

Lecture 18. Ancient and Modern Humans

Michael Schatz

April 13, 2017

JHU 600.649: Applied Comparative Genomics





Part I:

Clustering Refresher

Clustering Refresher

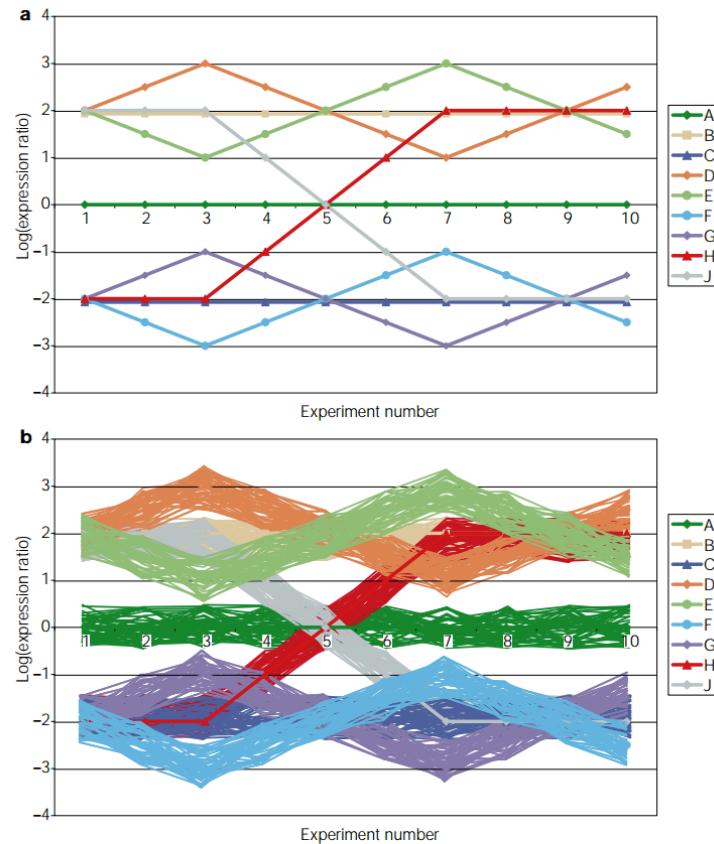
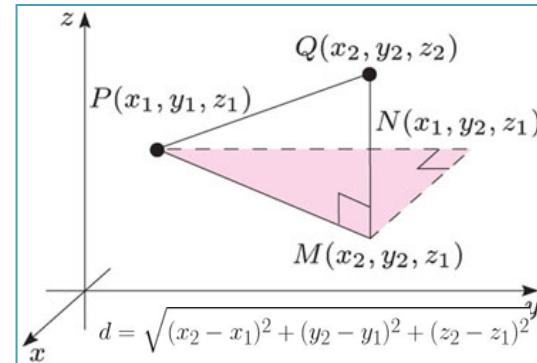
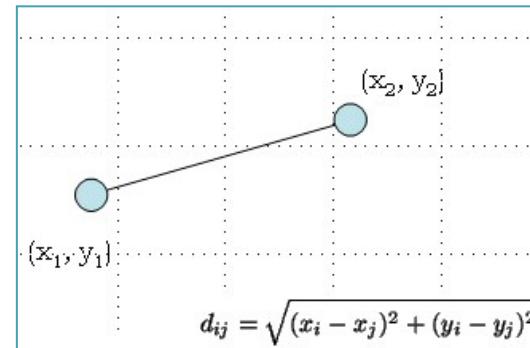


Figure 2 | A synthetic gene-expression data set. This data set provides an opportunity to evaluate how various clustering algorithms reveal different features of the data. **a** | Nine distinct gene-expression patterns were created with $\log_2(\text{ratio})$ expression measures defined for ten experiments. **b** | For each expression pattern, 50 additional genes were generated, representing variations on the basic patterns.

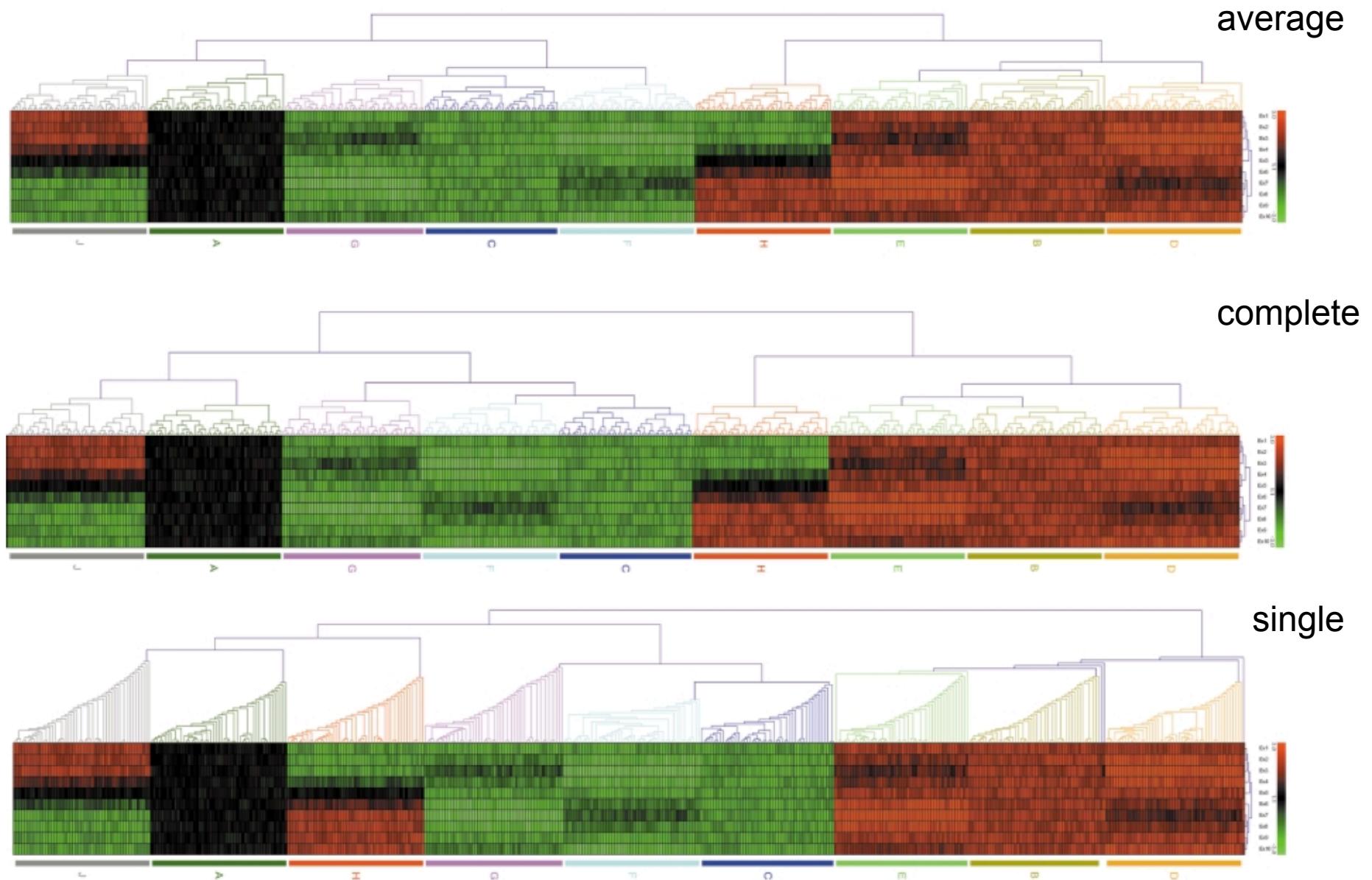
Euclidean Distance



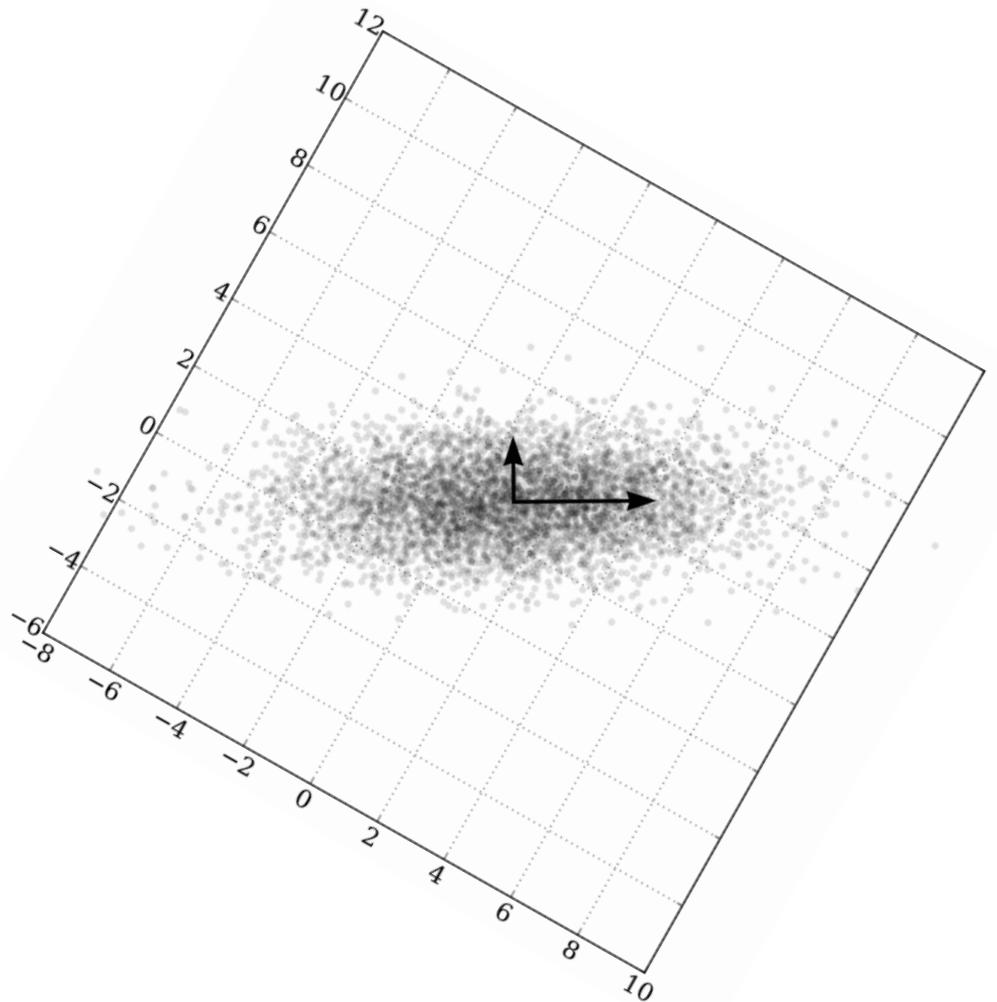
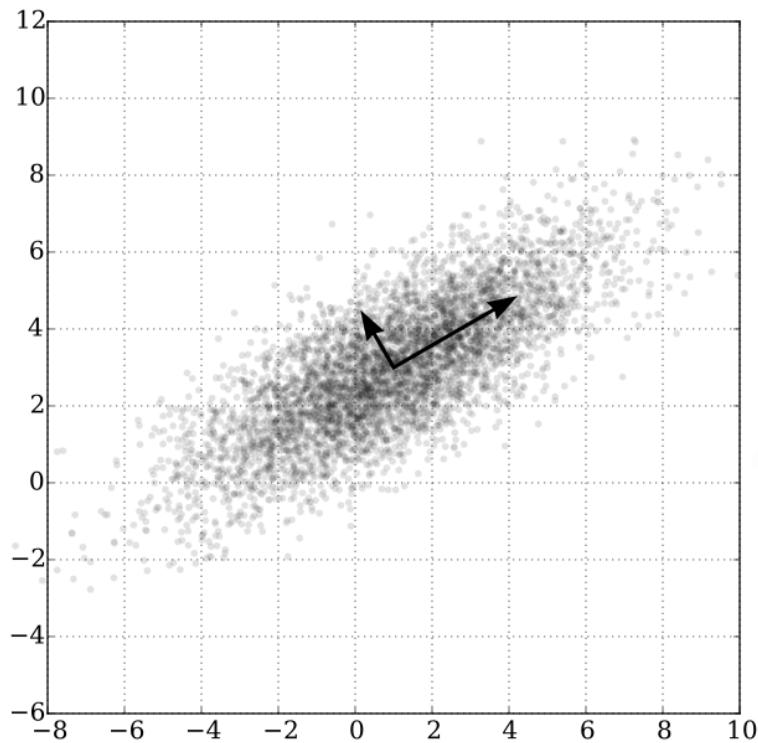
$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}.$$

Computational genetics: Computational analysis of microarray data
Quackenbush (2001) *Nature Reviews Genetics*. doi:10.1038/35076576

Hierarchical Clustering



Principle Components Analysis (PCA)



PC1: “New X”- The dimension with the most variability
PC2: “New Y”- The dimension with the second most variability

Principle Components Analysis (PCA)

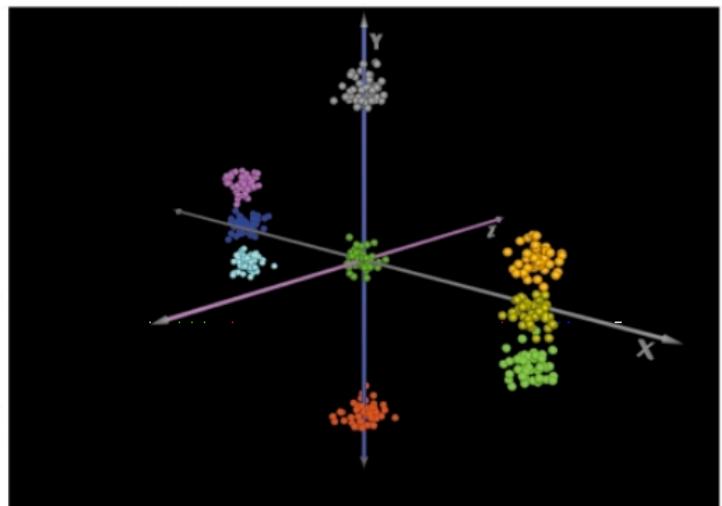
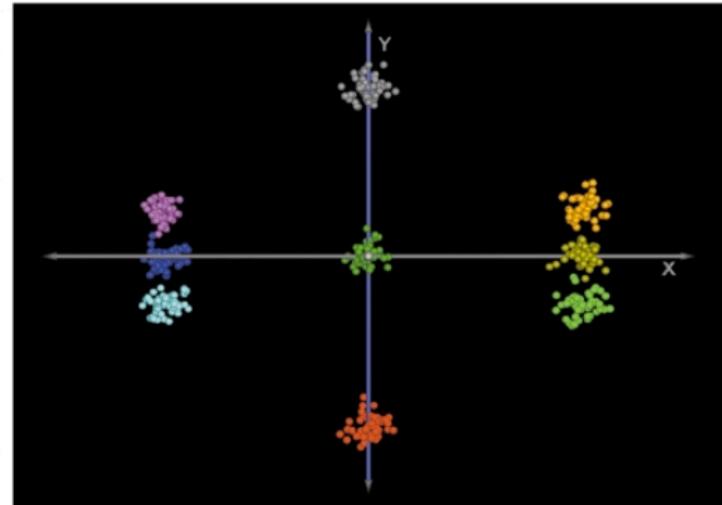
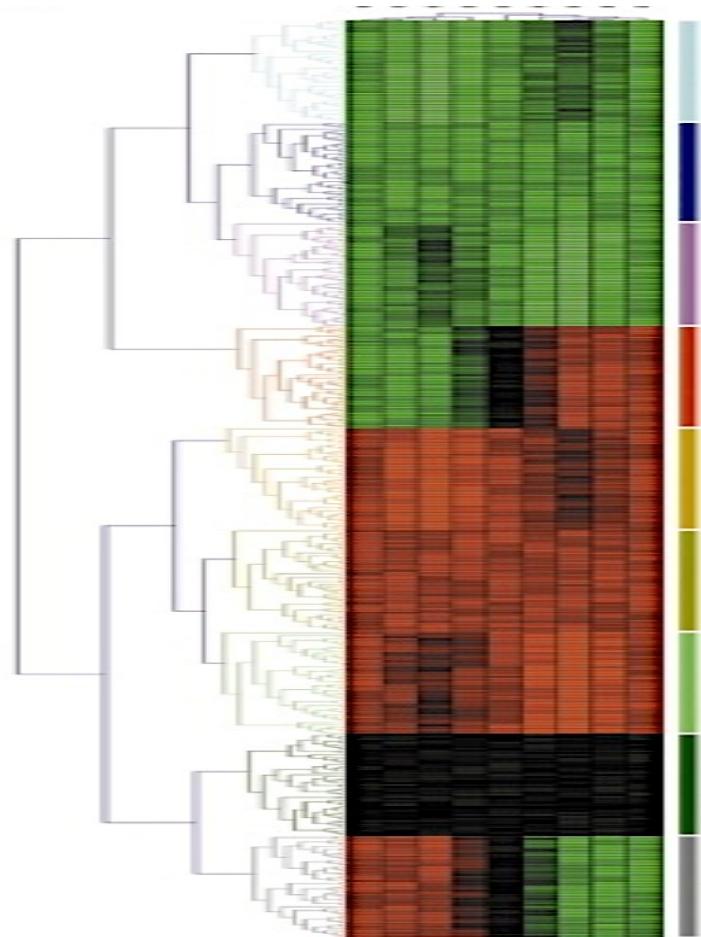
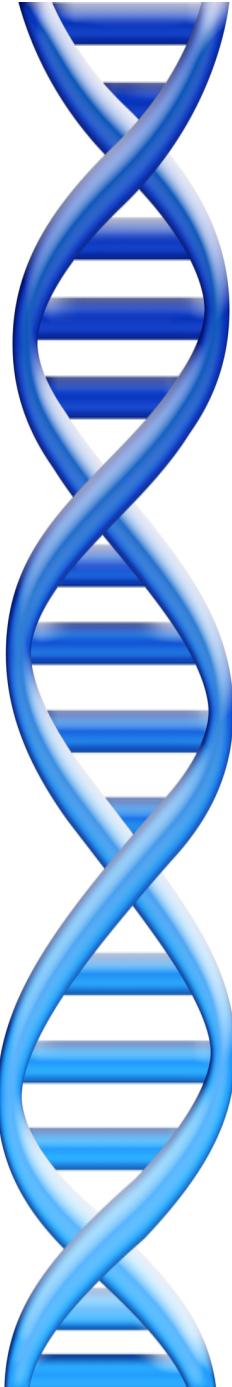


Figure 4 | **Principal component analysis.** The same demonstration data set was analysed using **a** | hierarchical (average-linkage) clustering and **b** | principal component analysis using Euclidean distance, to show how each treats the data, with genes colour coded on the basis of hierarchical clustering results for comparison.



