

Lecture 18. Ancient and Modern Humans

Michael Schatz

April 8, 2019

JHU 600.749: Applied Comparative Genomics



Preliminary Project Report

Assignment Date: April 8, 2019

Due Date: Monday, April 15, 2017 @ 11:59pm

Each team should submit a PDF of your preliminary project proposal (2 to 3 pages) to GradeScope by 11:59pm on Monday April 15.

The preliminary report should have at least:

- Title of your project
- List of team members and email addresses
- 1 paragraph abstract summarizing the project
- 1+ paragraph of Introduction
- 1+ paragraph of Methods that you are using
- 1+ paragraph of Results, describing the data evaluated and any any preliminary results
- 1+ paragraph of Dicsussion (what you have seen or expect to see)
- 1+ figure showing a preliminary result
- 5+ References to relevant papers and data

The preliminary report should use the Bioinformatics style template. Word and LaTeX templates are available at https://academic.oup.com/bioinformatics/pages/submission_online

Later, you will present your project in class starting the week of April 24. You will also submit your final written report (5-7 pages) of your project by May 15

Please use Piazza if you have any general questions!

601.749 Computational Genomics: Applied Comparative Genomics Midterm Exam

Michael C. Schatz
mschatz@cs.jhu.edu

April 3, 2019
Time: 75 Minutes

Start here: Please fill in the following important information using a permanent pen before you do anything else! Your exam will not be graded if you use a pencil or erasable ink on this page.

Name (print): _____

Email (print): _____

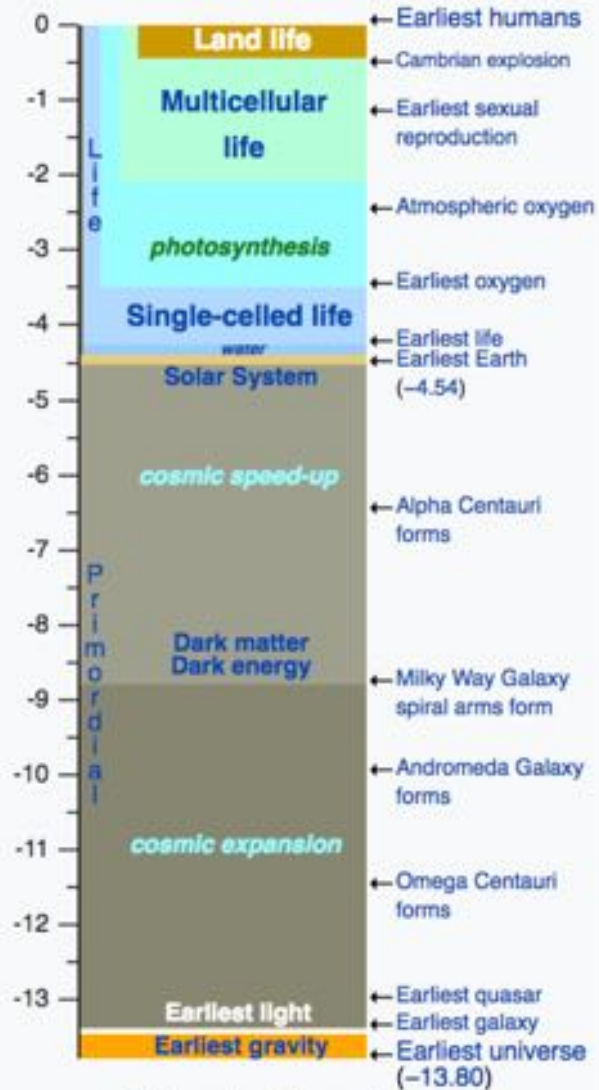


Part I: Ancient Hominds

Our Origins

Nature timeline

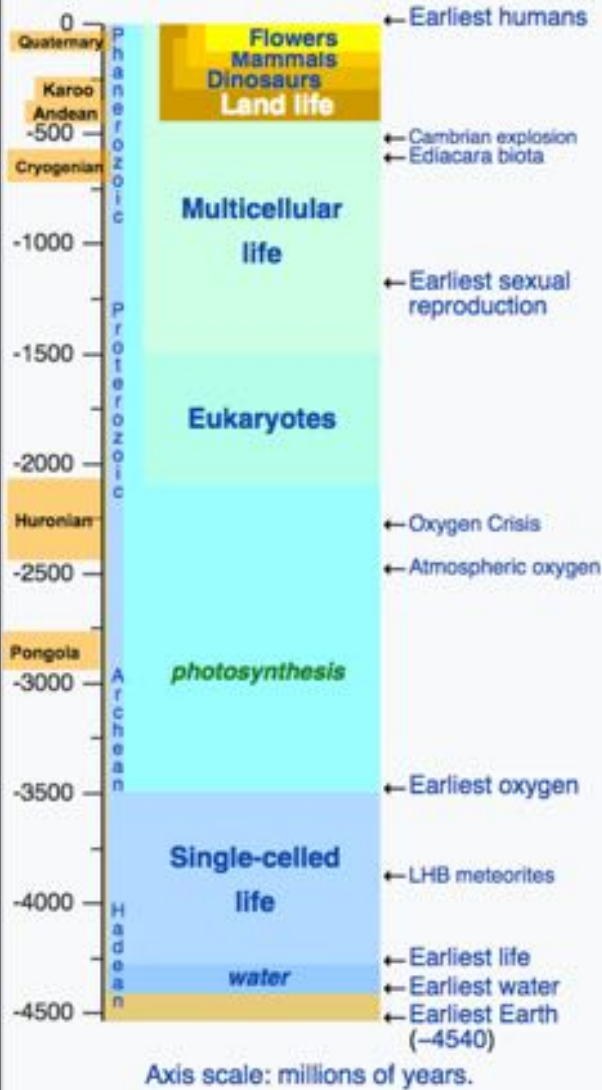
[view](#) • [discuss](#) • [edit](#)



Also see: [Human timeline](#) and [Life timeline](#)

