotgun and mtDNA sequencing information for TY1301 and TY1305

Bone ID	Library ID	Powder (mg)	Unique fragments aligned to hg19	Endogenous (%)	Average length (bp)	Unique fragments aligned to mtDNA	Coverage of mtDNA	Coverage of mtDNA/mg
TY1301-1	B3071	37	2,529	0.02	64	4,949	20	0.54
TY1301-2	B3073	84	2,993	0.03	52	9,883	35.6	0.42
TY1305-1	B3072	61	2,263	0.01	59	653	2.4	0.04
TY1305-2	B3074	93	725	0.01	42	625	1.7	0.02

Table S2. Distance matrix between all possible pairs of genotypes

	AA	AB	BB	A	B
AA	0	1	2	0	2
AB	1	0	1	1	1
BB	2	1	0	2	0
A	0	1	2	0	2
B	2	1	0	2	0

A, reference allele; B, alternative allele; double letters indicate diploid calls; single letters indicate haploid calls.