Part 2:

Modern Humans

ARTICLE

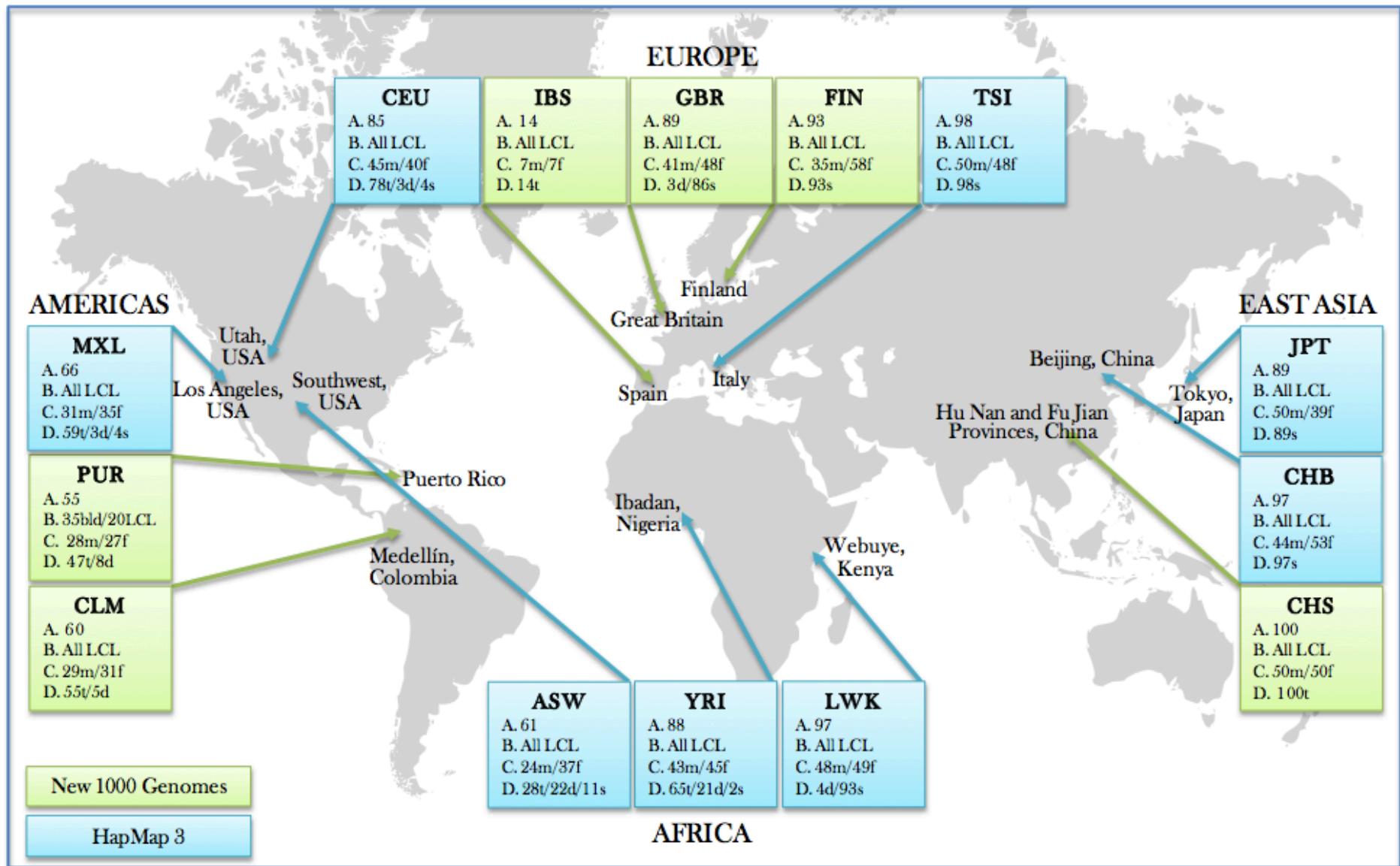
doi:10.1038/nature11632

An integrated map of genetic variation from 1,092 human genomes

The 1000 Genomes Project Consortium*

By characterizing the geographic and functional spectrum of human genetic variation, the 1000 Genomes Project aims to build a resource to help to understand the genetic contribution to disease. Here we describe the genomes of 1,092 individuals from 14 populations, constructed using a combination of low-coverage whole-genome and exome sequencing. By developing methods to integrate information across several algorithms and diverse data sources, we provide a validated haplotype map of 38 million single nucleotide polymorphisms, 1.4 million short insertions and deletions, and more than 14,000 larger deletions. We show that individuals from different populations carry different profiles of rare and common variants, and that low-frequency variants show substantial geographic differentiation, which is further increased by the action of purifying selection. We show that evolutionary conservation and coding consequence are key determinants of the strength of purifying selection, that rare-variant load varies substantially across biological pathways, and that each individual contains hundreds of rare non-coding variants at conserved sites, such as motif-disrupting changes in transcription-factor-binding sites. This resource, which captures up to 98% of accessible single nucleotide polymorphisms at a frequency of 1% in related populations, enables analysis of common and low-frequency variants in individuals from diverse, including admixed, populations.

1000 Genomes Populations



1000 Genomes Populations

Population	DNA sequenced from blood	Offspring Samples from Trios Available	Pilot Samples	Phase 1 Samples	Final Phase Discovery Sample	Final Release Sample	Total
Chinese Dai in Xishuangbanna, China (CDX)	no	yes	0	0	99	93	99
Han Chinese in Beijing, China (CHB)	no	no	91	97	103	103	106
Japanese in Tokyo, Japan (JPT)	no	no	94	89	104	104	105
Kinh in Ho Chi Minh City, Vietnam (KHV)	yes	yes	0	0	101	99	101
Southern Han Chinese, China (CHS)	no	yes	0	100	108	105	112
Total East Asian Ancestry (EAS)			185	286	515	504	523
Bengali in Bangladesh (BEB)	no	yes	0	0	86	86	86
Gujarati Indian in Houston, TX (GIH)	no	yes	0	0	106	103	106
Indian Telugu in the UK (ITU)	yes	yes	0	0	103	102	103
Punjabi in Lahore, Pakistan (PJL)	yes	yes	0	0	96	96	96
Sri Lankan Tamil in the UK (STU)	yes	yes	0	0	103	102	103
Total South Asian Ancestry (SAS)			0	0	494	489	494
African Ancestry in Southwest US (ASW)	no	yes	0	61	66	61	66
African Caribbean in Barbados (ACB)	yes	yes	0	0	96	96	96
Esan in Nigeria (ESN)	no	yes	0	0	99	99	99
Gambian in Western Division, The Gambia (GWD)	no	yes	0	0	113	113	113
Luhya in Webuye, Kenya (LWK)	no	yes	102	97	101	99	116
Mende in Sierra Leone (MSL)	no	yes	0	0	85	85	85
Yoruba in Ibadan, Nigeria (YRI)	no	yes	106	88	109	108	116
Total African Ancestry (AFR)			208	246	669	661	691
British in England and Scotland (GBR)	no	yes	0	89	92	91	94
Finnish in Finland (FIN)	no	no	0	93	99	99	100
Iberian populations in Spain (IBS)	no	yes	0	14	107	107	107
Toscani in Italy (TSI)	no	no	66	98	108	107	110
Utah residents with Northern and Western European ancestry (CEU)	no	yes	94	85	99	99	103
Total European Ancestry (EUR)			160	379	505	503	514
Colombian in Medellin, Colombia (CLM)	no	yes	0	60	94	94	95
Mexican Ancestry in Los Angeles, California (MXL)	no	yes	0	66	67	64	69
Peruvian in Lima, Peru (PEL)	yes	yes	0	0	86	85	86
Puerto Rican in Puerto Rico (PUR)	yes	yes	0	55	105	104	105
Total Americas Ancestry (AMR)			181	352	347	347	355
Total			553	1092	2535	2504	2577

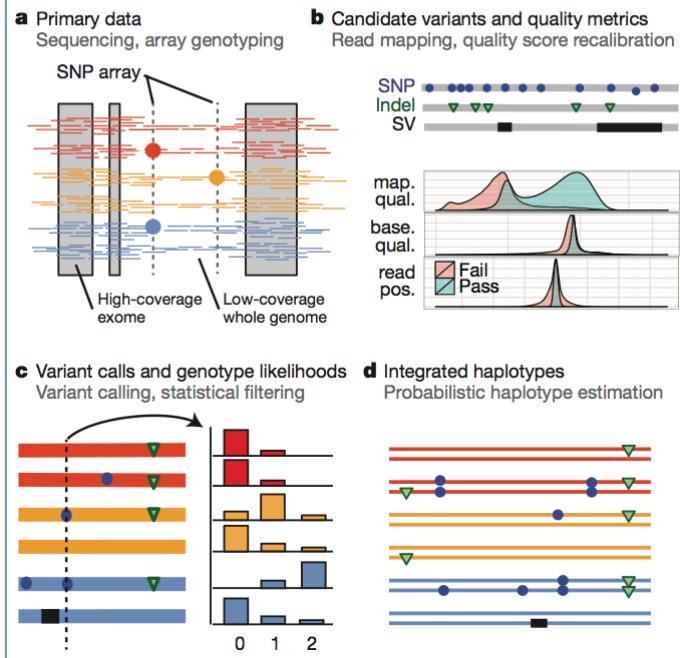
26 populations from 5 major population groups

1000 Genomes: Human Mutation Rate

- Phase I Release
 - 1092 individuals from 14 populations
 - Combination of low coverage WGS, deep coverage WES, and SNP genotype data
- Overall SNP rate between any two people is ~1/1200bp to ~1/1300
 - ~3M SNPs between me and you (.1%)
 - ~30M SNPs between human to Chimpanzees (1%)
- De novo mutation rate ~1/100,000,000
 - ~100 de novo mutations from generation to generation
 - ~1-2 de novo mutations within the protein coding genes

Constructing an integrated map of variation

The 1,092 haplotype-resolved genomes released as phase I by the 1000 Genomes Project are the result of integrating diverse data from multiple technologies generated by several centres between 2008 and 2010. The Box 1 Figure describes the process leading from primary data production to integrated haplotypes.



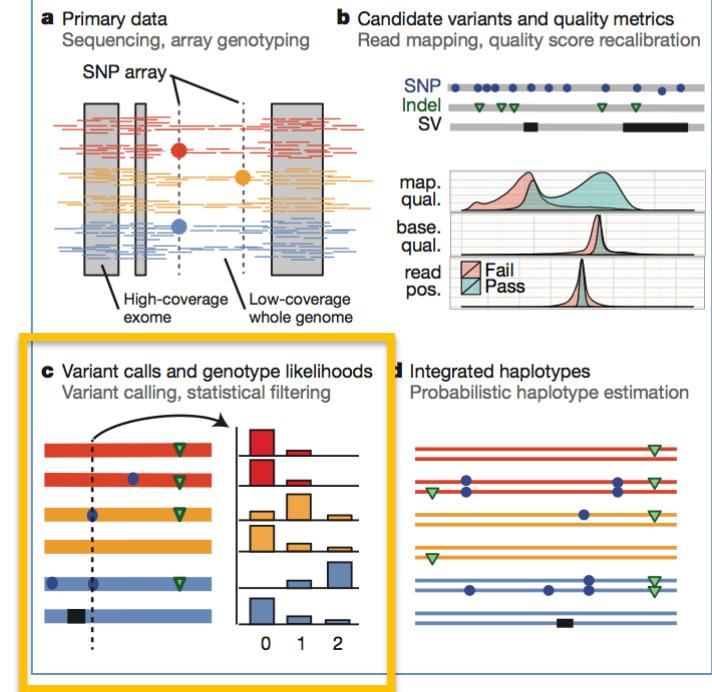
An integrated map of genetic variation from 1,092 human genomes
1000 genomes project (2012) *Nature*. doi:10.1038/nature11632

1000 Genomes: Human Mutation Rate

- Phase I Release
 - 1092 individuals from 14 populations
 - Combination of low coverage WGS, deep coverage WES, and SNP genotype data
- Overall SNP rate between any two people is ~1/1200bp to ~1/1300
 - ~3M SNPs between me and you (.1%)
 - ~30M SNPs between human to Chimpanzees (1%)
- De novo mutation rate ~1/100,000,000
 - ~100 de novo mutations from generation to generation
 - ~1-2 de novo mutations within the protein coding genes

Constructing an integrated map of variation

The 1,092 haplotype-resolved genomes released as phase I by the 1000 Genomes Project are the result of integrating diverse data from multiple technologies generated by several centres between 2008 and 2010. The Box 1 Figure describes the process leading from primary data production to integrated haplotypes.



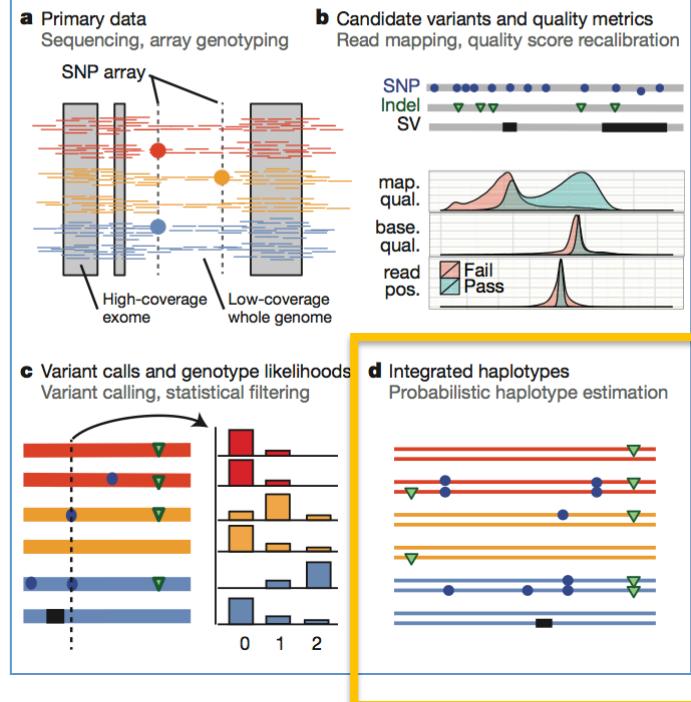
An integrated map of genetic variation from 1,092 human genomes
1000 genomes project (2012) *Nature*. doi:10.1038/nature11632

1000 Genomes: Human Mutation Rate

- Phase I Release
 - 1092 individuals from 14 populations
 - Combination of low coverage WGS, deep coverage WES, and SNP genotype data
- Overall SNP rate between any two people is ~1/1200bp to ~1/1300
 - ~3M SNPs between me and you (.1%)
 - ~30M SNPs between human to Chimpanzees (1%)
- De novo mutation rate ~1/100,000,000
 - ~100 de novo mutations from generation to generation
 - ~1-2 de novo mutations within the protein coding genes

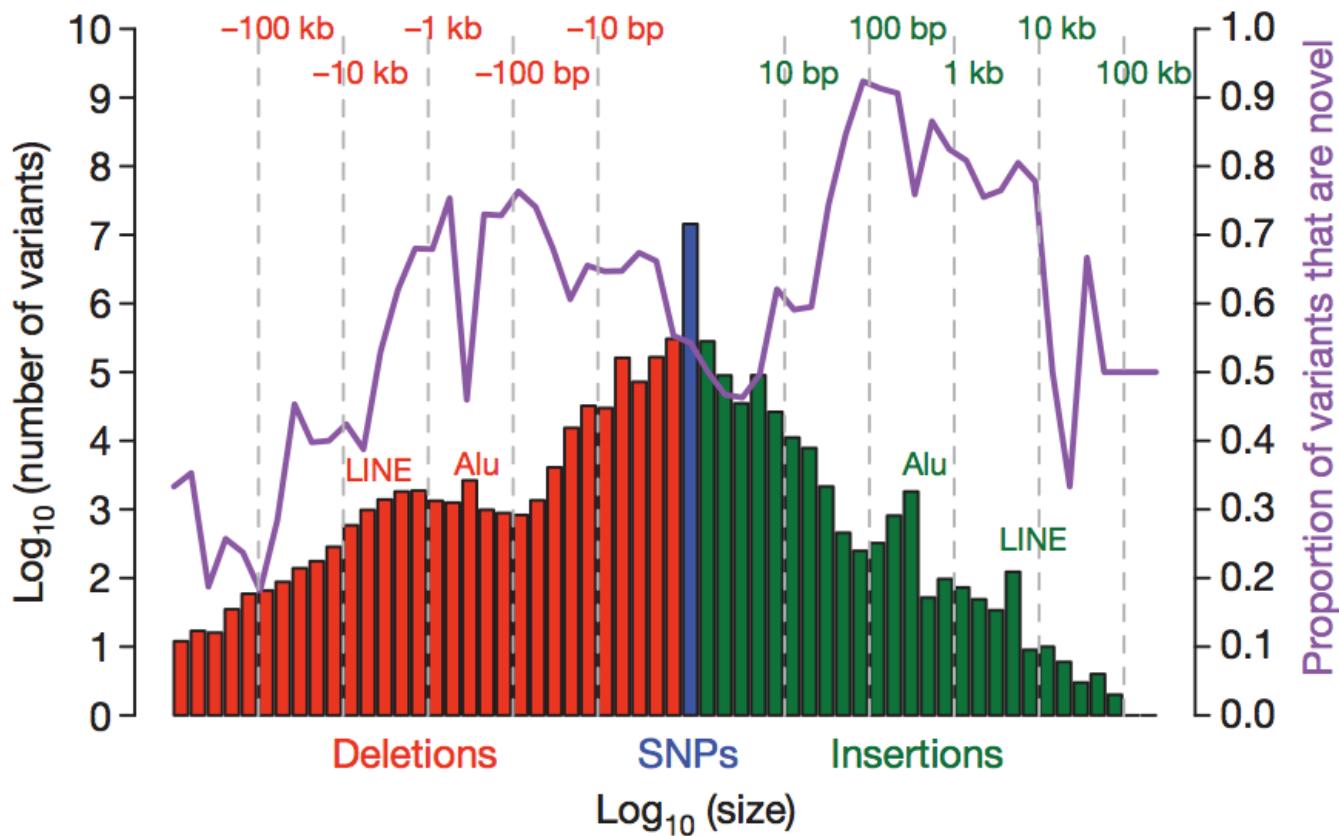
Constructing an integrated map of variation

The 1,092 haplotype-resolved genomes released as phase I by the 1000 Genomes Project are the result of integrating diverse data from multiple technologies generated by several centres between 2008 and 2010. The Box 1 Figure describes the process leading from primary data production to integrated haplotypes.



An integrated map of genetic variation from 1,092 human genomes
1000 genomes project (2012) *Nature*. doi:10.1038/nature11632

Human Mutation Types



- Mutations follows a “log-normal” frequency distribution
 - Most mutations are SNPs followed by small indels followed by larger events

A map of human genome variation from population-scale sequencing
1000 genomes project (2010) *Nature*. doi:10.1038/nature09534

Copy Number Variations

Large-Scale Copy Number Polymorphism in the Human Genome

Jonathan Sebat,¹ B. Lakshmi,¹ Jennifer Troge,¹ Joan Alexander,¹ Janet Young,² Pär Lundin,³ Susanne Månér,³ Hillary Massa,² Megan Walker,² Maoyen Chi,¹ Nicholas Navin,¹ Robert Lucito,¹ John Healy,¹ James Hicks,¹ Kenny Ye,⁴ Andrew Reiner,¹ T. Conrad Gilliam,⁵ Barbara Trask,² Nick Patterson,⁶ Anders Zetterberg,³ Michael Wigler^{1*}

The extent to which large duplications and deletions contribute to human genetic variation and diversity is unknown. Here, we show that large-scale copy number polymorphisms (CNPs) (about 100 kilobases and greater) contribute substantially to genomic variation between normal humans. Representational oligonucleotide microarray analysis of 20 individuals revealed a total of 221 copy number differences representing 76 unique CNPs. On average, individuals differed by 11 CNPs, and the average length of a CNP interval was 465 kilobases. We observed copy number variation of 70 different genes within CNP intervals, including genes involved in neurological function, regulation of cell growth, regulation of metabolism, and several genes known to be associated with disease.

Many of the genetic differences between humans and other primates are a result of large duplications and deletions (1–3). From these observations, it is reasonable to expect that differences in gene copy number could be a significant source of genetic variation between humans. A few examples of large duplication polymorphisms have been reported (4). However, because of previous limitations in the power to determine DNA copy number at high resolution throughout the genome, the extent to which copy number polymorphisms (CNPs) contribute to human genetic diversity is unknown.

In our previous studies of human cancer with the use of representational oligonucleotide microarray analysis (ROMA), we have detected many genomic amplifications and deletions in tumor genomes when analyzed in comparison to an unrelated normal genome (5), but some of these genetic differences could be due to germline CNPs. To correctly interpret genomic data relating to cancer and other diseases, we must distinguish abnormal genetic lesions from normal CNPs.

We used ROMA to investigate the extent of copy number variation between normal

Strong Association of De Novo Copy Number Mutations with Autism

Jonathan Sebat,^{1,*} B. Lakshmi,¹ Dheeraj Malhotra,^{1*} Jennifer Troge,^{1*} Christa Lese-Martin,² Tom Walsh,³ Boris Yamrom,¹ Seungtai Yoon,¹ Alex Krasnitz,¹ Jude Kendall,¹ Anthony Leotta,¹ Deepa Pai,¹ Ray Zhang,¹ Yoon-Ha Lee,¹ James Hicks,¹ Sarah J. Spence,¹ Annette T. Lee,⁵ Kaija Puura,⁶ Terho Lehtimäki,⁷ David Ledbetter,⁷ Peter K. Gregersen,⁵ Joel Bregman,⁸ James S. Sutcliffe,⁹ Vaidehi Jobanputra,¹⁰ Wendy Chung,¹⁰ Dorothy Warburton,¹⁰ Mary-Claire King,³ David Skuse,¹¹ Daniel H. Geschwind,¹² T. Conrad Gilliam,¹³ Kenny Ye,¹⁴ Michael Wigler^{1†}

We tested the hypothesis that de novo copy number variation (CNV) is associated with autism spectrum disorders (ASDs). We performed comparative genomic hybridization (CGH) on the genomic DNA of patients and unaffected subjects to detect copy number variants not present in their respective parents. Candidate genomic regions were validated by higher-resolution CGH, paternity testing, cytogenetics, fluorescence in situ hybridization, and microsatellite genotyping. Confirmed de novo CNVs were significantly associated with autism ($P = 0.0005$). Such CNVs were identified in 12 out of 118 (10%) of patients with sporadic autism, in 2 out of 77 (3%) of patients with an affected first-degree relative, and in 2 out of 196 (1%) of controls. Most de novo CNVs were smaller than microscopic resolution. Affected genomic regions were highly heterogeneous and included mutations of single genes. These findings establish de novo germline mutation as a more significant risk factor for ASD than previously recognized.