Life timeline

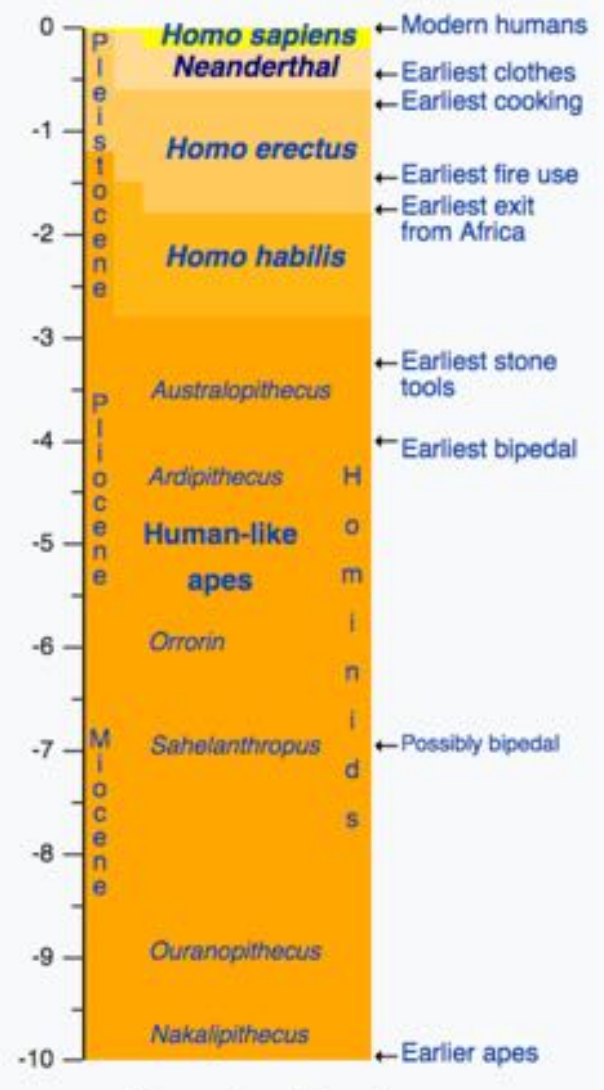
[view](#) • [discuss](#) • [edit](#)



Orange labels: known ice ages.
Also see: [Human timeline](#) and [Nature timeline](#)

Human timeline

[view](#) • [discuss](#) • [edit](#)