Autism spectrum disorders (ASDs) [Mendelian Inheritance in Man (MIM) 209850] are characterized by language impairments, social deficits, and repetitive behaviors. The onset of symptoms occurs by the age of 3 and usually requires extensive support for the lifetime of the afflicted. The prevalence of ASD is estimated to be 1 in 166 (1), making it a major burden to society.

Genetics plays a major role in the etiology of autism. The concordance rates in monozygotic twins are 70% for autism and 90% for ASD, whereas the concordance rates in dizygotic twins are 5% and 10%, respectively. Previous studies suggest autism displays a high degree of genetic heterogeneity. Efforts to map disease genes using linkage analysis have found evidence for autism loci on 20 different chromosomes. Regions implicated by multiple studies include 1p, 5q, 7q, 15q, 16p, 17q, 19p, and Xq (2). Moreover, microscopy studies have identified cytogenetic abnormalities in >5% of affected children, involving many different loci on all chromosomes (3). In some rare syndromic forms of autism, such as Rett syndrome (4) and tuberous

sclerosis (5), mutations in a single gene have been identified. Otherwise, neither linkage nor cytogenetics has unambiguously identified specific genes involved.

Genetic heterogeneity poses a considerable challenge to traditional approaches for gene mapping (6). Some of these limitations are overcome by methods that rely on the direct detection of functional variants, which in most cases are de novo events. New array-based technologies can detect differences in DNA copy number at much higher resolution than cytogenetic methods (7) and, hence, might reveal spontaneous mutations that were previously unidentified. These techniques have shown an abundance of copy number variants (CNVs) in humans (8, 9), and the same methods have been used to find de novo chromosome aberrations below the resolution of microscopy in children with mental retardation and dysmorphic features (10–14), including patients with syndromic forms of autism (15). Yet, the association of spontaneous CNVs in idiopathic autism has not been systematically investigated. Thus, a large-scale study of genome copy number variation in

While fewer numbers of CNVs occur per person, the total number of bases involved can be much greater and have profound effect.

A Systematic Survey of Loss-of-Function Variants in Human Protein-Coding Genes

Daniel G. MacArthur,^{1,2*} Suganthi Balasubramanian,^{3,4} Adam Frankish,¹ Ni Huang,¹ James Morris,¹ Klaudia Walter,¹ Luke Jostins,¹ Lukas Habegger,^{3,4} Joseph K. Pickrell,⁵ Stephen B. Montgomery,^{6,7} Cornelis A. Albers,^{1,8} Zhengdong D. Zhang,⁹ Donald F. Conrad,¹⁰ Gerton Lunter,¹¹ Hancheng Zheng,¹² Qasim Ayub,¹ Mark A. DePristo,¹³ Eric Banks,¹³ Min Hu,¹ Robert E. Handsaker,^{13,14} Jeffrey A. Rosenfeld,¹⁵ Menachem Fromer,¹³ Mike Jin,³ Xinmeng Jasmine Mu,^{3,4} Ekta Khurana,^{3,4} Kai Ye,¹⁶ Mike Kay,¹ Gary Ian Saunders,¹ Marie-Marthe Suner,¹ Toby Hunt,¹ If H. A. Barnes,¹ Clara Amid,^{1,17} Denise R. Carvalho-Silva,¹ Alexandra H. Bignell,¹ Catherine Snow,¹ Bryndis Yngvadottir,¹ Suzannah Bumpstead,¹ David N. Cooper,¹⁸ Yali Xue,¹ Irene Gallego Romero,^{1,5} 1000 Genomes Project Consortium, Jun Wang,¹² Yingrui Li,¹² Richard A. Gibbs,¹⁹ Steven A. McCarroll,^{13,14} Emmanouil T. Dermitzakis,⁷ Jonathan K. Pritchard,^{5,20} Jeffrey C. Barrett,¹ Jennifer Harrow,¹ Matthew E. Hurles,¹ Mark B. Gerstein,^{3,4,21†} Chris Tyler-Smith^{1†}

Genome-sequencing studies indicate that all humans carry many genetic variants predicted to cause loss of function (LoF) of protein-coding genes, suggesting unexpected redundancy in the human genome. Here we apply stringent filters to 2951 putative LoF variants obtained from 185 human genomes to determine their true prevalence and properties. We estimate that human genomes typically contain ~100 genuine LoF variants with ~20 genes completely inactivated. We identify rare and likely deleterious LoF alleles, including 26 known and 21 predicted severe disease-causing variants, as well as common LoF variants in nonessential genes. We describe functional and evolutionary differences between LoF-tolerant and recessive disease genes and a method for using these differences to prioritize candidate genes found in clinical sequencing studies.

Homozygous LoF Mutations

LETTER

doi:10.1038/nature22034

Human knockouts and phenotypic analysis in a cohort with a high rate of consanguinity

Danish Saleheen^{1,2*}, Pradeep Natarajan^{3,4*}, Irina M. Armean^{4,5}, Wei Zhao¹, Asif Rasheed², Sumeet A. Khetarpal⁶, Hong-Hee Won⁷, Konrad J. Karczewski^{1,5}, Anne H. O'Donnell-Luria^{4,5,8}, Caitlin E. Samocha^{4,5}, Benjamin Weisburd^{4,5}, Namrata Gupta^{1*}, Mozzam Zaidi⁹, Maria Samuel⁹, Atif Imran², Shahid Abbas⁹, Faisal Majeed⁹, Madhiha Ishaq², Saba Akhtar², Kevin Trindade⁶, Megan Mucksavage⁶, Nadeem Qamar¹⁰, Khan Shah Zaman¹⁰, Zia Yaqoob¹⁰, Tahir Saghir¹⁰, Syed Nadeem Hasan Rizvi¹⁰, Anis Memon¹⁰, Nadeem Hayyat Mallick¹¹, Mohammad Ishaq¹², Syed Zahed Rasheed¹², Fazal-ur-Rehman Memon¹³, Khalid Mahmood¹⁴, Naveeduddin Ahmed¹⁵, Ron Do^{16,17}, Ronald M. Krauss¹⁸, Daniel G. MacArthur^{1,5}, Stacey Gabriel¹⁴, Eric S. Lander⁴, Mark J. Daly^{4,5}, Philippe Frossard^{2,8}, John Danesh^{19,20§}, Daniel J. Rader^{2,21§} & Sekar Kathiresan^{3,4§}

A major goal of biomedicine is to understand the function of every gene in the human genome¹. Loss-of-function mutations can disrupt both copies of a given gene in humans and phenotypic analysis of such ‘human knockouts’ can provide insight into gene function. Consanguineous unions are more likely to result in offspring carrying homozygous loss-of-function mutations. In Pakistan, consanguinity rates are notably high². Here we sequence the protein-coding regions of 10,503 adult participants in the Pakistan Risk of Myocardial Infarction Study (PROMIS), designed to understand the determinants of cardiometabolic diseases in individuals from South Asia³. We identified individuals carrying homozygous predicted loss-of-function (pLoF) mutations, and performed phenotypic analysis involving more than 200 biochemical and disease traits. We enumerated 49,138 rare (<1% minor allele frequency) pLoF mutations. These pLoF mutations are estimated to knock out 1,317 genes, each in at least one participant. Homozygosity for pLoF mutations at *PLA2G7* was associated with absent enzymatic activity of soluble lipoprotein-associated phospholipase A2; at *CYP2F1*, with higher plasma interleukin-8 concentrations; at *TREH*, with lower concentrations of apoB-containing lipoprotein subfractions; at either *A3GALT2* or *NRG4*, with markedly reduced plasma insulin C-peptide concentrations; and at *SLC9A3RL*, with mediators of calcium and phosphate signalling. Heterozygous deficiency of *APOC3* has been shown to protect against coronary heart disease^{4,5}; we identified *APOC3* homozygous pLoF carriers in our cohort. We recruited these human knockouts and challenged them with an oral fat load. Compared with family members lacking the mutation, individuals with *APOC3* knocked out displayed marked blunting of the usual post-prandial rise in plasma triglycerides. Overall, these observations provide a roadmap for a ‘human knockout project’, a systematic effort to understand the phenotypic consequences of complete disruption of genes in humans.

Across all participants (Table 1), exome sequencing yielded 1,639,223 exonic and splice-site sequence variants in 19,026 autosomal genes that passed initial quality control metrics. Of these, 57,137 mutations

across 14,345 autosomal genes were annotated as pLoF mutations (that is, nonsense, frameshift, or canonical splice-site mutations predicted to inactivate a gene). To increase the probability that mutations are correctly annotated as pLoF by automated algorithms, we removed nonsense and frameshift mutations occurring within the last 5% of the transcript and within exons flanked by non-canonical splice sites, splice-site mutations at small (<15 bp) introns, at non-canonical splice sites, and where the purported pLoF allele is observed across primates. Common pLoF alleles are less likely to exert strong functional effects as they are less constrained by purifying selection; thus, we define pLoF mutations in the rest of the manuscript as variants with a minor allele frequency (MAF) of <1% and passing the aforementioned bioinformatic filters. Applying these criteria, we generated a set of 49,138 pLoF mutations across 13,074 autosomal genes. The site-frequency spectrum for these pLoF mutations revealed that the majority was seen only in one or a few individuals (Extended Data Fig. 1).

Across all 10,503 PROMIS participants, both copies of 1,317 distinct genes were predicted to be inactivated owing to pLoF mutations. A full listing of all 1,317 genes knocked out, the number of knockout participants for each gene, and the specific pLoF mutation(s) are provided in Supplementary Table 1. 891 (67.7%) of the genes were knocked out only in one participant (Fig. 1a). Nearly 1 in 5 of the participants that were sequenced (1,843 individuals, 17.5%) had at least one gene knocked out by a homozygous pLoF mutation. 1,504 of these 1,843 individuals (81.6%) were homozygous pLoF carriers for just one gene, but the minority of participants had more than one gene knocked out and one participant had six genes with homozygous pLoF genotypes.

We compared the coefficient of inbreeding (*F* coefficient) in PROMIS participants with that of 15,249 individuals from outbred populations of European or African American ancestry. The *F* coefficient estimates the excess homozygosity compared with an outbred ancestor. PROMIS participants had a fourfold higher median inbreeding coefficient compared to outbred populations (0.016 versus 0.0041; $P < 2 \times 10^{-16}$) (Fig. 1b). Additionally, those in PROMIS who reported that their parents were closely related had even higher median inbreeding coefficients than

- **Homozygous LoF mutations are rare in most people, but enriched in people born from consanguineous relationships**
- **Sequence the exomes of many such people, find their homozygous LoFs, relate to 200 biochemical or disease traits**
- **A “natural” experiment to understand what genes do: people with both copies of *APOC3* disabled can clear fat from their bloodstream much faster than others, suggests we should develop compounds to prevent heart attacks**

Variation across populations

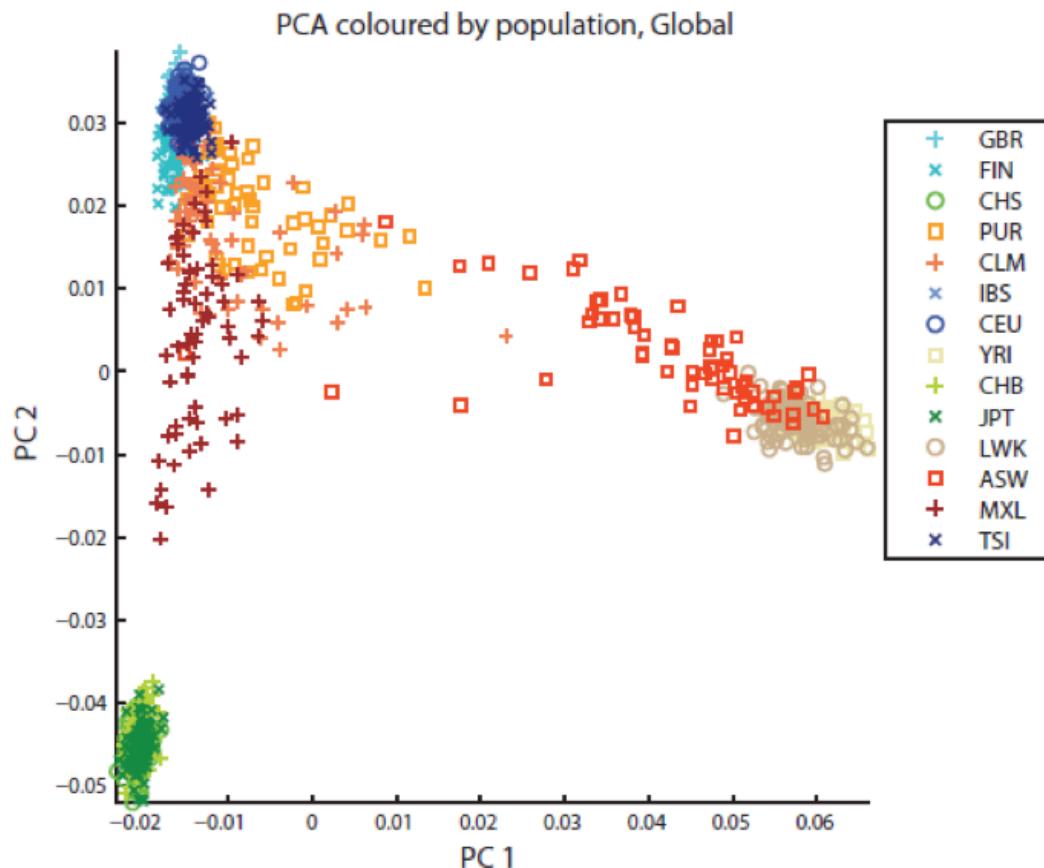


Table S12A Summary of sites showing high levels of population differentiation

- Not a single variant 100% unique to a given population
- 17% of low-frequency variants (.5-5% pop. freq) observed in a single ancestry group
- 50% of rare variants (<.5%) observed in a single population

Variation across populations

Europeans

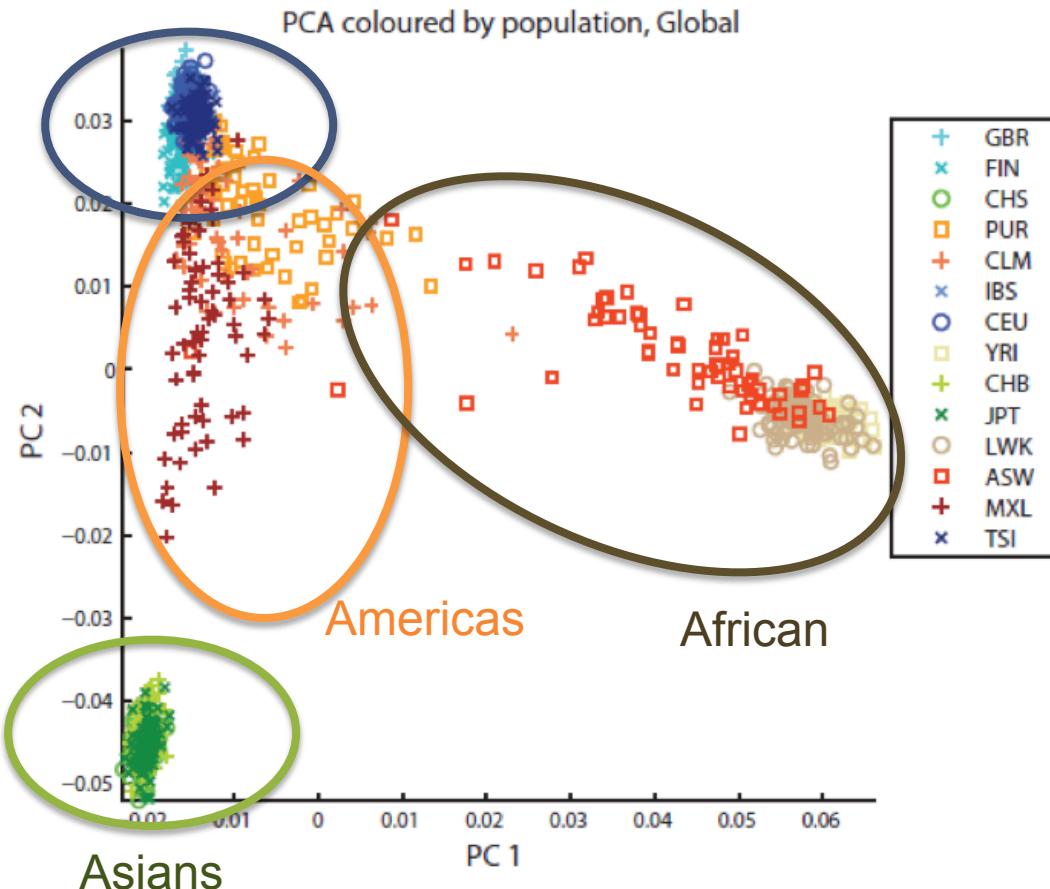
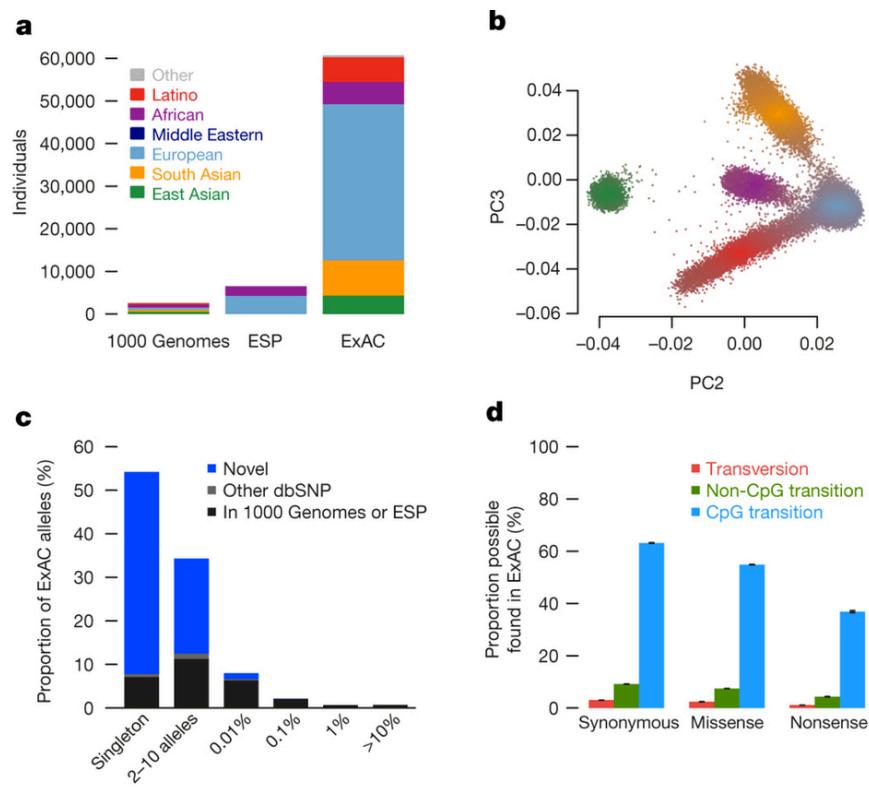


Table S12A Summary of sites showing high levels of population differentiation

- Not a single variant 100% unique to a given population
- 17% of low-frequency variants (.5-5% pop. freq) observed in a single ancestry group
- 50% of rare variants (<.5%) observed in a single population

ExAC: Exome Aggregation Consortium



- The aggregation and analysis of high-quality exome (protein-coding region) DNA sequence data for **60,706 individuals**
- This catalogue of human genetic diversity contains an average of **one variant every eight bases of the exome**
- We have used this catalogue to calculate objective metrics of pathogenicity for sequence variants, and to identify genes subject to strong selection against various classes of mutation; **identifying 3,230 genes with near-complete depletion of predicted protein-truncating**

Analysis of protein-coding genetic variation in 60,706 humans

Lek et al (2016) Nature. doi:10.1038/nature19057

dbSNP

NCBI

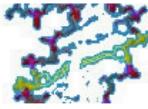
PubMed Nucleotide Protein Genome Structure PopSet Taxonomy OMIM Books SNP

Search for SNP on NCBI Reference Assembly
Search Entrez SNP for Go

Have a question about dbSNP? Try searching the SNP FAQ Archive!
Go

GENERAL RSS Feed Contact Us Site Map dbSNP Homepage Announcements dbSNP Summary FTP Download HUMAN VARIATION SNP SUBMISSION DOCUMENTATION SEARCH RELATED SITES

dbSNP Short Genetic Variations



dbSNP Summary

RELEASE: NCBI dbSNP Build 141

dbSNP Component Availability Dates:

Component	Date available
dbSNP web query for build 141:	May 21, 2014
ftp data for build 141:	May 21, 2014
Entrez Indexing for build 141:	May 21, 2014
BLAST database for build 141:	May 21, 2014

- The complete data for build 141 are available at <ftp://ftp.ncbi.nlm.nih.gov/snp/> in multiple formats.
- All formats and conventions are described in <ftp://ftp.ncbi.nlm.nih.gov/snp/00readme.txt>.
- Please address any questions or comments regarding the data to snp-admin@ncbi.nlm.nih.gov.

New Submission since previous build:

Organism	Current Build	New Submissions (ss#s)	New RefSNP Clusters (rs#s) (# validated)	New ss# with Genotype	New ss# with Frequency
Homo sapiens	141	20,708,470	137 (0)		4
Total: 1 Organisms		20,708,470	137 (0)		4

**Submissions received after reclustering of current build will appear as new rs# clusters in the next build.*

BUILD STATISTICS:

Organism	dbSNP Build	Genome Build	Number of Submissions (ss#s)	Number of RefSNP Clusters (rs#s) (# validated)	Number of (rs#s) in gene	Number of (ss#s) with genotype	Number of (ss#s) with frequency
Homo sapiens	141	38.1	260,570,204	62,387,983 (43,737,321)	29,901,117	73,909,256	35,997,943
Total: 1 Organisms		0 genomes	260,570,204	62,387,983 (43,737,321)	29,901,117	73,909,256	35,997,943

- Periodic release of databases of known variants and their population frequencies
- Generally assumed to be non-disease related
- However, as catalog grows, almost certainly to contain some medically relevant SNPs.

ClinVar

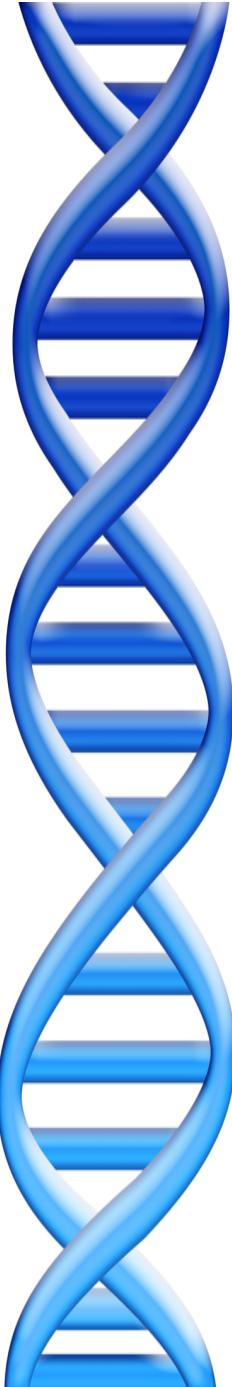
The screenshot shows the ClinVar website at https://www.ncbi.nlm.nih.gov/clinvar/. The page has a dark header with the NCBI logo and a search bar. Below the header is a navigation bar with links for Home, About, Access, Help, Submit, Statistics, and FTP. The main content area features a large sequence snippet on the left and a central box with the ClinVar logo and a description of the database's purpose. On the right, there are three columns: 'Using ClinVar' (About ClinVar, Data Dictionary, Downloads/FTP site, FAQ, Contact Us, RSS feed/What's new?, Factsheet), 'Tools' (ACMG Recommendations for Reporting of Incidental Findings, ClinVar Submission Portal, Submissions, Variation Viewer, Clinical Remapping - Between assemblies and RefSeqGenes, RefSeqGene/LRG), and 'Related Sites' (ClinGen, GeneReviews®, GTR®, MedGen, OMIM®, Variation). A 'Submitter highlights' section at the bottom left lists acknowledgments, RSS feed information, and a link to submitter details.

- ClinVar is a freely accessible, public archive of reports of the relationships among human variations and phenotypes, with supporting evidence
- Currently has 295k mutations
- Most (179k) variants have uncertain affect, only 23 have “4 stars” of significance

OMIM

The screenshot shows the OMIM homepage. At the top left is a banner for '50 YEARS OF OMIM: Human Genetics Knowledge for the World'. Below it, the title 'OMIM®' is displayed in large, bold letters, followed by 'Online Mendelian Inheritance in Man®'. A subtitle 'An Online Catalog of Human Genes and Genetic Disorders' is present, along with a note that it was 'Updated April 7, 2017'. A search bar at the top right contains the placeholder text 'Search OMIM for clinical features, phenotypes, genes, and more...' with a magnifying glass icon. Below the search bar are links for 'Advanced Search', 'Need help?', and 'Mirror site'. A note below the search bar states: 'OMIM is supported by a grant from NHGRI, licensing fees, and generous contributions from people like you.' A 'Make a donation!' button is located near the bottom left. Logos for 'McKUSICK-NATHANS Institute of Genetic Medicine' and 'JOHNS HOPKINS MEDICINE' are visible. A link to 'Follow us on Twitter' is also present. At the bottom, a note reads: 'NOTE: OMIM is intended for use primarily by physicians and other professionals concerned with genetic disorders, by genetics researchers, and by advanced students in science and medicine. While the OMIM database is open to the public, users seeking information about a personal medical or genetic condition are urged to consult with a qualified physician for diagnosis and for answers to personal questions.' It also mentions that 'OMIM® and Online Mendelian Inheritance in Man® are registered trademarks of the Johns Hopkins University.' and 'Copyright © 1966-2017 Johns Hopkins University.'

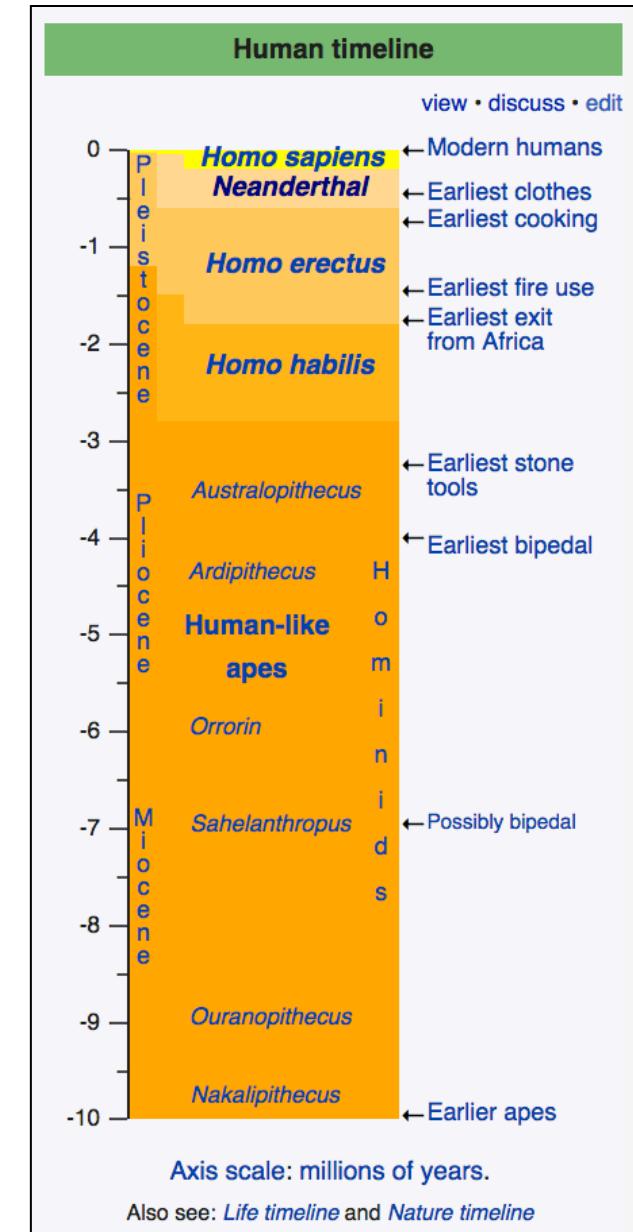
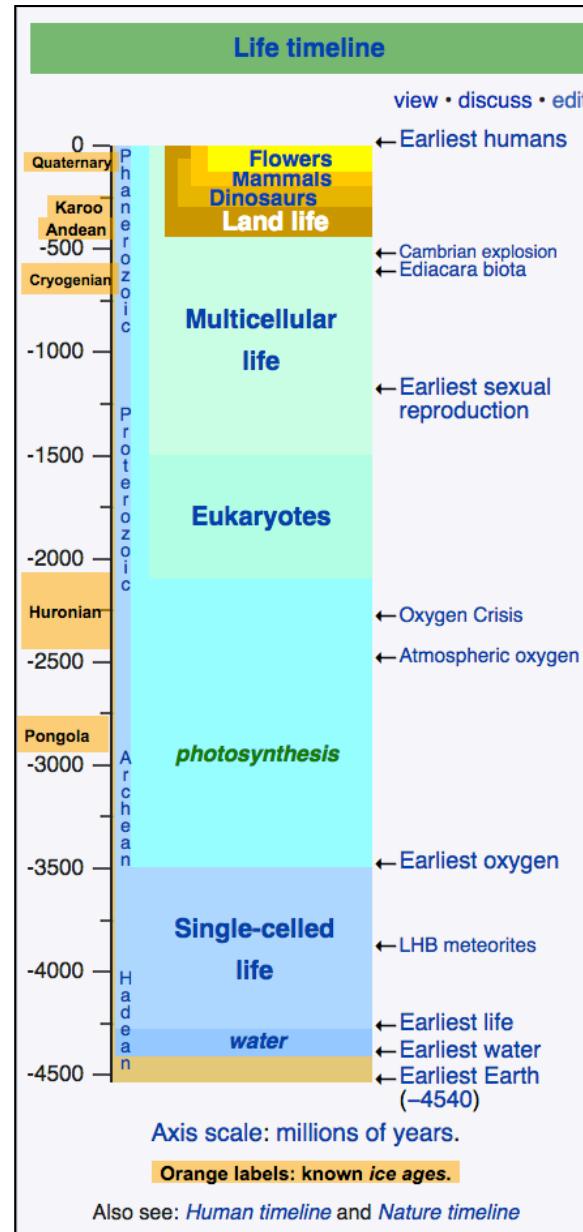
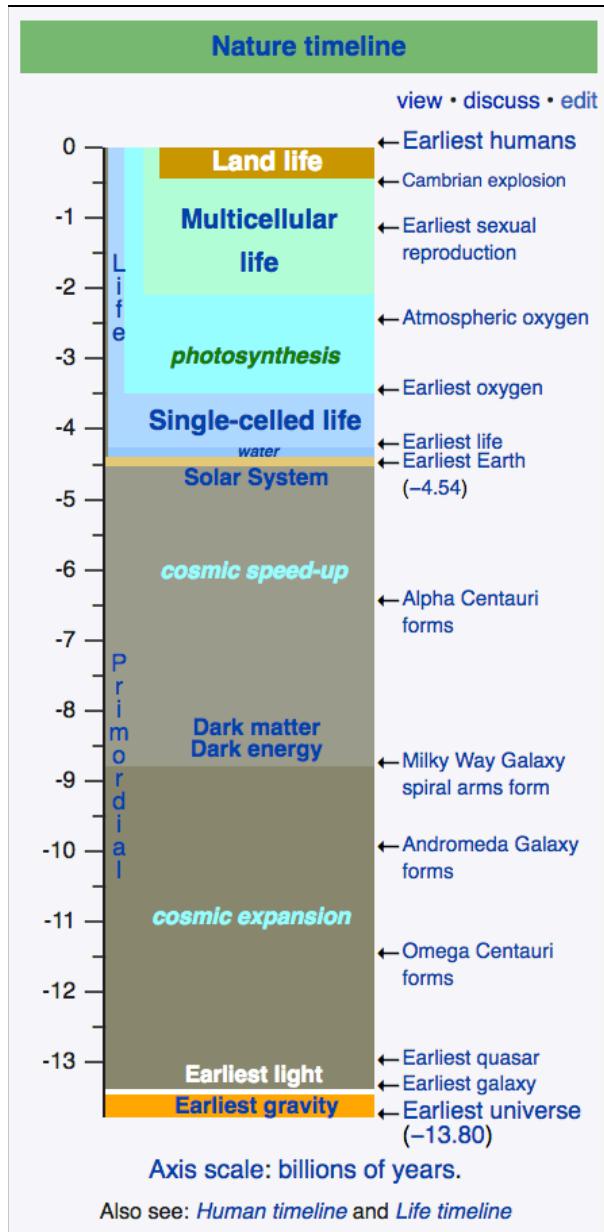
- For many different diseases and phenotypes, lists what are all of the known genetic associations
- Has records for nearly all genes, ~5k different conditions with known molecular basis, ~1k with unknown basis, ~1k with questionable basis
- Started at JHU 50 years ago 😊



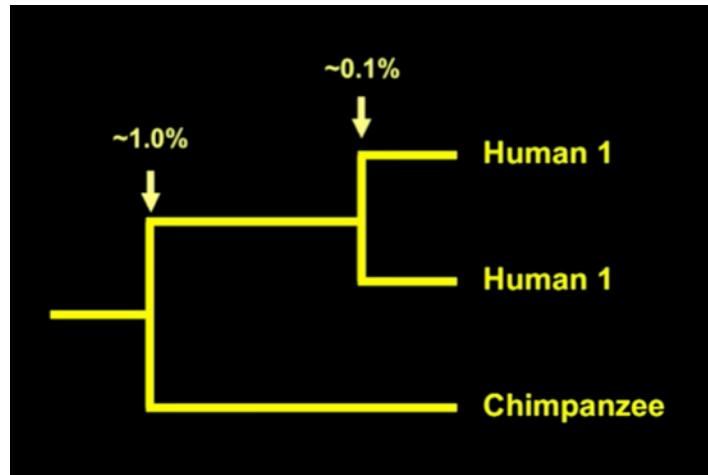
Part 3:

Ancient Hominds

Our Origins



Mutation Rates and Evolutionary Time



Since mutation occur as a function of time we can use the number of mutation to age when different populations split

Interestingly, there is much more variability within Africa than outside of Africa despite the much smaller population

We see “African” alleles all around the world

- Zero SNPs occur exclusively in Africa
- Only 12 SNPs across the entire genome ‘unique’ to Africa (allowing 95% tolerance)
- We are all African (either currently living in Africa or recent exiles)!

Open question if/how early modern humans interacted with earlier hominid

DNA clues to our inner neanderthal

Svante Pääbo (2011). *TED Global*.

https://www.ted.com/talks/svante_paaebo_dna_clues_to_our_inner_neanderthal

A
A
C
G T G C G
A T C g G
G A c C A c G
T T G A T G
G T A T T G A
T C A C G A G
G T G A C C A C
T A T C T G T A
A C A T C T G T G
C T A T G G A T A C C G
A C A C G G A T C T G
C T A G A T A A G A C T
T C T G C T G A T C A C
A T G T C T G A C T G A
T C A G C T G A G c
T A T G C T T G
T A T G C T A C A
G T

Sequencing ancient genomes

Janet Kelso

Max-Planck Institute



Homo neanderthalensis

- Proto-Neanderthals emerge around 600k years ago
- “True” Neanderthals emerge around 200k years ago
- Died out approximately 40,000 years ago
- Known for their robust physique
- Made advanced tools, probably had a language (the nature of which is debated and likely unknowable) and lived in complex social groups



Homo sapiens sapiens

- Apparently emerged from earlier hominids in Africa around 50k years ago
- Capable of amazing intellectual and social behaviors
- Mostly Harmless ☺



A Draft Sequence of the Neandertal Genome
Richard E. Green, et al.
Science **328**, 710 (2010);
DOI: 10.1126/science.1188021

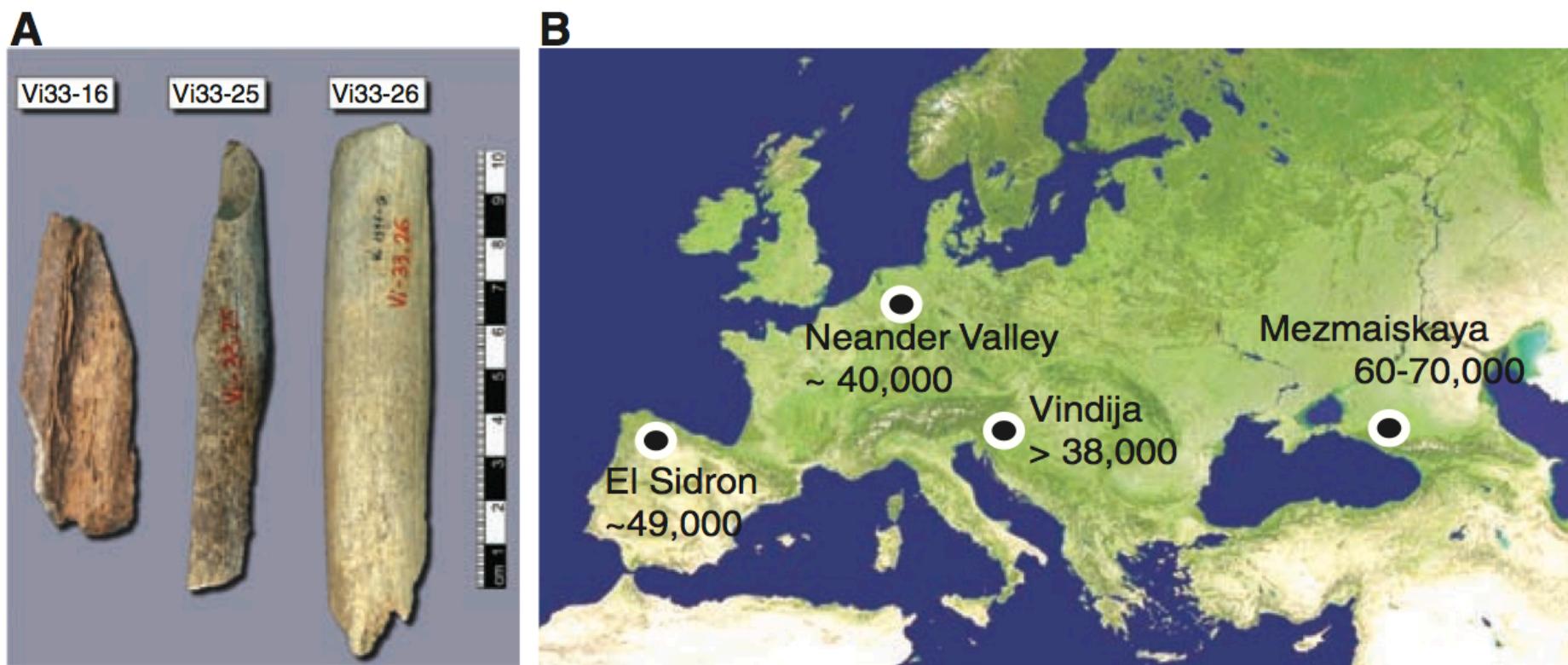
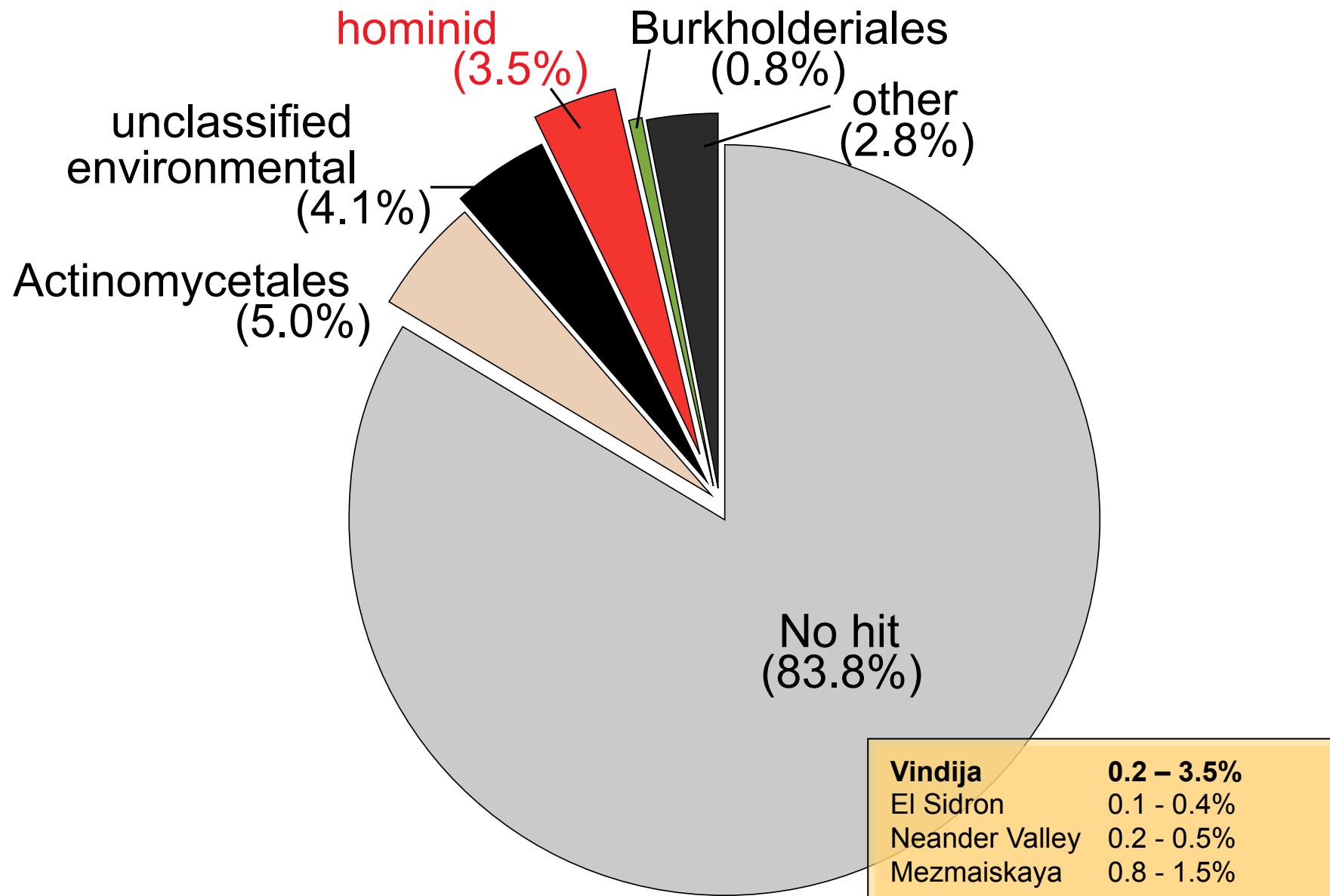


Fig. 1. Samples and sites from which DNA was retrieved. **(A)** The three bones from Vindija from which Neandertal DNA was sequenced. **(B)** Map showing the four archaeological sites from which bones were used and their approximate dates (years B.P.).