Also see: [Life timeline](#) and [Nature timeline](#)

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

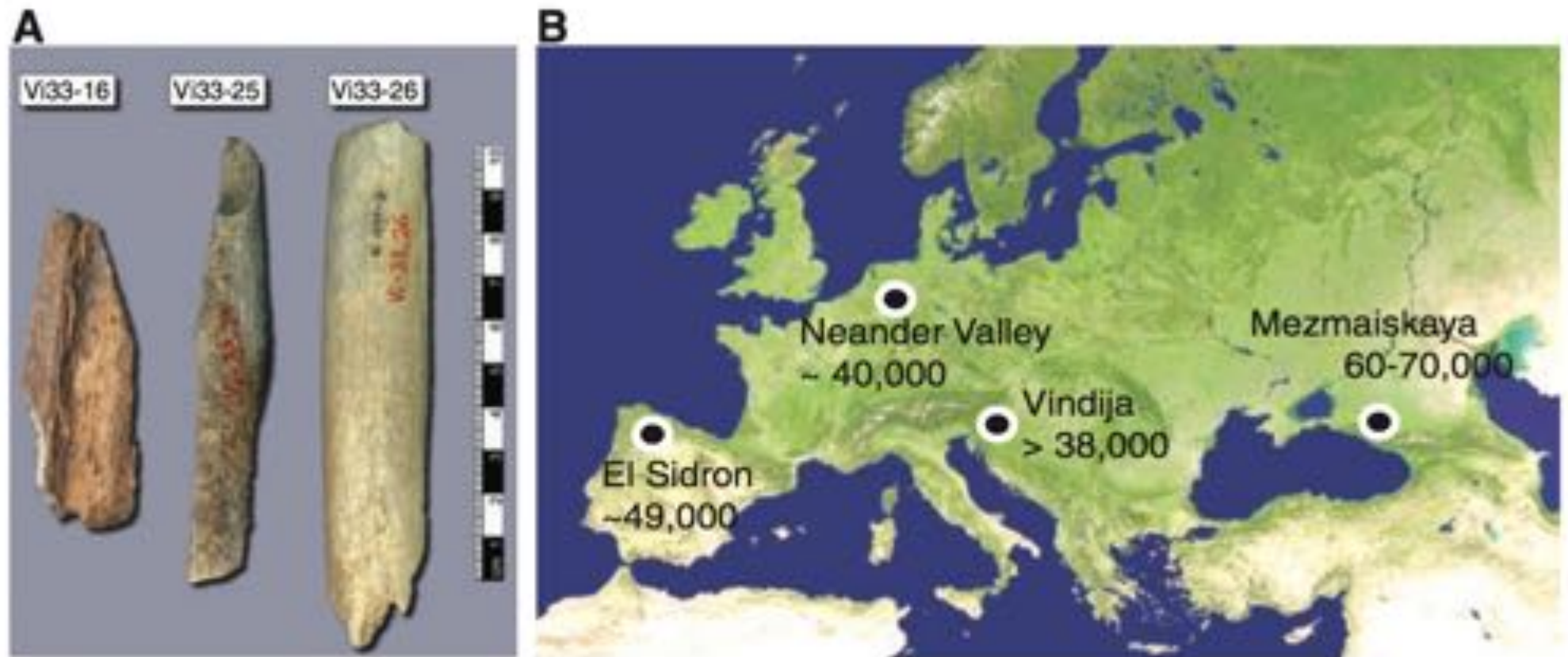
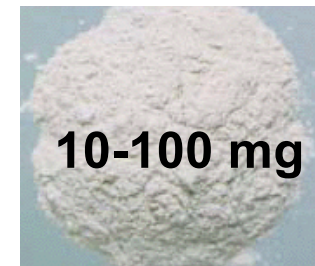
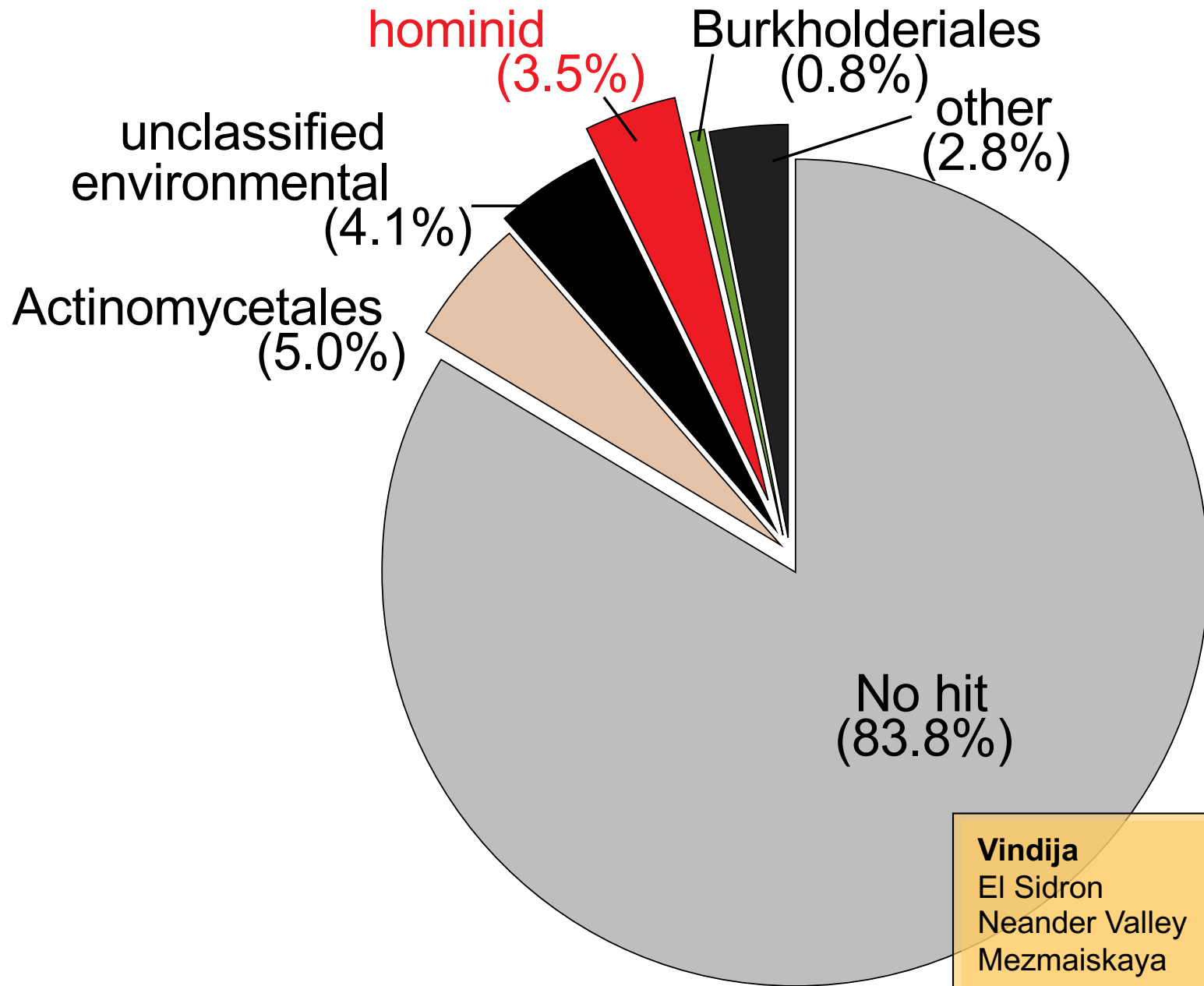


Fig. 1. Samples and sites from which DNA was retrieved. (A) The three bones from Vindija from which Neandertal DNA was sequenced. (B) Map showing the four archaeological sites from which bones were used and their approximate dates (years B.P.).

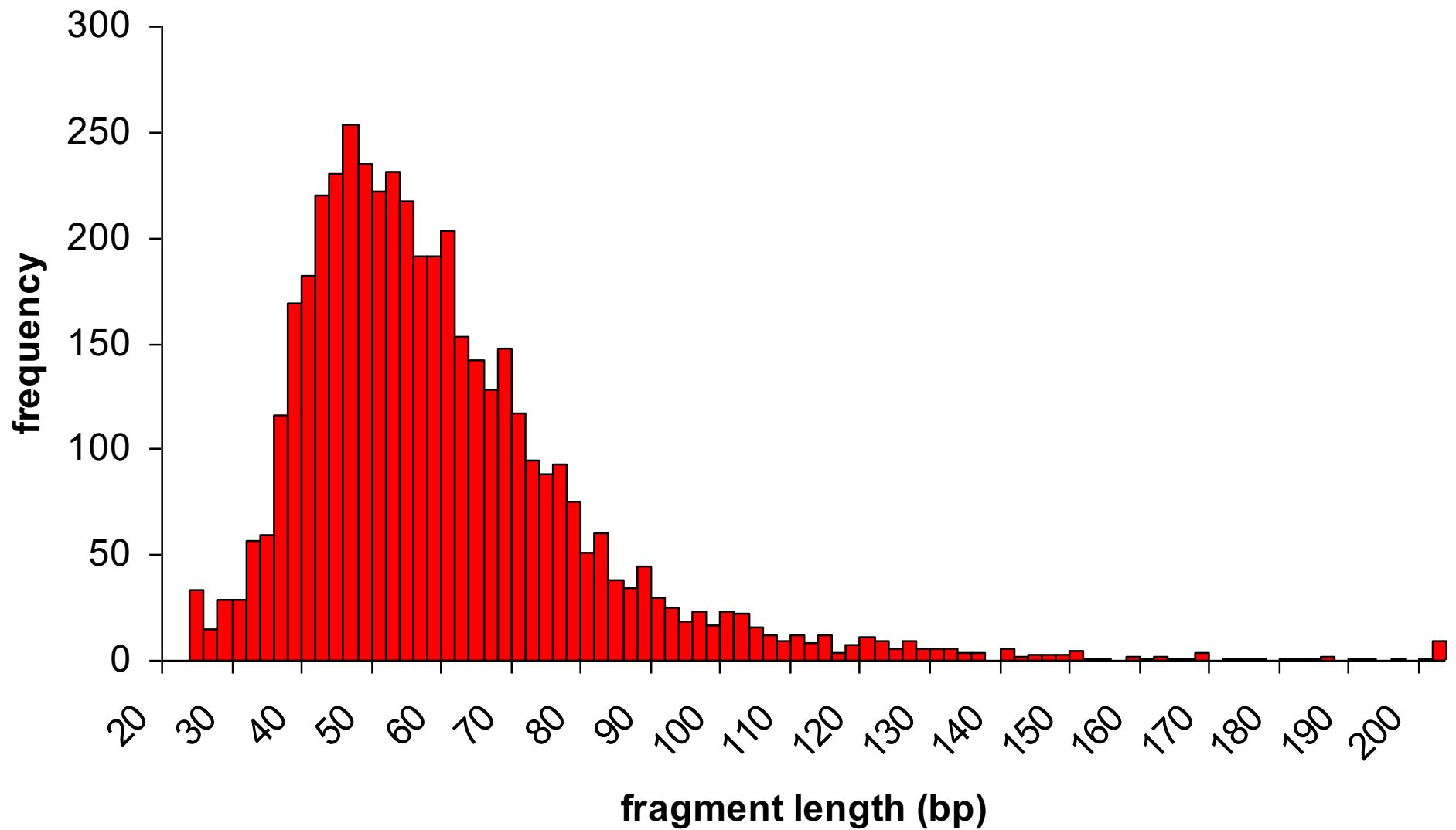
Extracting Ancient DNA



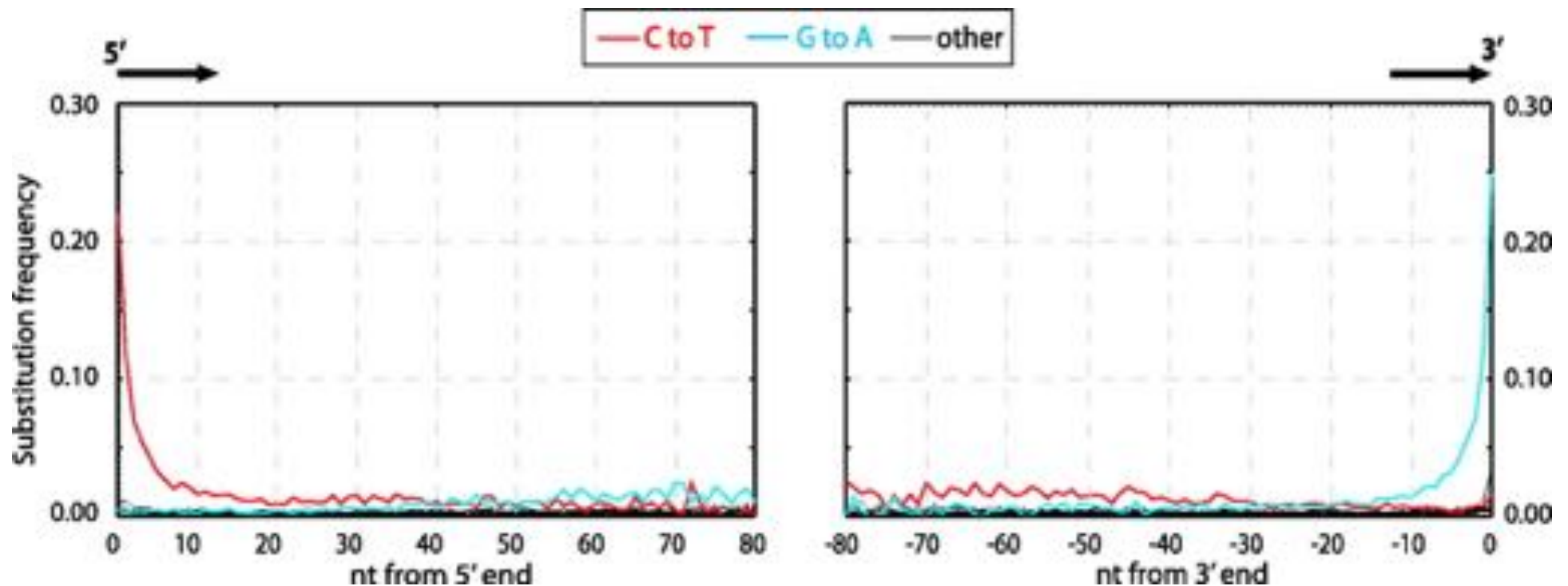
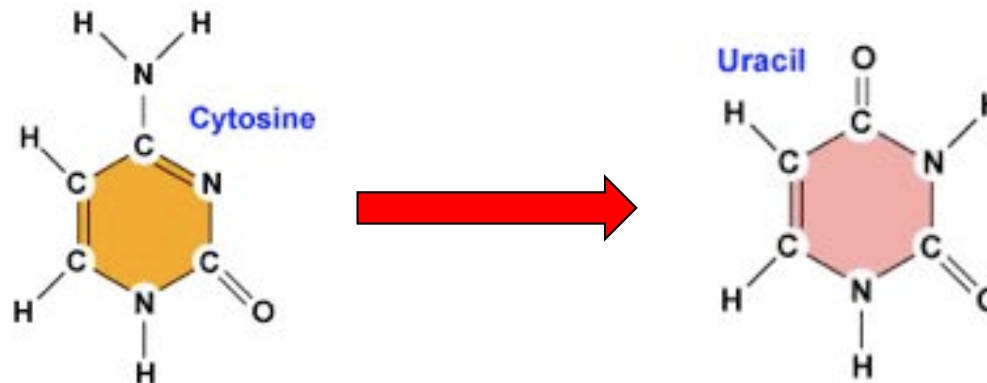
DNA is from mixed sources