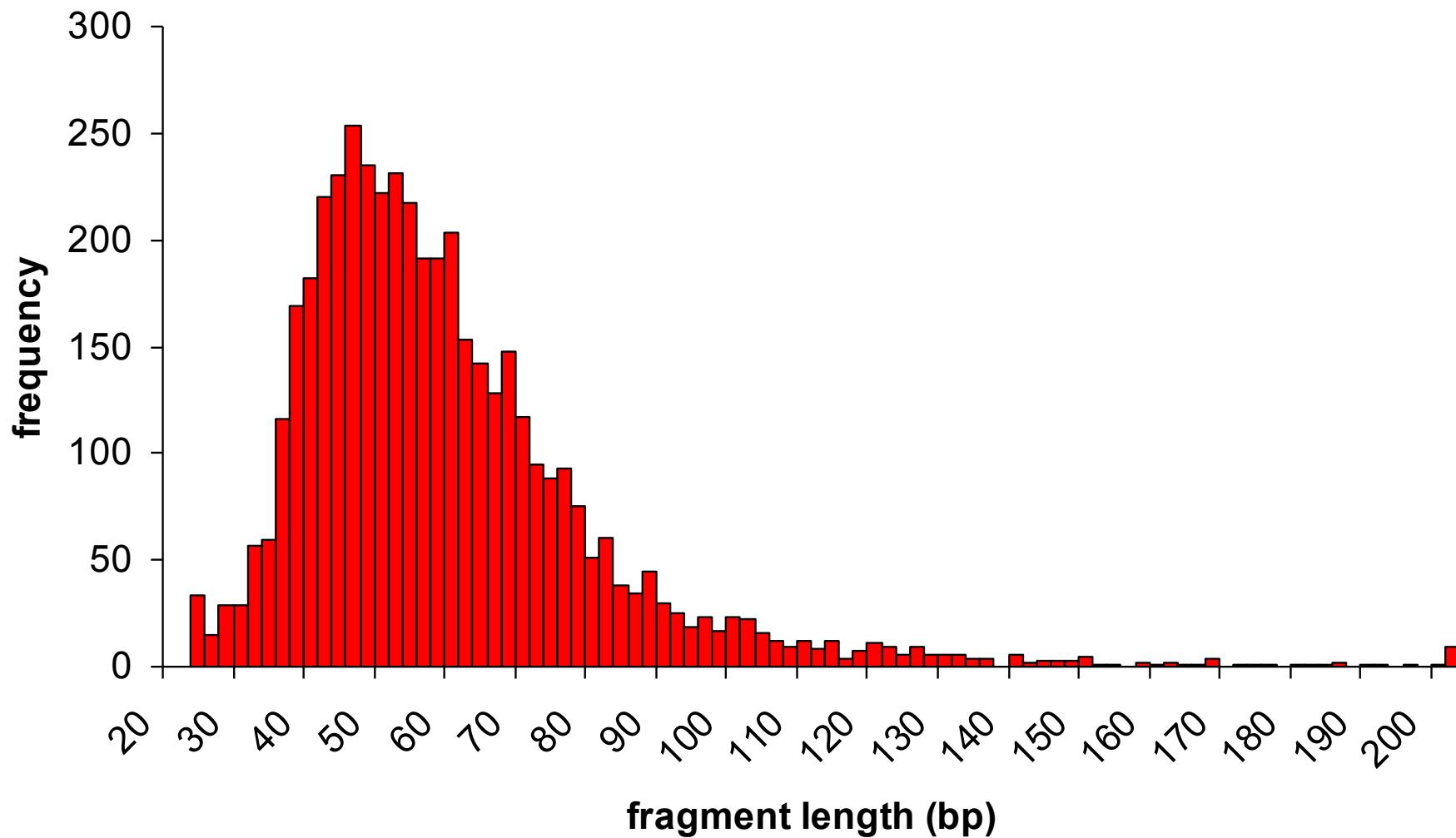
Extracting Ancient DNA



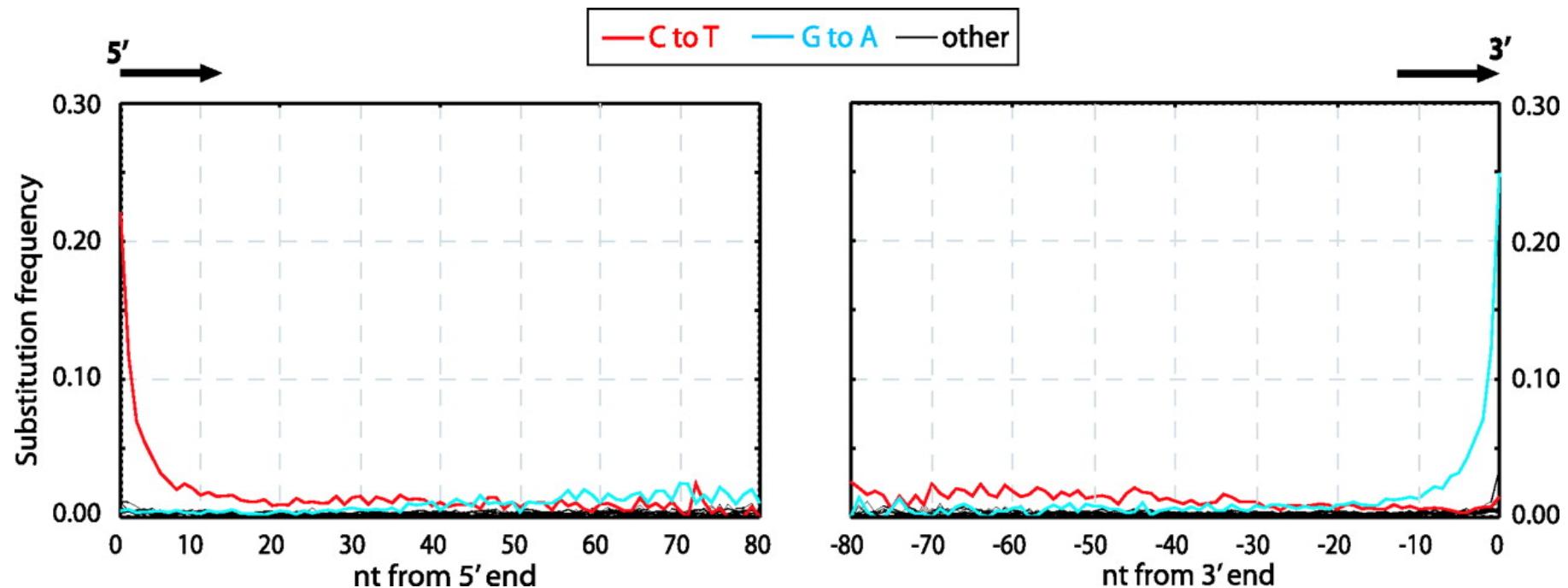
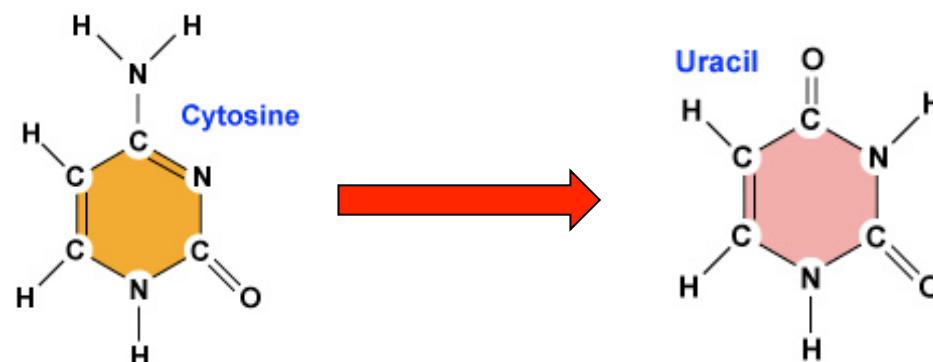
DNA is from mixed sources



DNA is degraded



DNA is chemically damaged





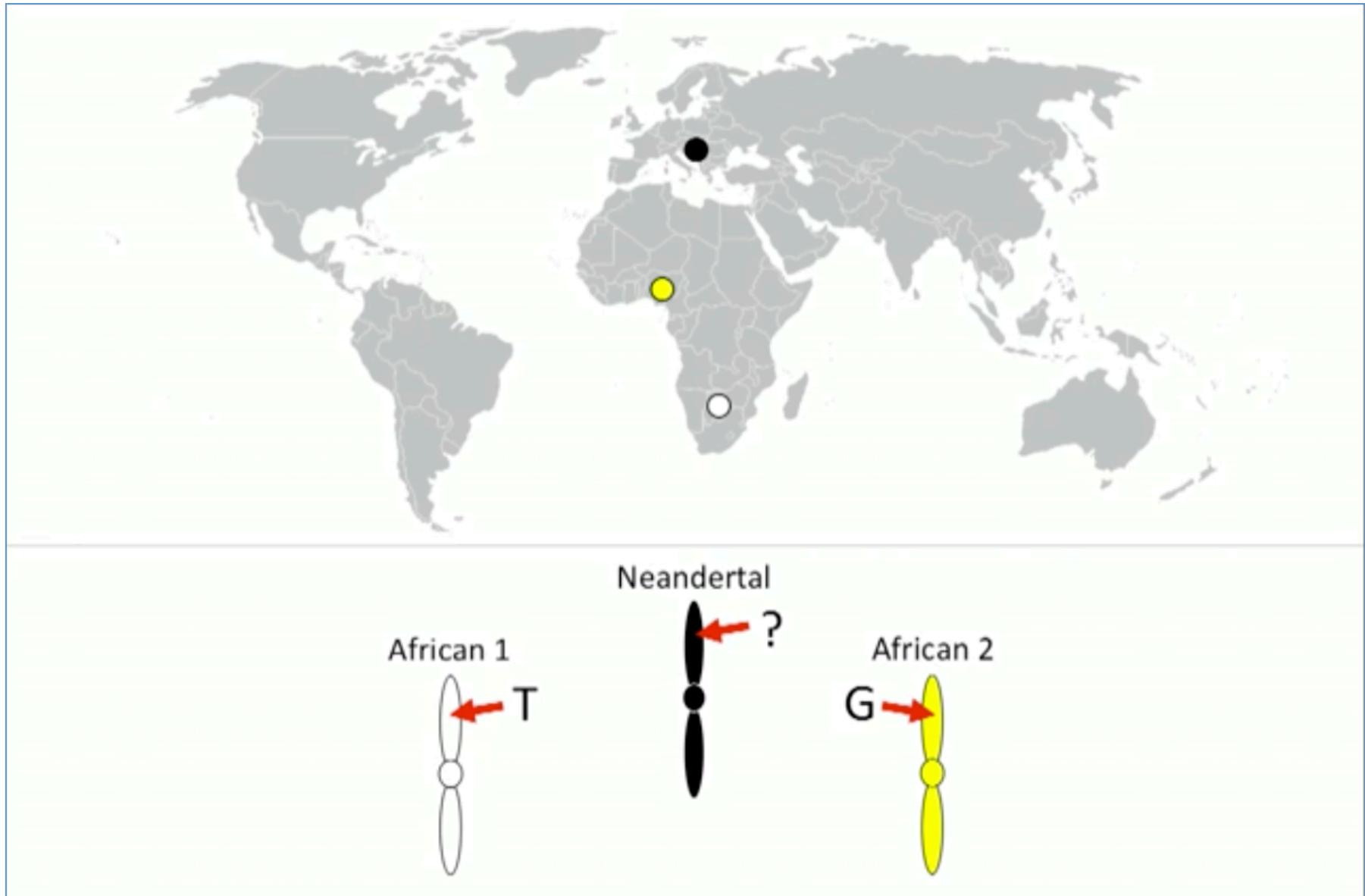
Green et al. 2010

Vindija	33.16	~1.2 Gb
	33.25	~1.3 Gb
	33.26	~1.5 Gb
El Sidron (1253)	~2.2 Mb	
Feldhofer 1	~2.2 Mb	
Mezmaiskaya 1	~56.4 Mb	

~35 Illumina flow cells

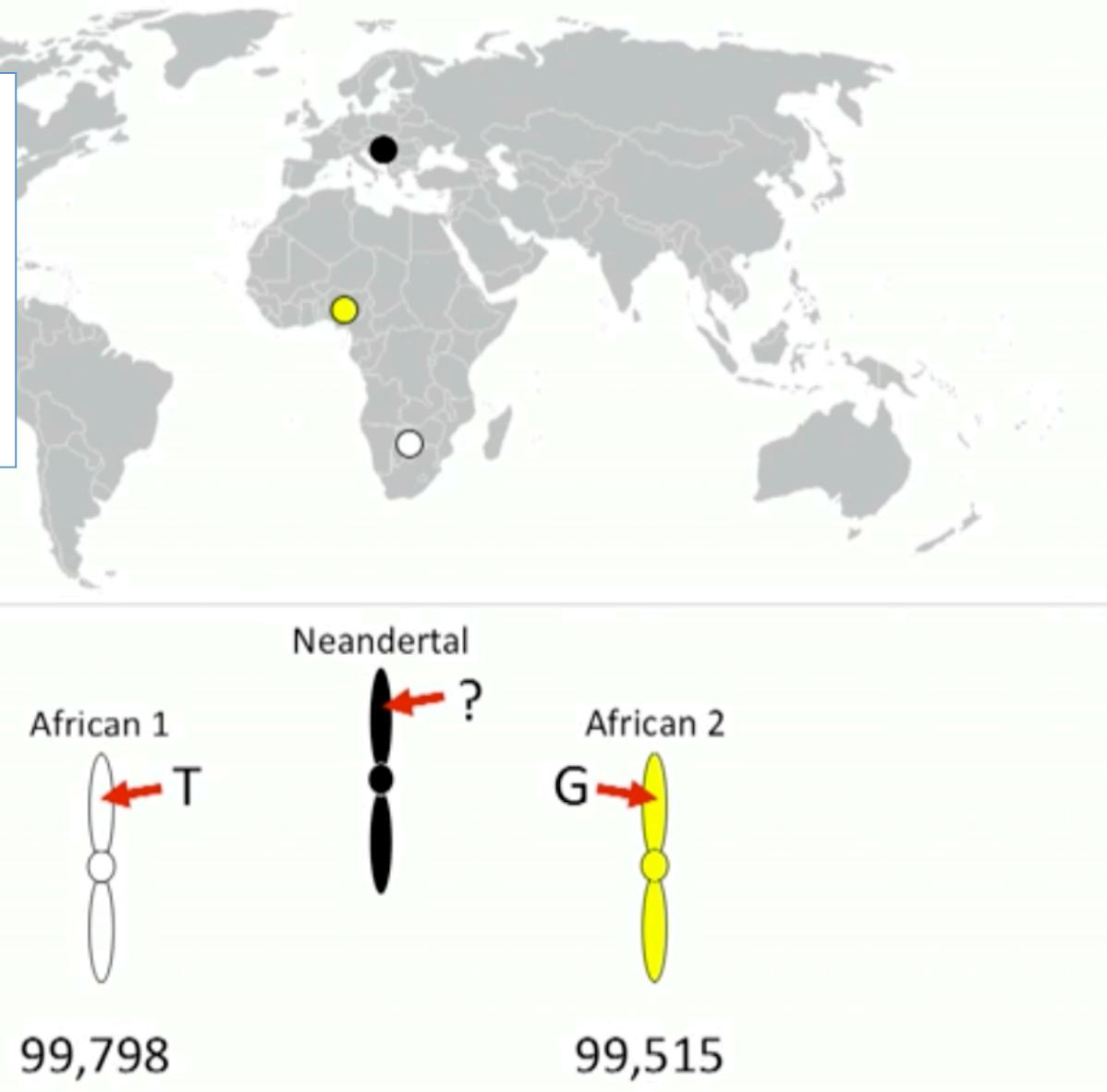
Genome coverage ~1.3 X

Did we mix?



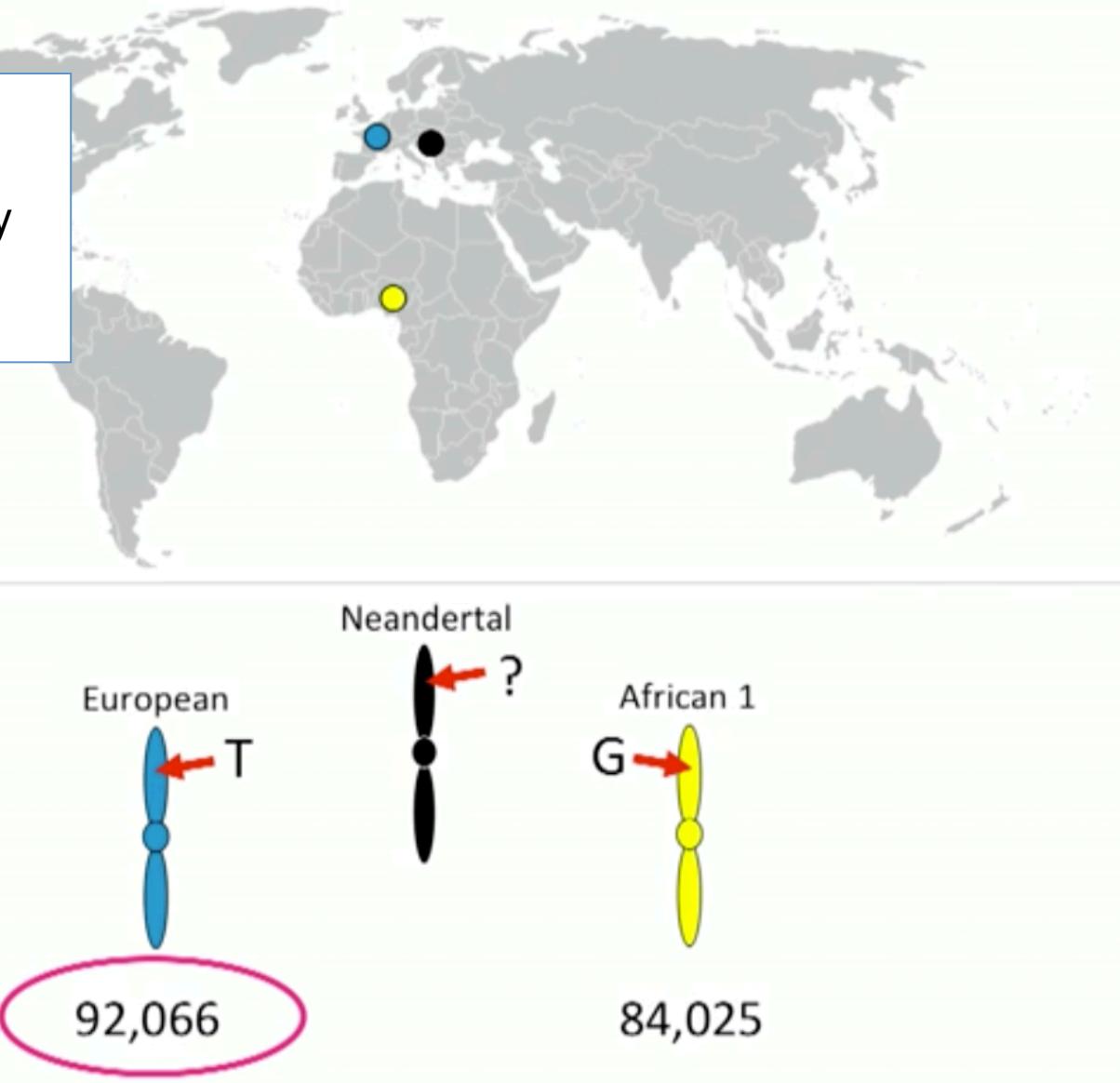
Did we mix?