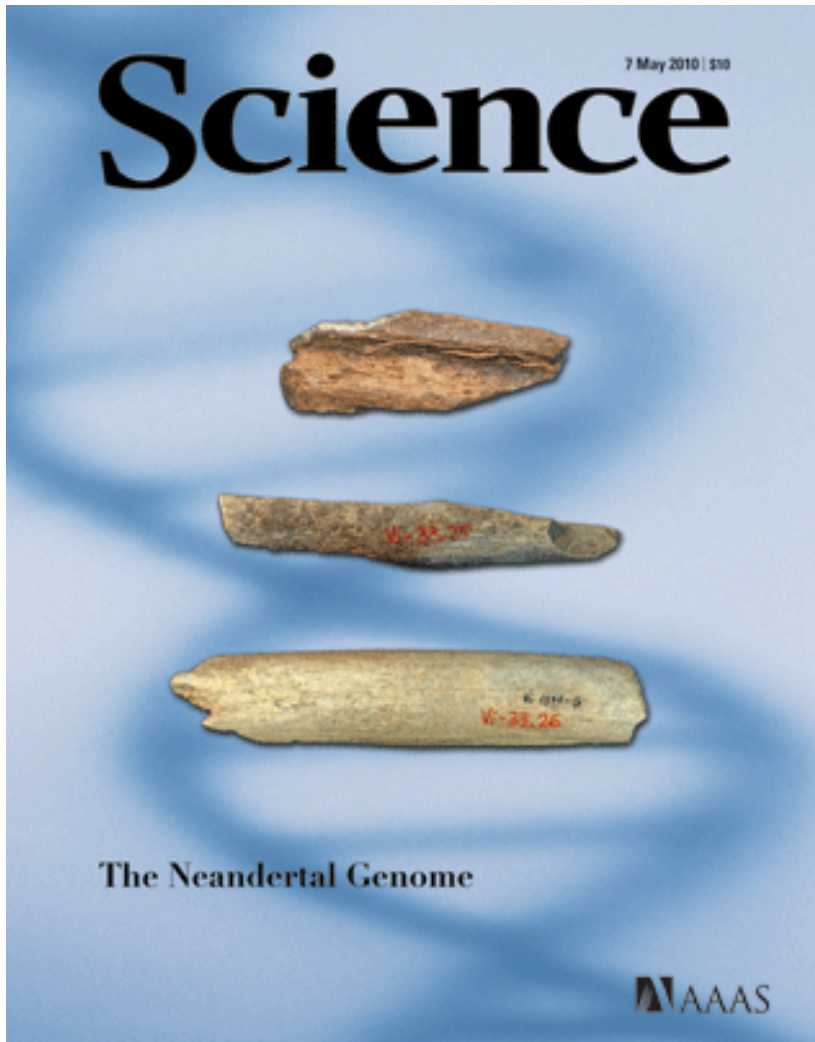


DNA is degraded



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

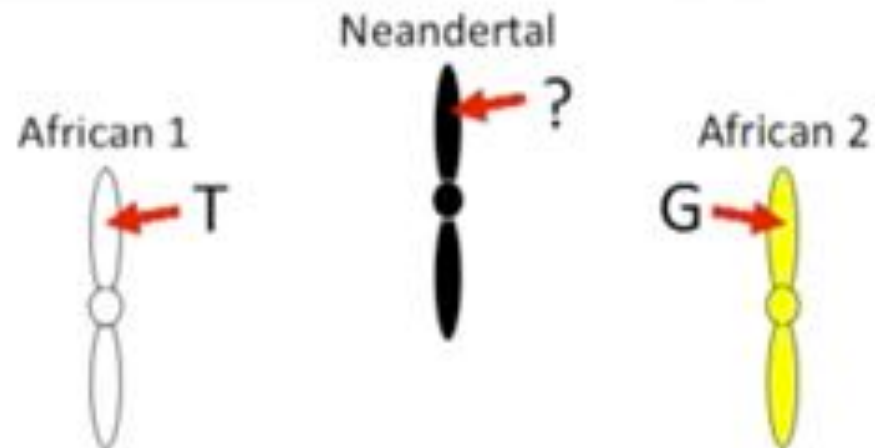
Feldhofer 1 ~2.2 Mb

Mezmaiskaya 1 ~56.4 Mb

~35 Illumina flow cells

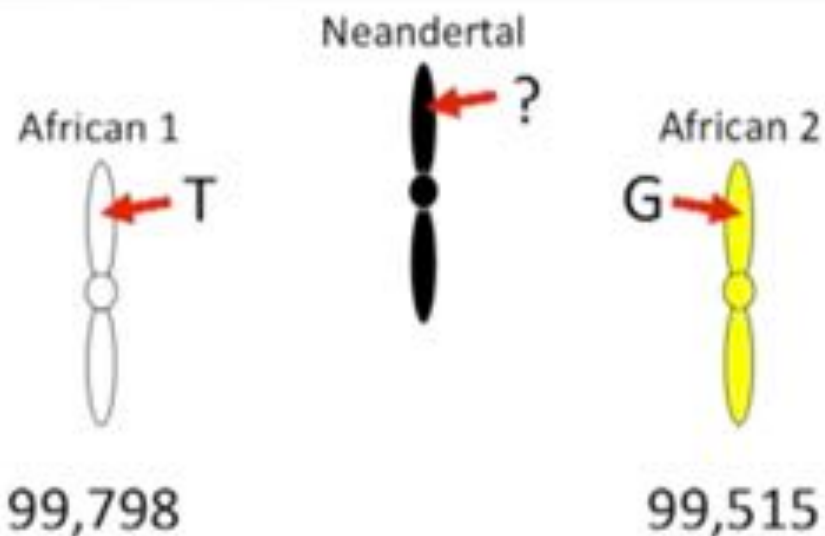
Genome coverage ~1.3 X

Did we mix?



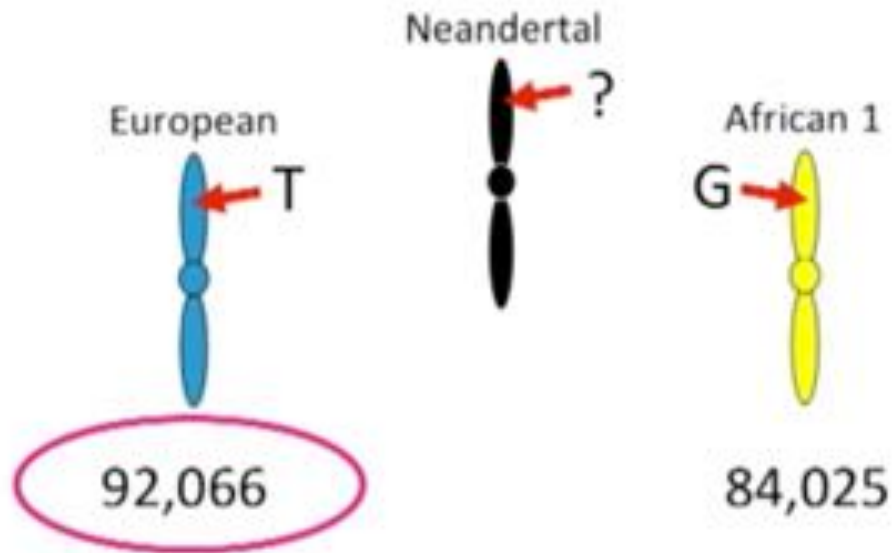
Did we mix?

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?

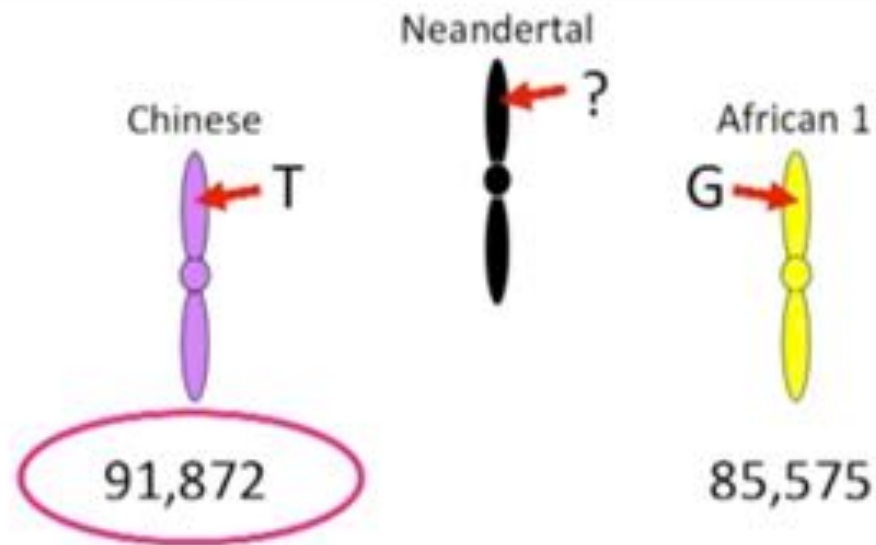
In contrast, we do see Neanderthals match Europeans significantly more frequently than Africans



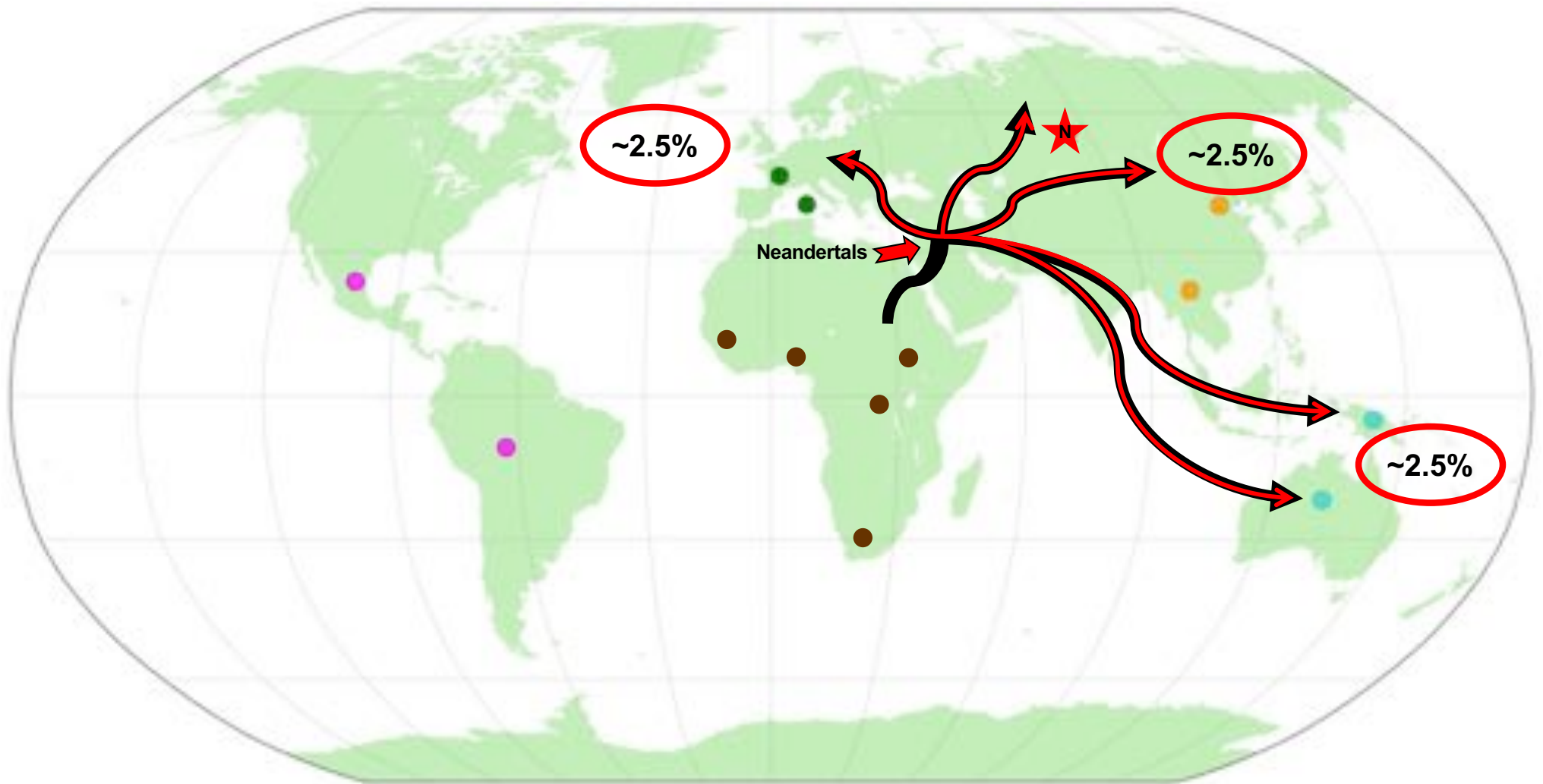
Did we mix?

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 Neanderthals 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

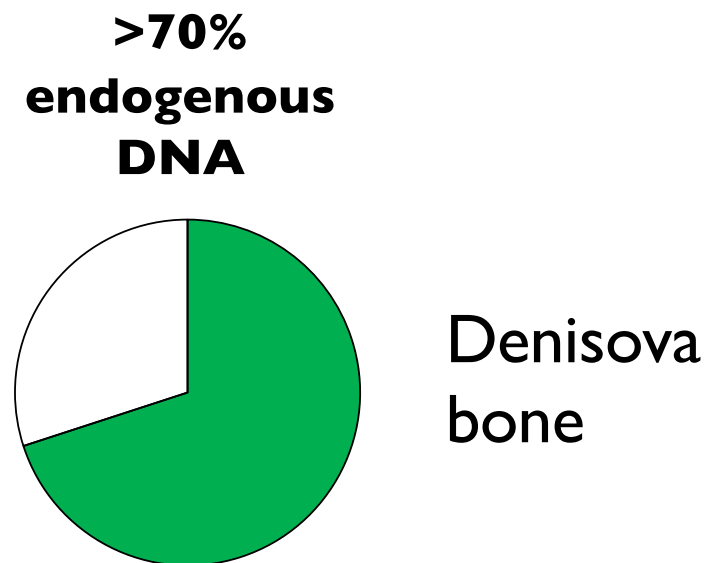
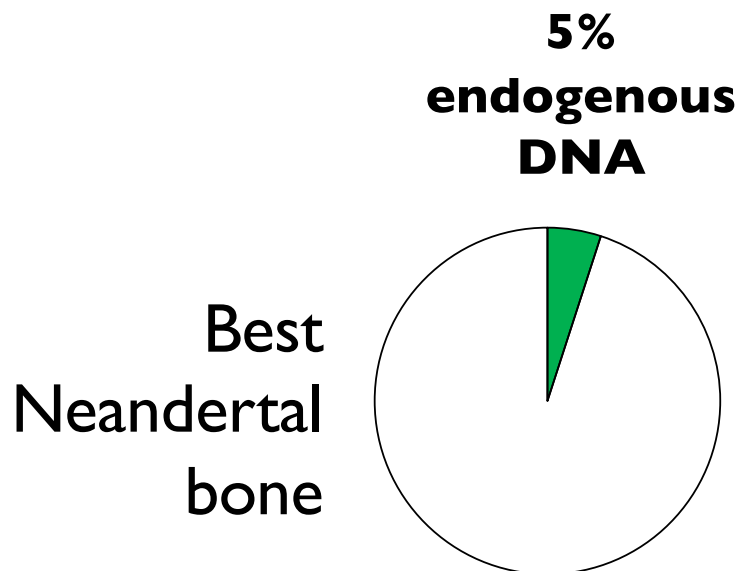
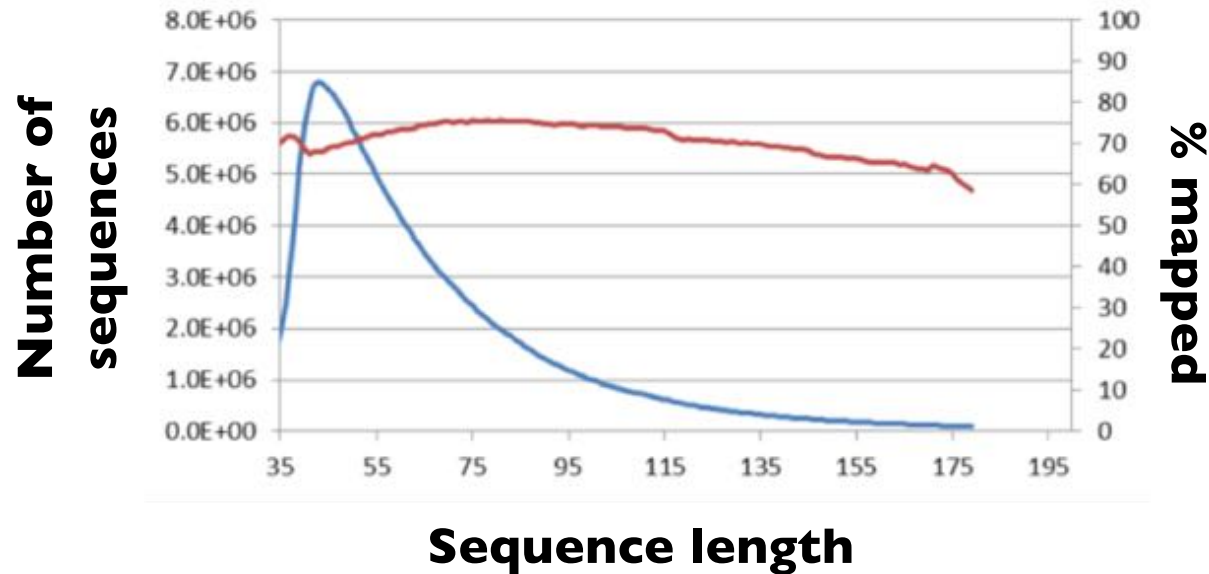