As far as we know,
Neanderthals were never
in Africa, and do not see
Neanderthal alleles to be
more common in one
African population over
another



Did we mix?

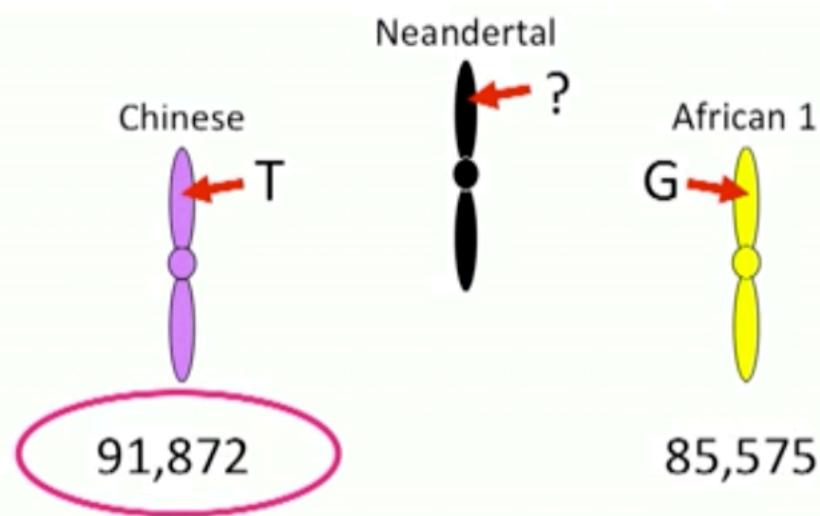
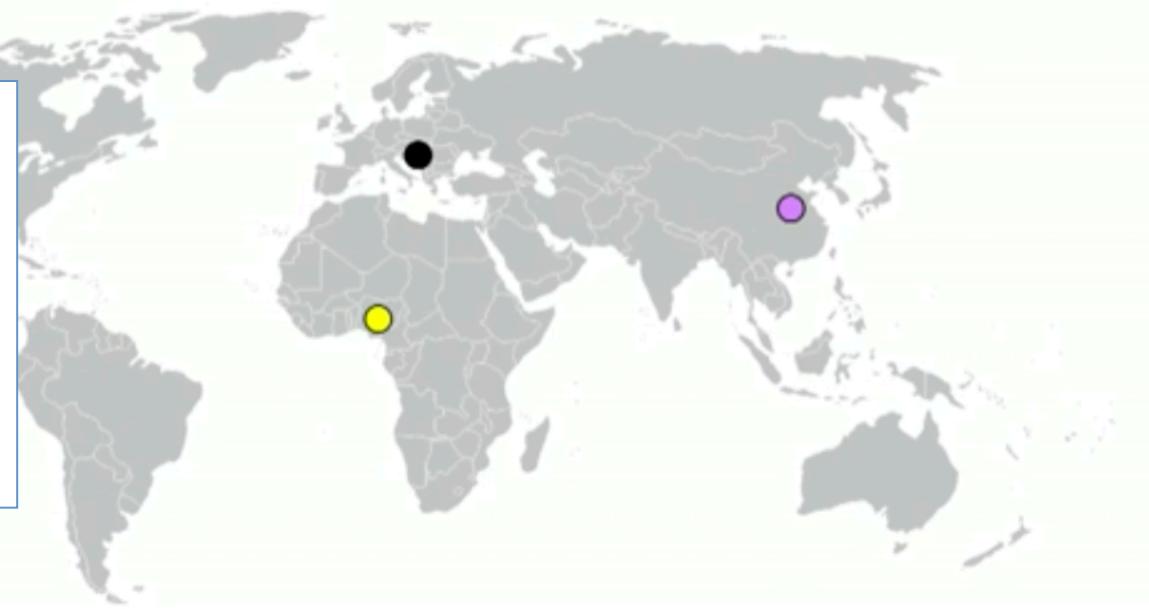
In contrast, we do see Neanderthals match Europeans significantly more frequently than Africans



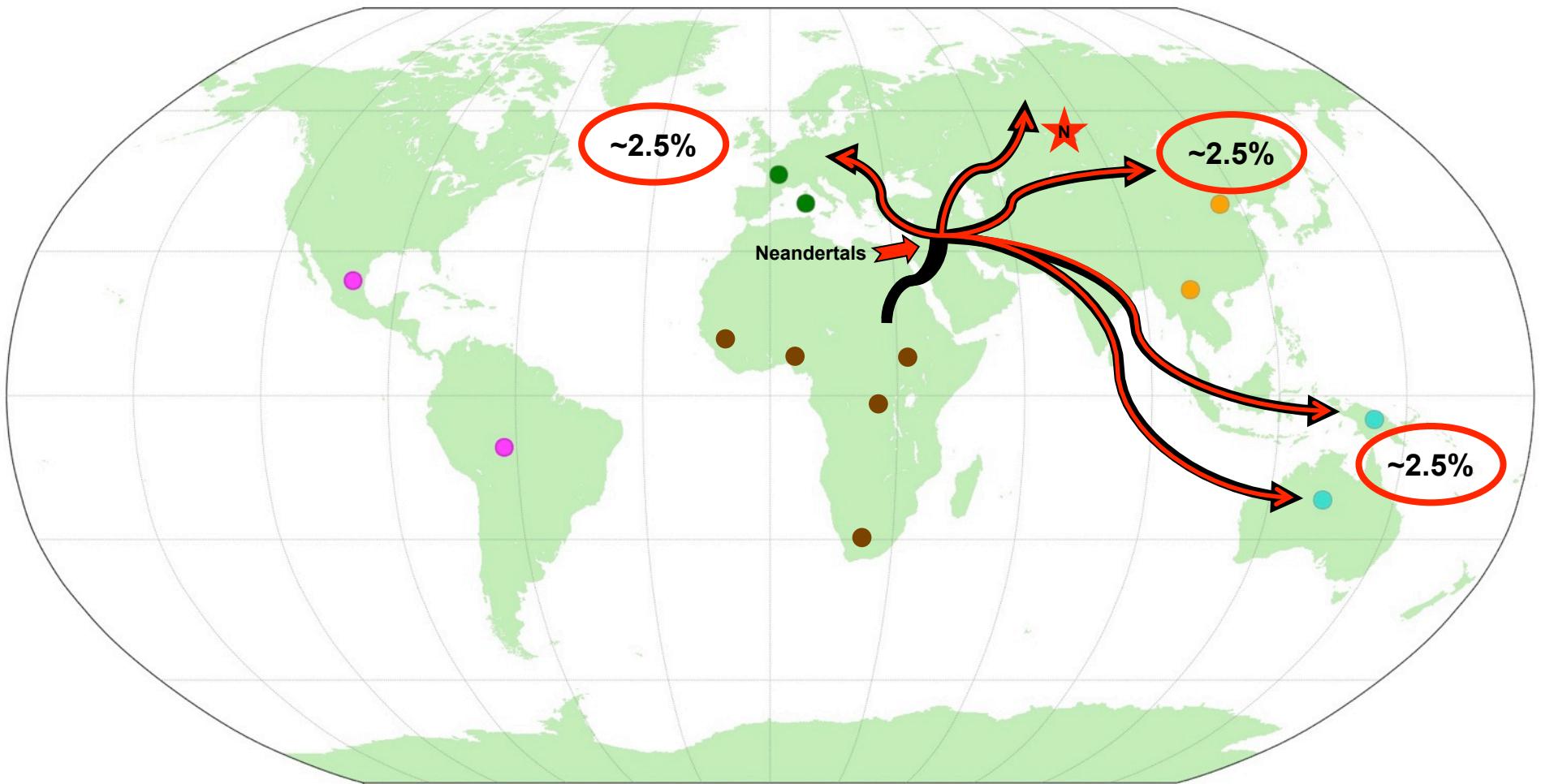
Did we mix?

Also see Neanderthals
match Chinese
significantly more
often...

... but Neanderthals
never lived in China!



Neanderthal Interbreeding

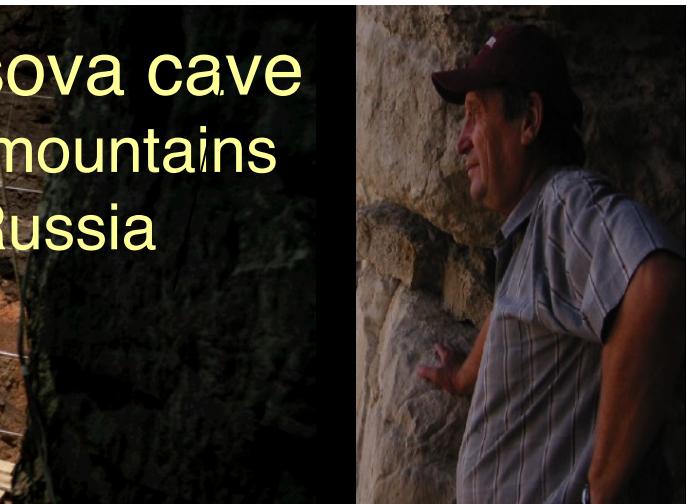


As modern humans migrated out of Africa, they apparently interbred with Neanderthal's so we see their alleles across the rest of the world and carry about 2.5% of their genome with us!

What about other ancient hominids?



Denisova cave
Altai mountains
Russia

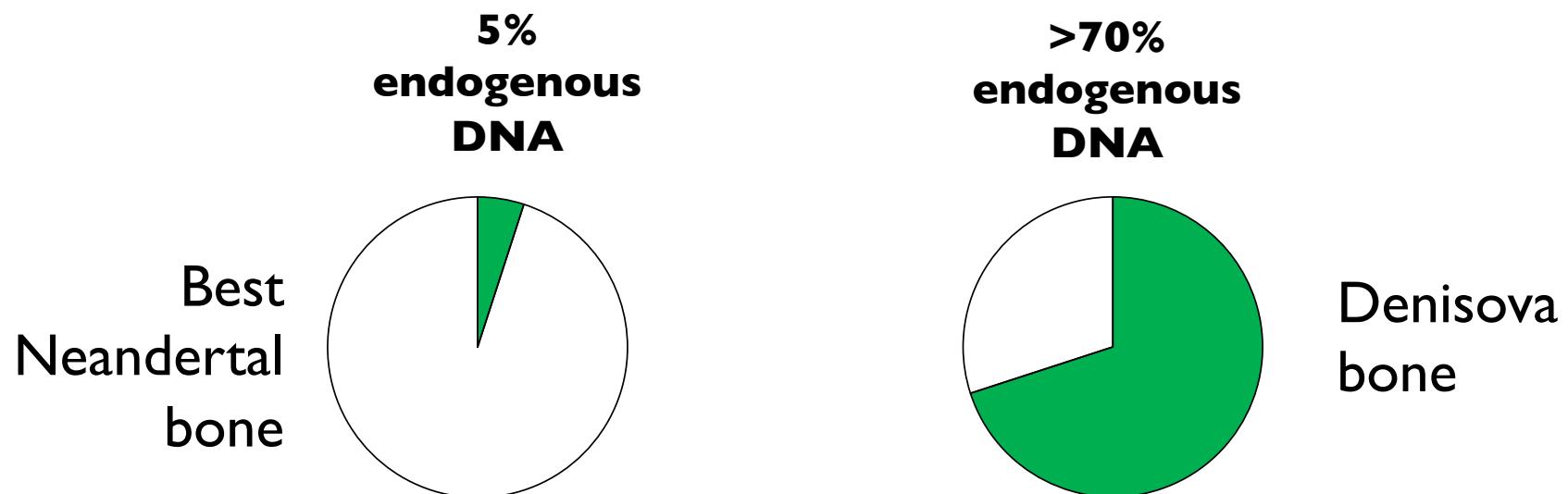
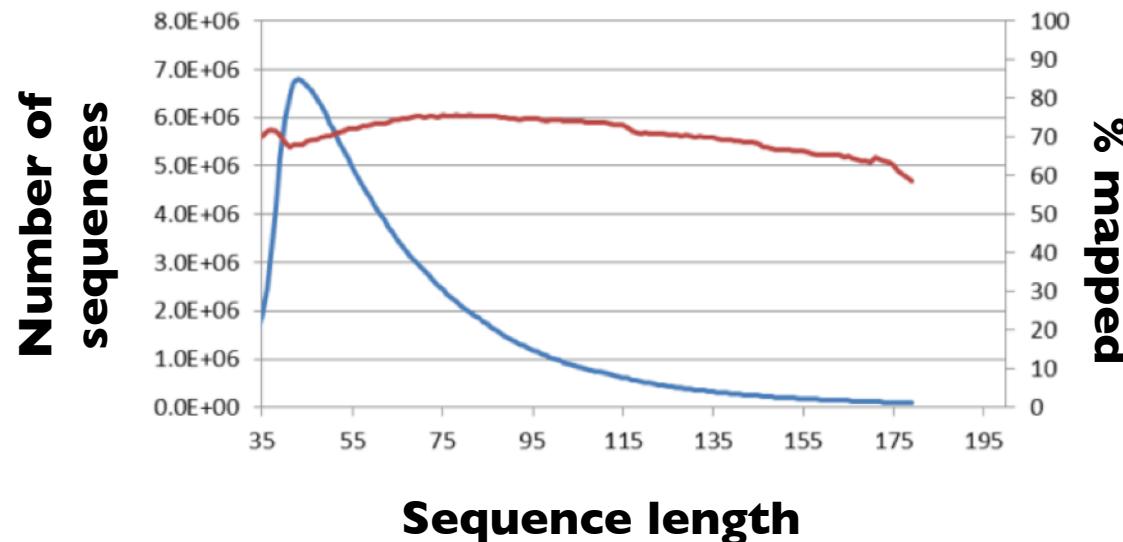


Academician A.P. Derevianko

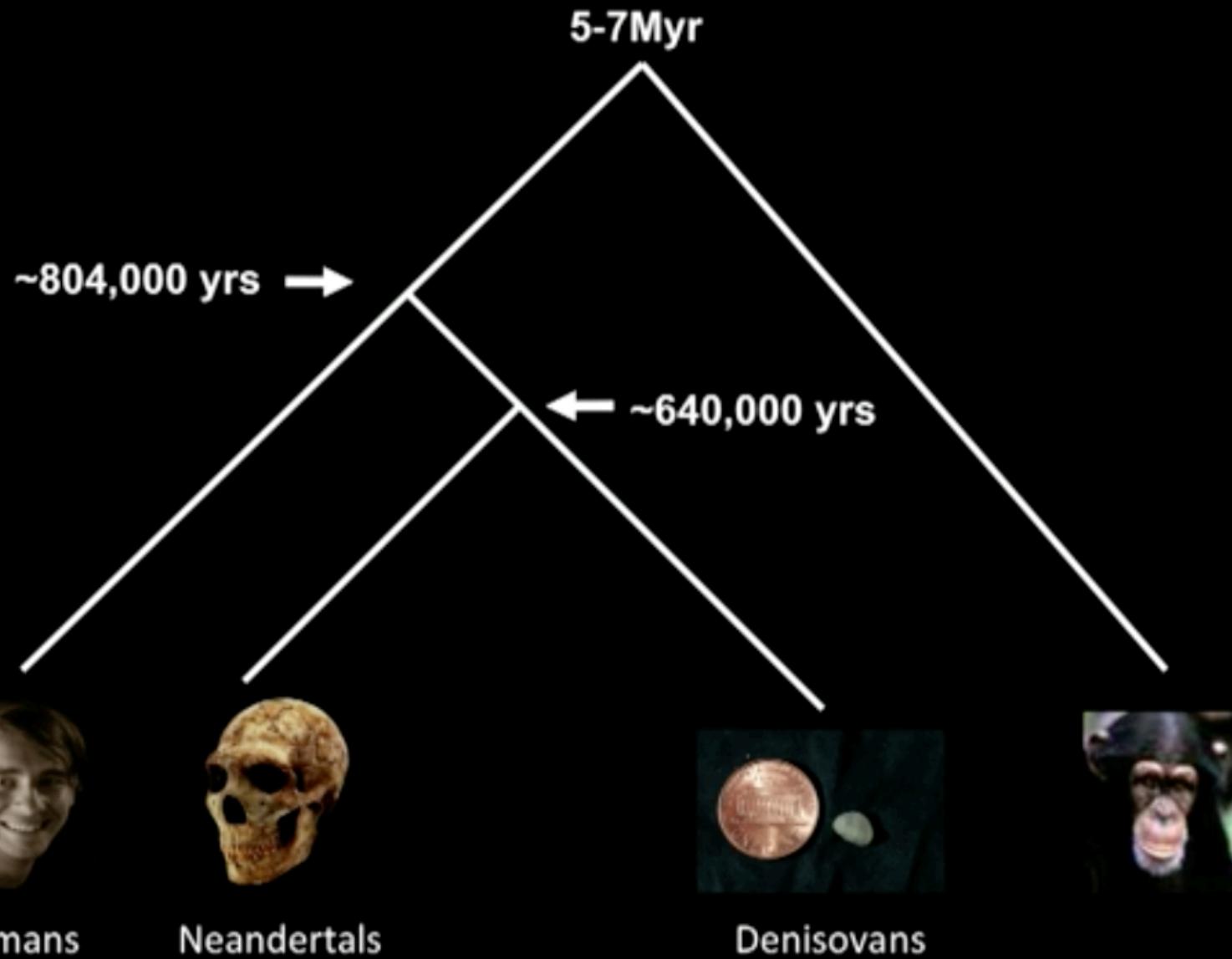




Extraordinary preservation



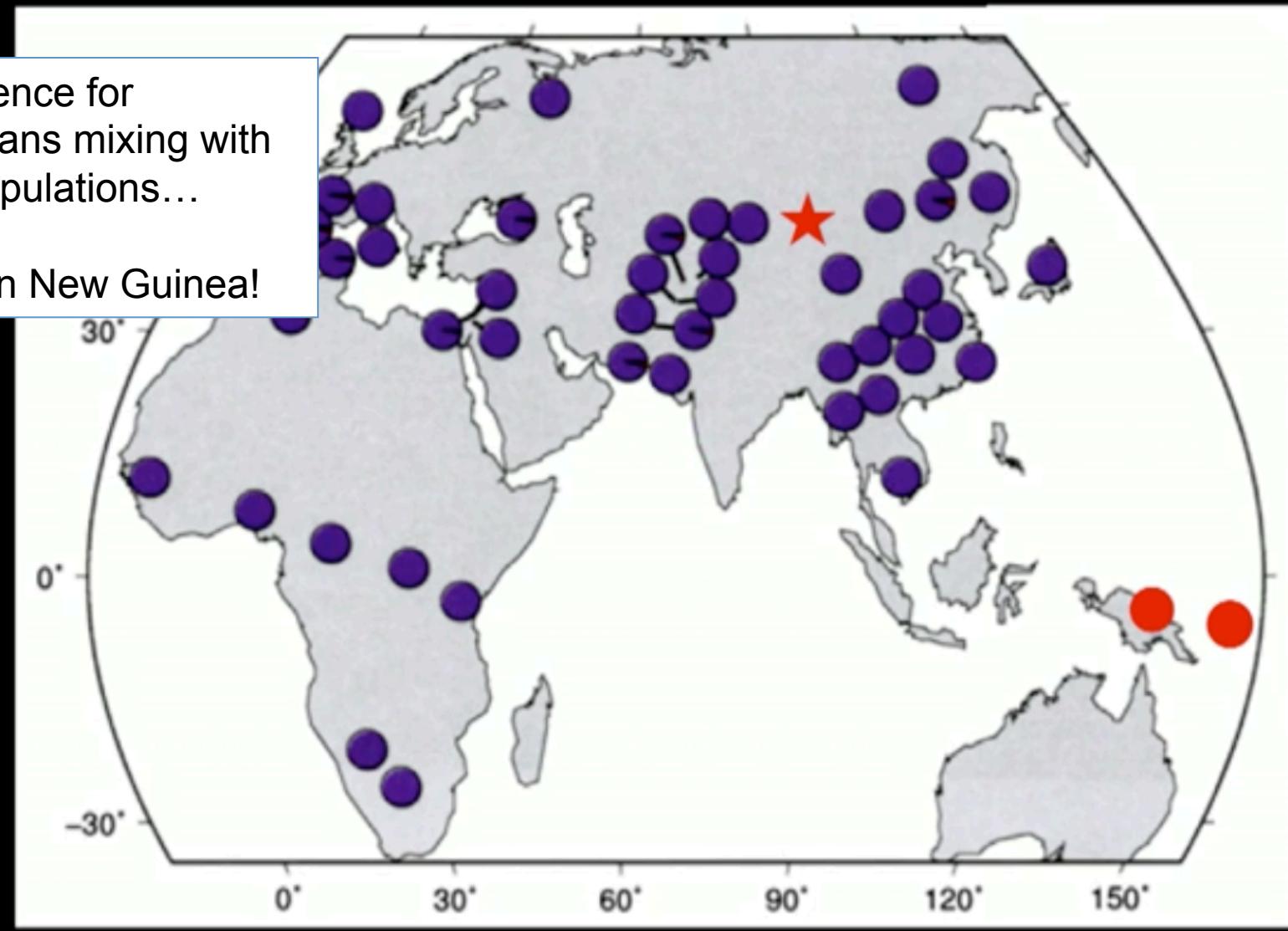
Denisovans & Neandertals



Did we mix?

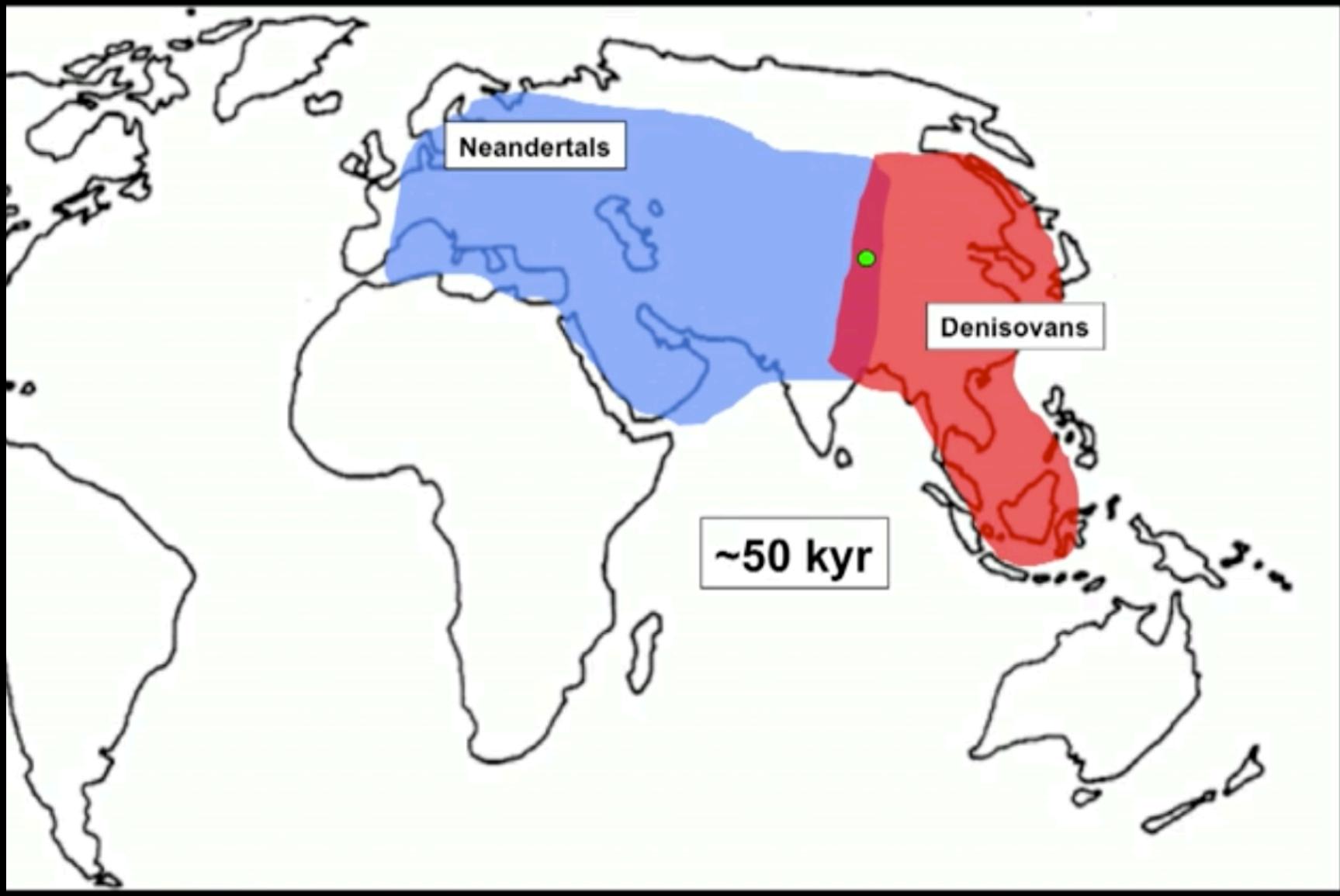
No evidence for
Denisovans mixing with
other populations...

Except in New Guinea!

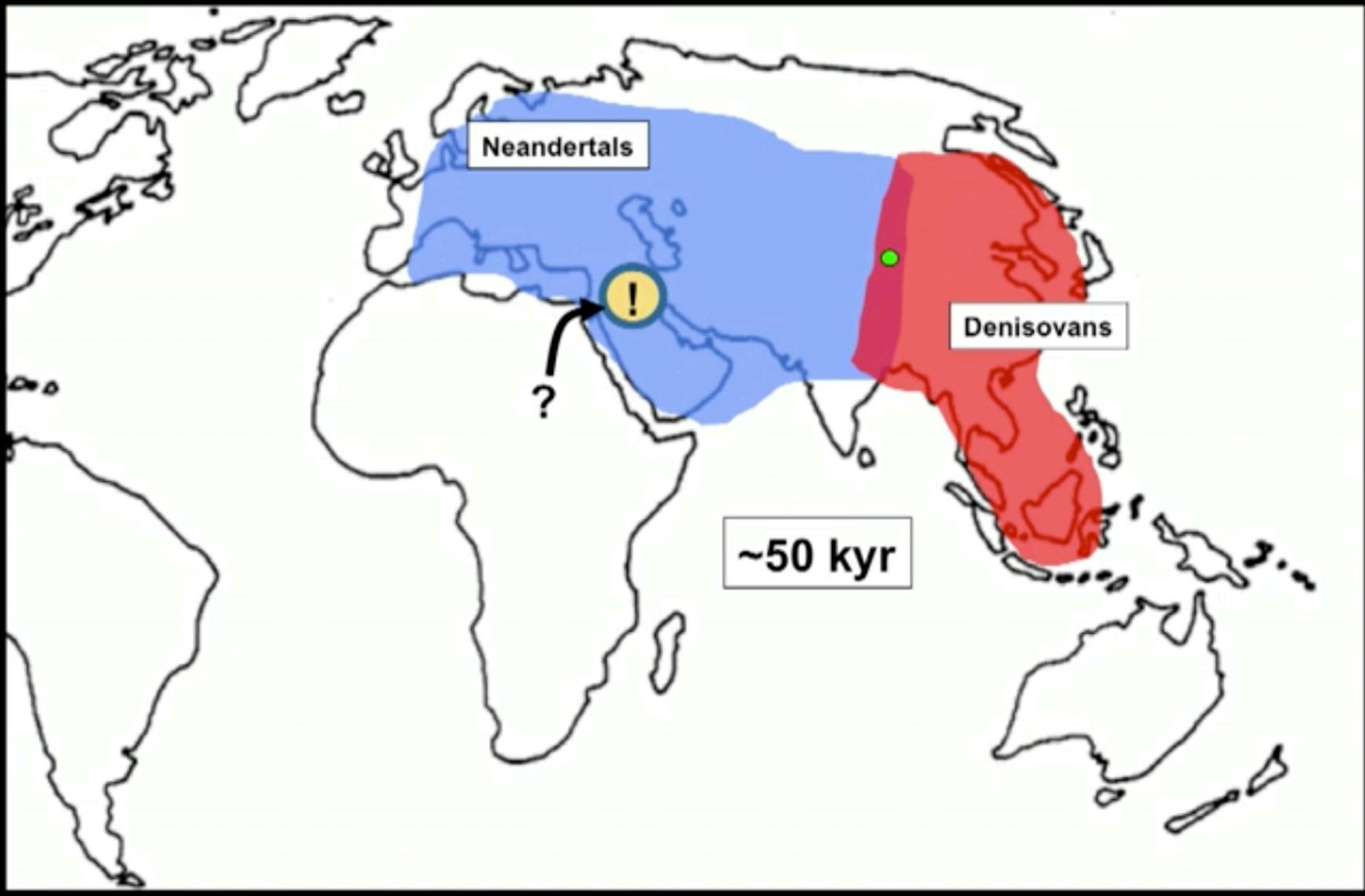


Map after Pickrell et al., 2009

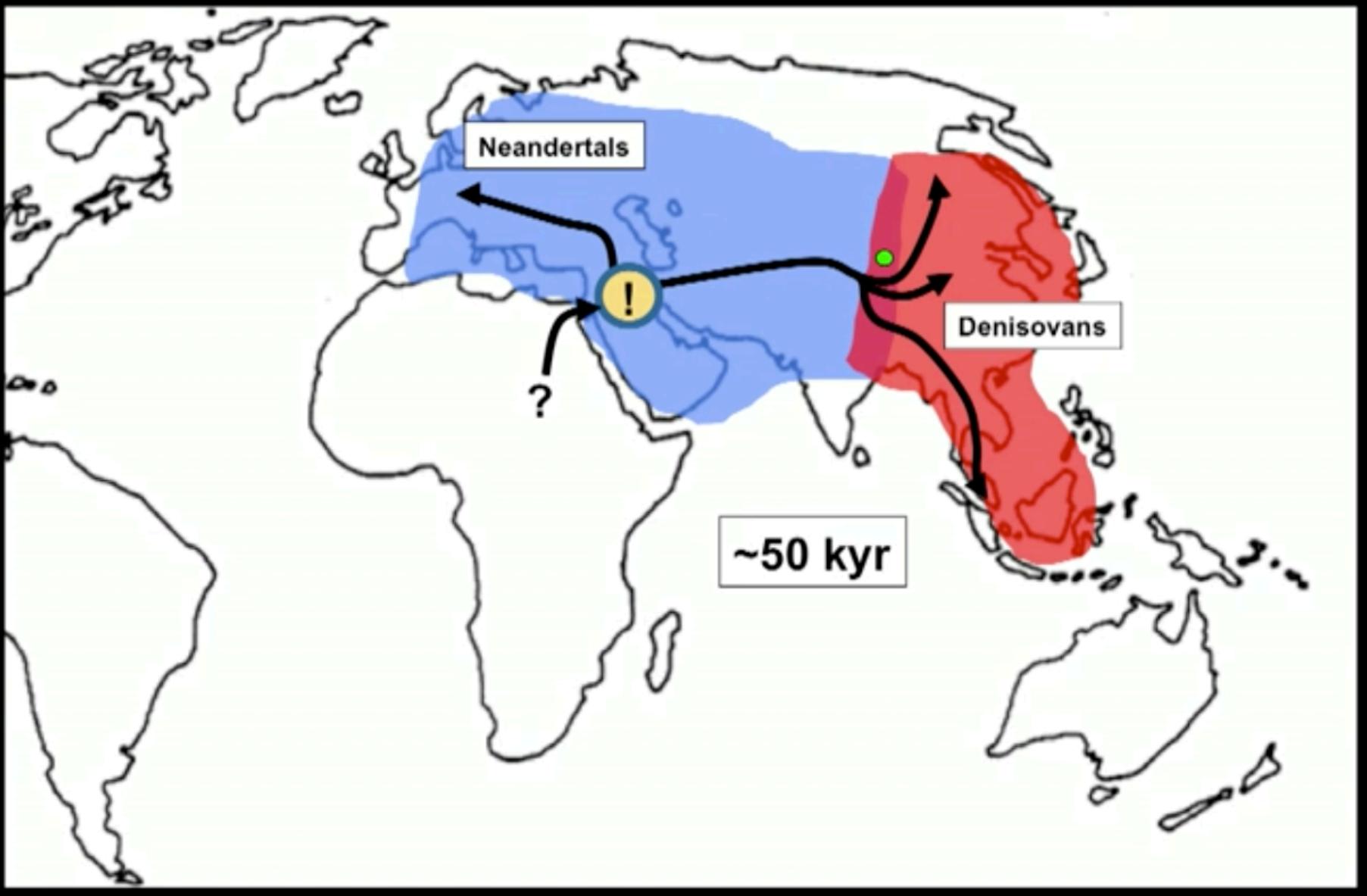
Timeline of ancient hominids



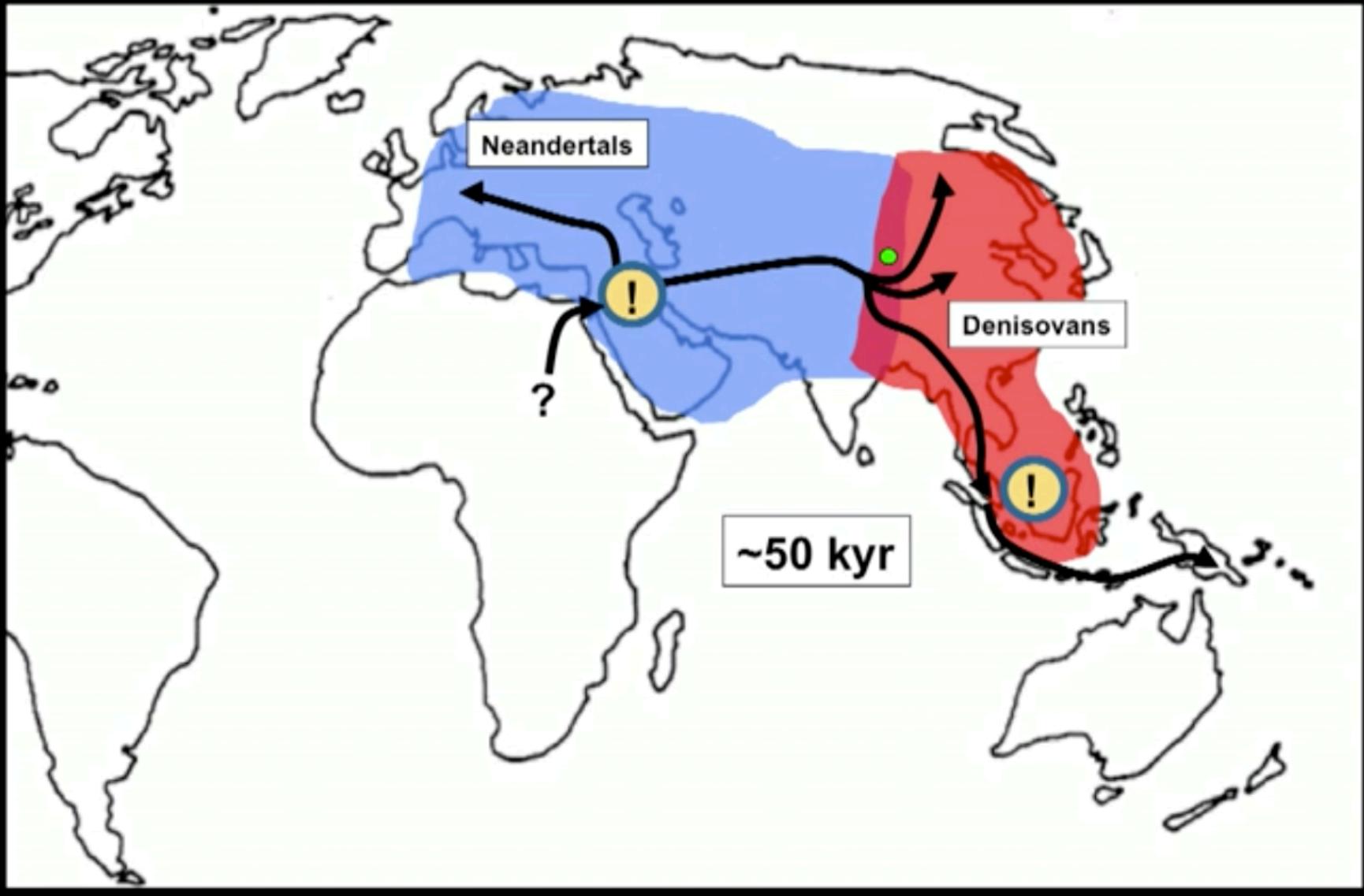
Timeline of ancient hominids



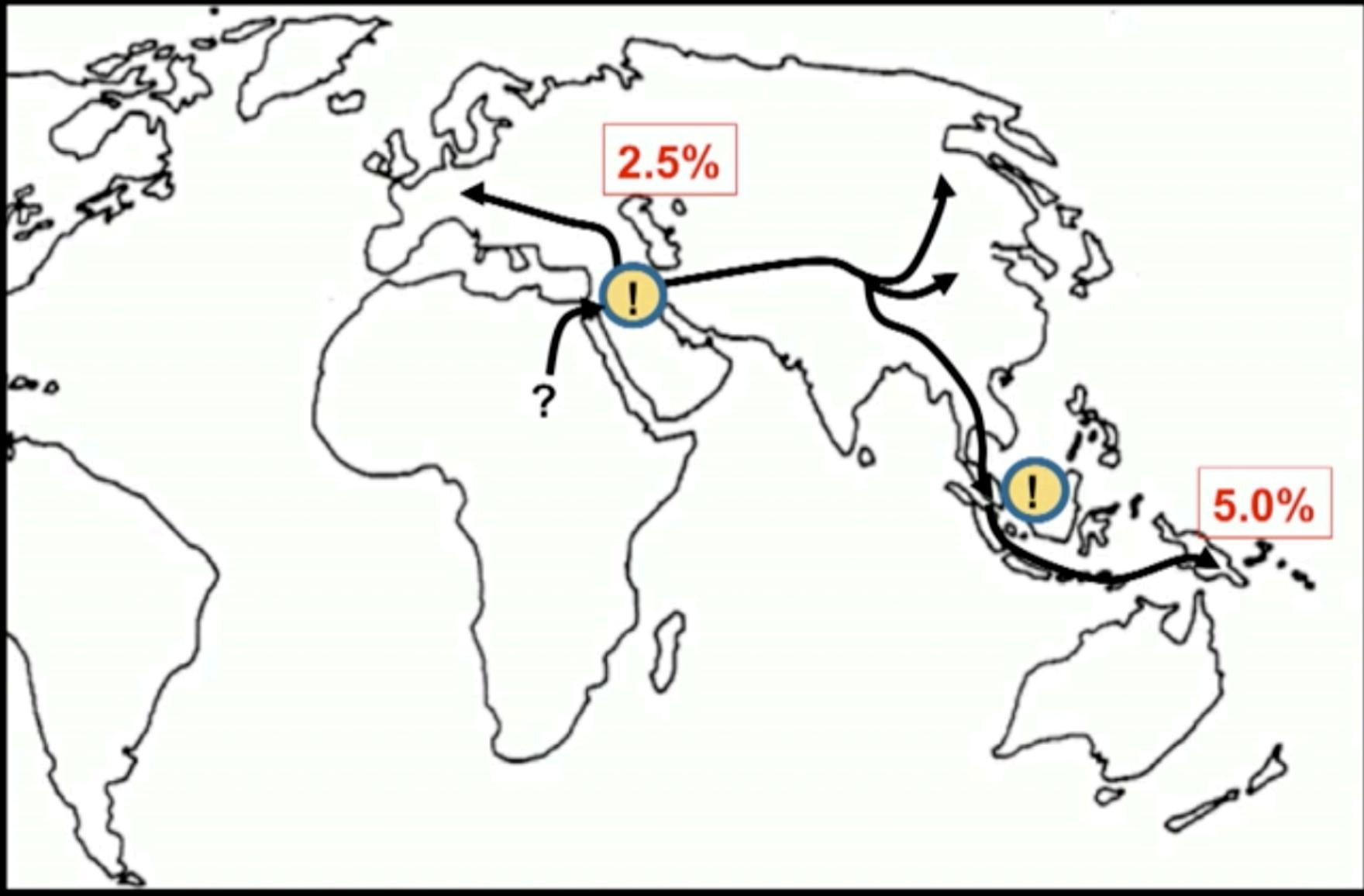
Timeline of ancient hominids



Timeline of ancient hominids



Timeline of ancient hominids



We have always mixed!

Cite as: B. Vernot *et al.*, *Science* 10.1126/science.aad9416 (2016).