Academician A.P. Derevianko

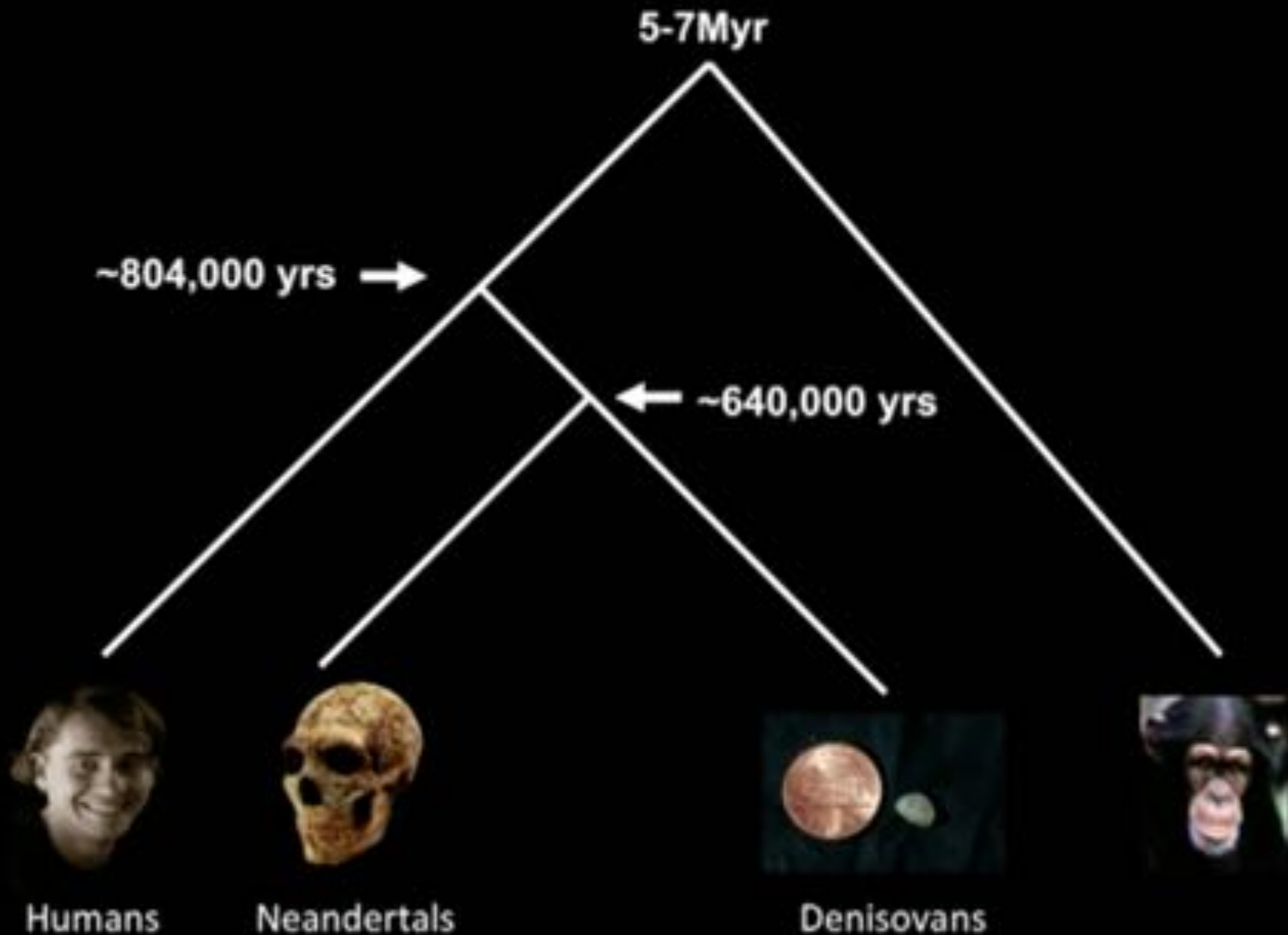




Extraordinary preservation



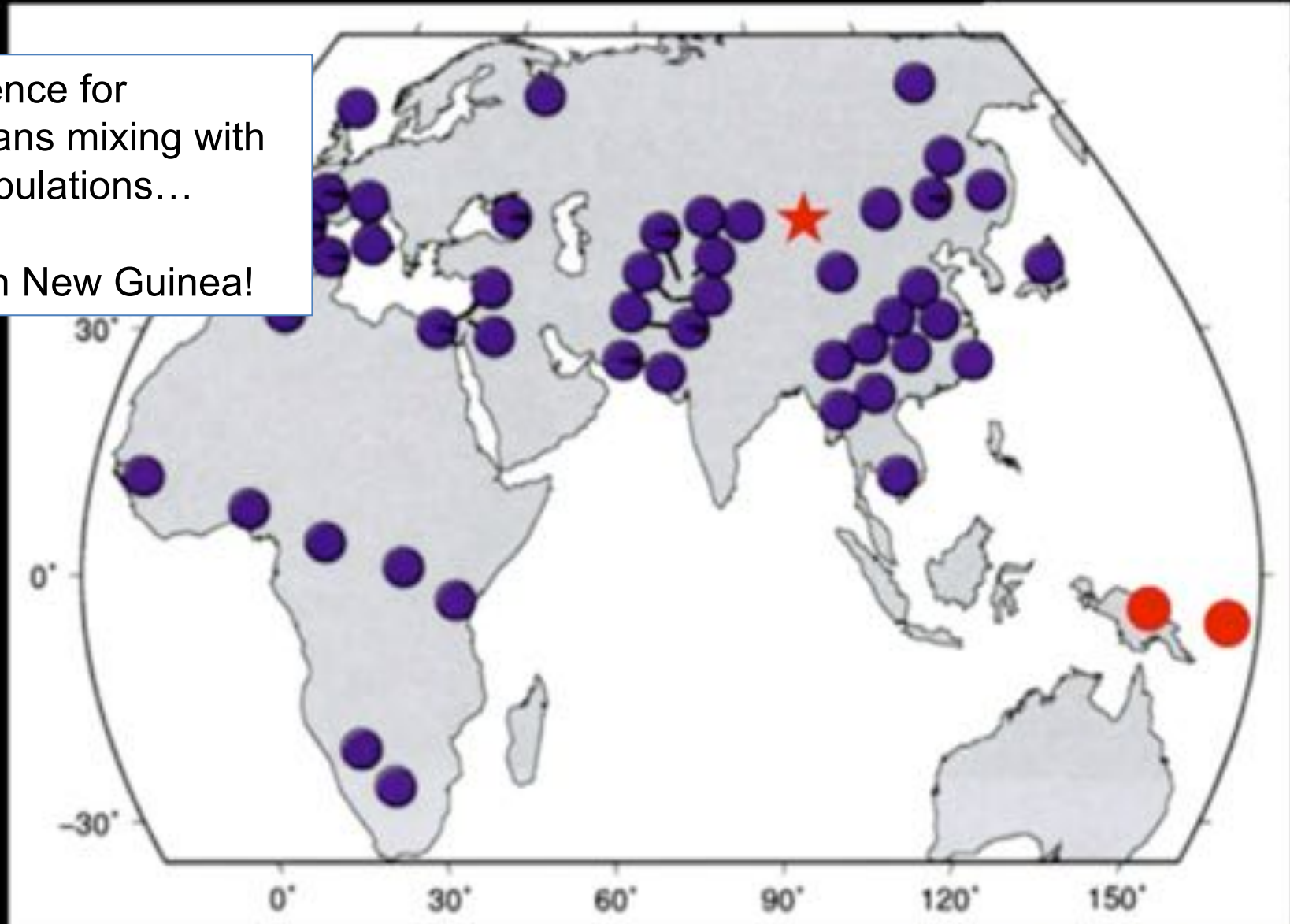