Excavating Neandertal and Denisovan DNA from the genomes of Melanesian individuals

Benjamin Vernot,¹ Serena Tucci,^{1,2} Janet Kelso,³ Joshua G. Schraiber,¹ Aaron B. Wolf,¹ Rachel M. Gittelman,¹ Michael Dannemann,³ Steffi Grote,³ Rajiv C. McCoy,¹ Heather Norton,⁴ Laura B. Scheinfeldt,⁵ David A. Merriwether,⁶ George Koki,⁷ Jonathan S. Friedlaender,⁸ Jon Wakefield,⁹ Svante Pääbo,^{2*} Joshua M. Akey^{1*}

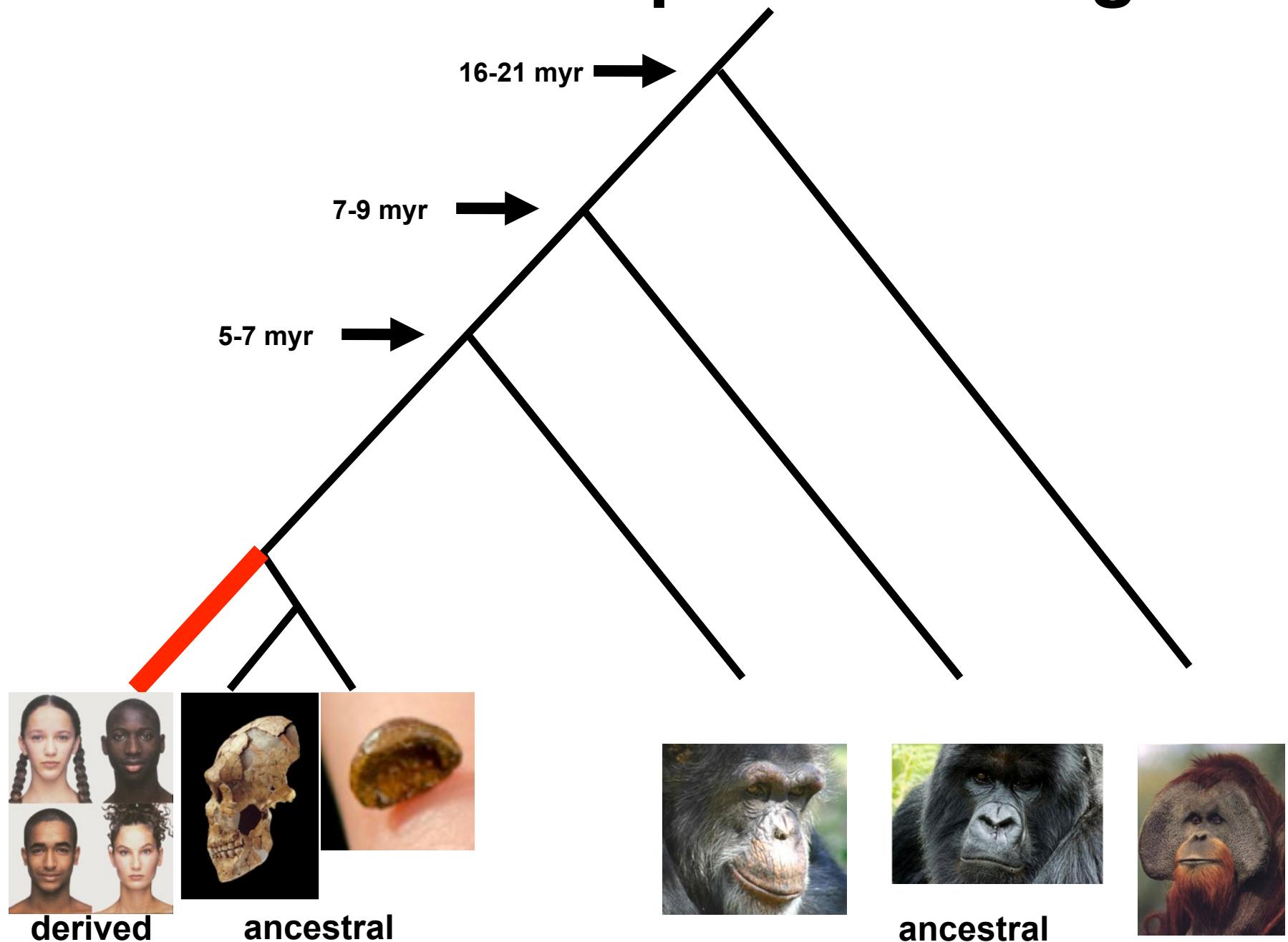
¹Department of Genome Sciences, University of Washington, Seattle, Washington, USA. ²Department of Life Sciences and Biotechnology, University of Ferrara, Italy.

³Department of Evolutionary Genetics, Max-Planck-Institute for Evolutionary Anthropology, Leipzig, Germany. ⁴Department of Anthropology, University of Cincinnati, Cincinnati, OH, USA. ⁵Coriell Institute for Medical Research, Camden, NJ, USA. ⁶Department of Anthropology, Binghamton University, Binghamton, NY, USA. ⁷Institute for Medical Research, Goroka, Eastern Highlands Province, Papua New Guinea. ⁸Department of Anthropology, Temple University, Philadelphia PA, USA. ⁹Department of Statistics, University of Washington, Seattle, Washington, USA.

*Corresponding author. E-mail: paabo@eva.mpg.de (S.P.); akeyj@uw.edu (J.M.A.)

Although Neandertal sequences that persist in the genomes of modern humans have been identified in Eurasians, comparable studies in people whose ancestors hybridized with both Neandertals and Denisovans are lacking. We developed an approach to identify DNA inherited from multiple archaic hominin ancestors and applied it to whole-genome sequences from 1523 geographically diverse individuals, including 35 new Island Melanesian genomes. In aggregate, we recovered 1.34 Gb and 303 Mb of the Neandertal and Denisovan genome, respectively. We leverage these maps of archaic sequence to show that Neandertal admixture occurred multiple times in different non-African populations, characterize genomic regions that are significantly depleted of archaic sequence, and identify signatures of adaptive introgression.

Modern human-specific changes



Recipe for a modern human

109,295 single nucleotide changes (SNCs)
7,944 insertions and deletions

Changes in protein coding genes

277 cause fixed amino acid substitutions
87 affect splice sites

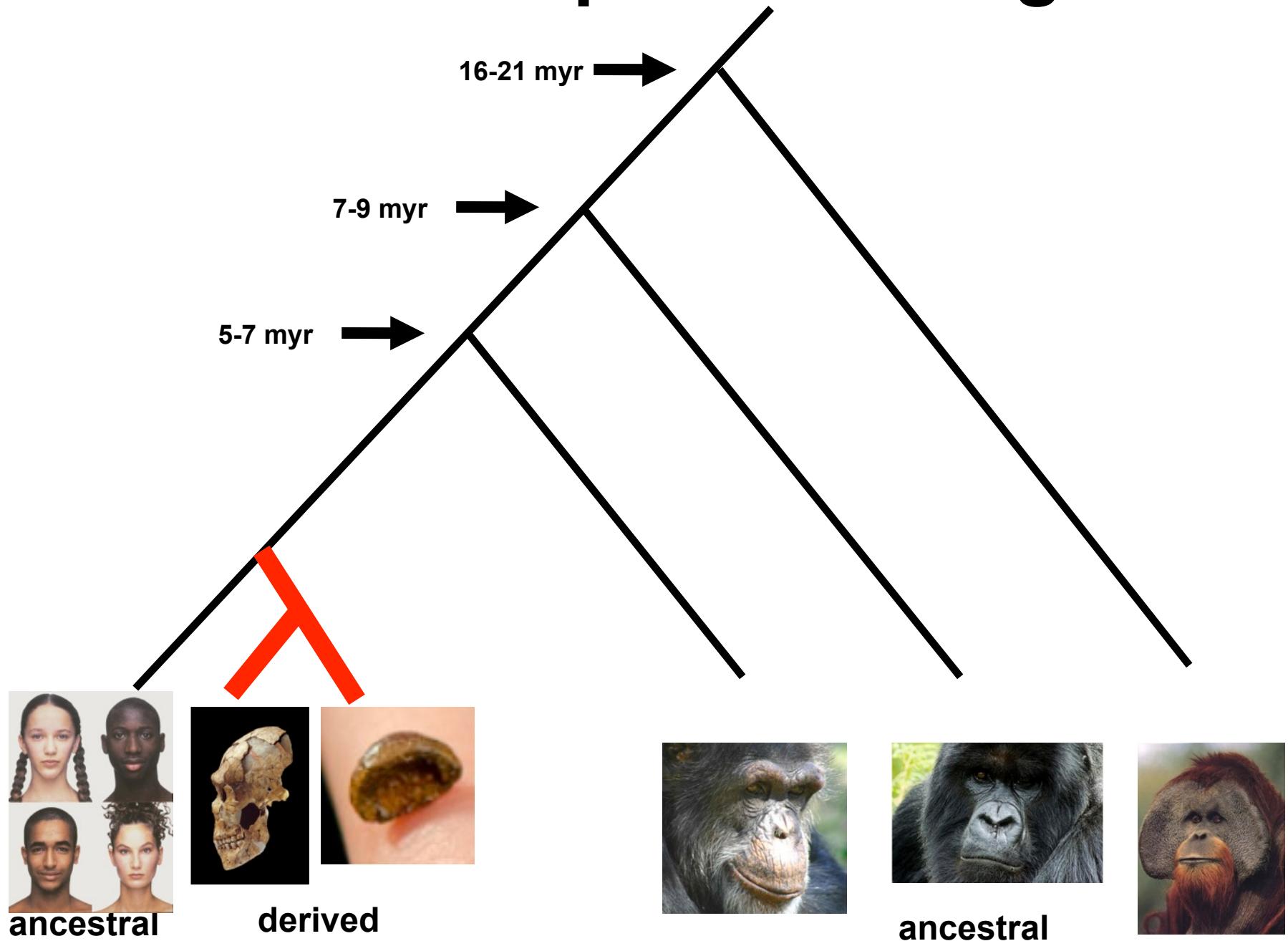
Changes in Non-coding & regulatory sequences

26 affect well-defined motifs inside
regulatory regions

Enrichment analysis

Nonsynonymous	None	- Giant melanosomes in melanocytes (p=6.77e-6; FWER=0.091;
Splice sites	skin pigmentation	
3' UTR	None	<ul style="list-style-type: none"> - 1-3 toe syndactyly (p=1.34288e-05; FWER=0.538; FDR=0.0887928) - 1-5 toe syndactyly (p=1.34288e-05; FWER=0.538; FDR=0.0887928) - Aplasia/Hypoplasia of the distal phalanx of the thumb (p=1.34288e-05; FWER=0.538; FDR=0.0887928) - Bifid or hypoplastic epiglottis (p=1.34288e-05; FWER=0.538; FDR=0.0887928) - Central polydactyly (feet) (p=1.34288e-05; FWER=0.538; FDR=0.0887928)
skeletal morphologies (limb length, digit development)		
<ul style="list-style-type: none"> - Distal urethral duplication (p=1.34288e-05; FWER=0.538; FDR=0.0887928) - Dysplastic distal thumb phalanges with a central hole (p=1.34288e-05; FWER=0.538; FDR=0.0887928) 		
morphologies of the larynx and the epiglottis		
<ul style="list-style-type: none"> - Laryngeal cleft (p=1.34288e-05; FWER=0.538; FDR=0.0887928) - Midline facial capillary hemangioma (p=1.34288e-05; FWER=0.538; FDR=0.0887928) - Preductal coarctation of the aorta (p=1.34288e-05; FWER=0.538; FDR=0.0887928) - Radial head subluxation (p=1.34288e-05; FWER=0.538; FDR=0.0887928) - Short distal phalanx of the thumb (p=1.34288e-05; FWER=0.538; FDR=0.0887928) 		

Neandertal-specific changes



Enrichment analysis

Nonsynonymous	None	- Abnormality of the thumb (p=3.01e-5; FWER=0.025; FDR=0.02) - Aplasia/Hypoplasia of the thumb (p=6.31e-5; FWER=0.054; FDR=0.024) - Facial cleft (p=0.0004; FWER=0.36; FDR=0.098) - Wide pubic symphysis (p=0.0004; FWER=0.36; FDR=0.098) - Abnormality of the frontal hairline (p=0.00042; FWER=0.39; FDR=0.096)
Skeletal and hair morphology		
		- Abnormality of the scalp (p=0.00042; FWER=0.42; FDR=0.08) - Abnormality of the finger (p=0.0005; FWER=0.44; FDR=0.08) - Brachydactyly syndrome (p=0.00062; FWER=0.48; FDR=0.088)

Protein	Ensembl ID	Protein position	Ancestral amino acid	Derived amino acid	Description
ABCA12	ENSP00000272895	199	W	C	ATP-binding cassette, sub-family A (ABC1)
FRAS1	ENSP00000264895	209	P	S	Fraser syndrome 1
GLI3	ENSP00000379258	1537	R	C	GLI family zinc finger 3
LAMB3	ENSP00000355997	926	A	D	Laminin, beta 3
MOGS	ENSP00000233616	495	R	Q	Mannosyl-oligosaccharide glucosidase

FOXP2 Analysis

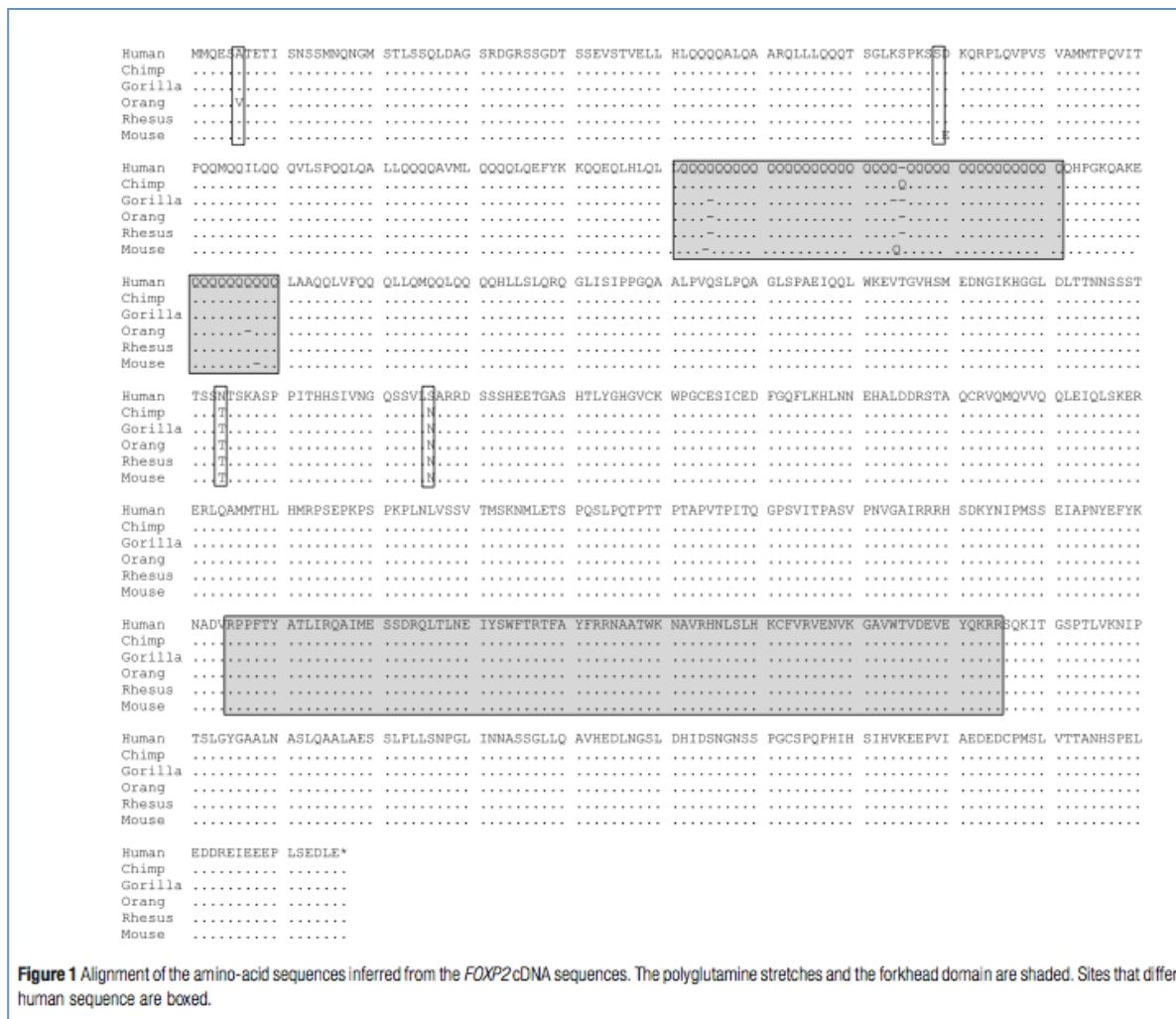


Figure 1 Alignment of the amino-acid sequences inferred from the *FOXP2* cDNA sequences. The polyglutamine stretches and the forkhead domain are shaded. Sites that differ from the human sequence are boxed.

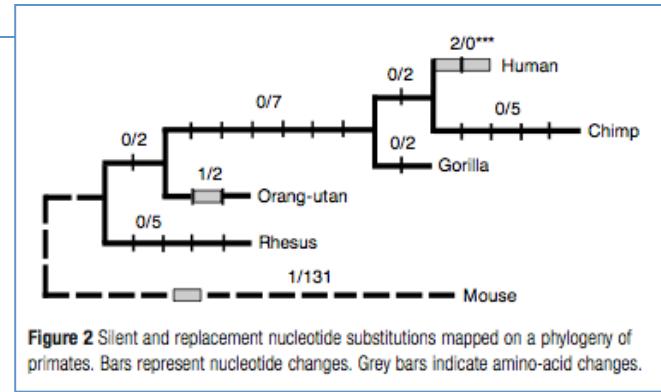


Figure 2 Silent and replacement nucleotide substitutions mapped on a phylogeny of primates. Bars represent nucleotide changes. Grey bars indicate amino-acid changes.

- Mutations of FOXP2 cause a severe speech and language disorder in people
 - Versions of FOXP2 exist in similar forms in distantly related vertebrates; functional studies of the gene in mice and in songbirds indicate that it is important for modulating plasticity of neural circuits.
 - Outside the brain FOXP2 has also been implicated in development of other tissues such as the lung and gut.

What makes us human?

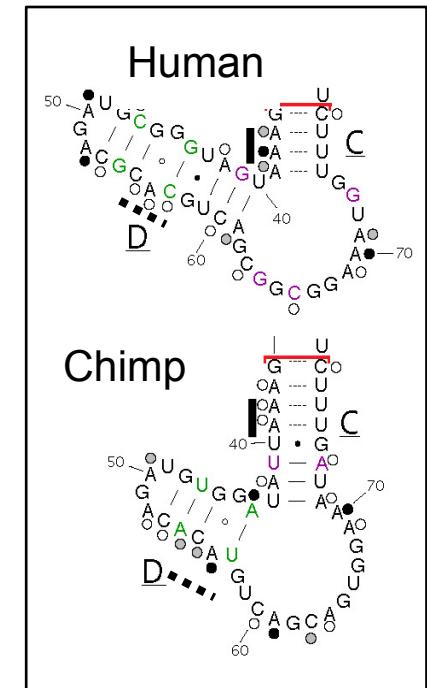
“Human Accelerated Regions”



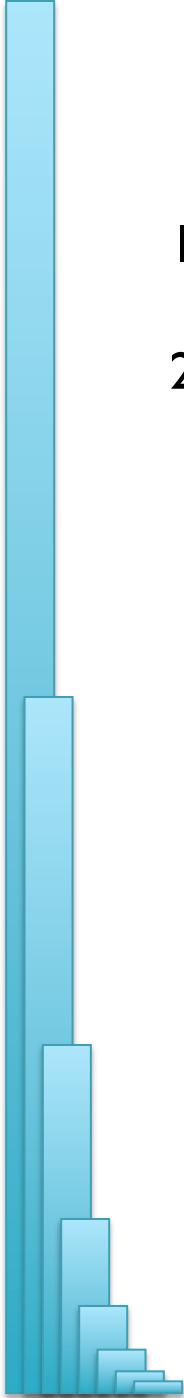
A sequence alignment of six species: human, chimp, dog, mouse, rat, and chicken. The sequences are shown as horizontal lines of DNA bases (A, T, C, G). Red vertical bars highlight specific positions where the human sequence differs from the others, indicating 'accelerated regions'. The human sequence is: TCGATGGTAGACCCACGTCAGCGGGAAATGGTTCTATCAAATCAAAGTCTTAGAGATTTCCCTCAAGTTCAATGA. The other species' sequences are very similar, with minor variations.

Systematic scan of recent human evolution identified the gene *HAR1F* as the most dramatic “human accelerated region”.

Follow up analysis found it was specifically expressed in Cajal-Retzius neurons in the human brain from 6 to 19 gestational weeks.



(Pollard et al., *Nature*, 2006)



Next Steps

1. Questions on project?
2. Check out the course webpage



Welcome to Applied Comparative Genomics

<https://github.com/schatzlab/appliedgenomics>

Questions?