Denisovans & Neandertals



Did we mix?

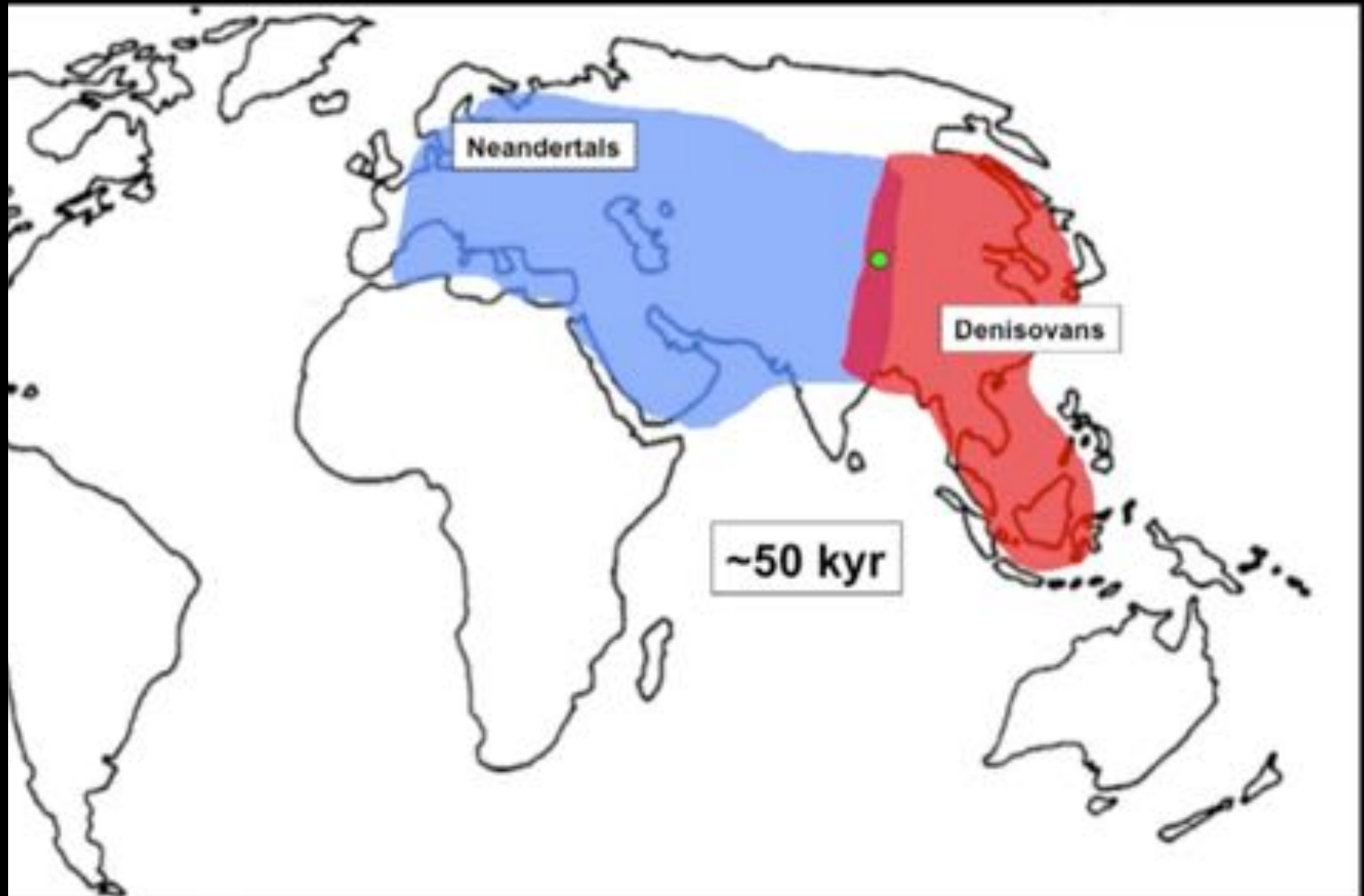
No evidence for
Denisovans mixing with
other populations...

Except in New Guinea!

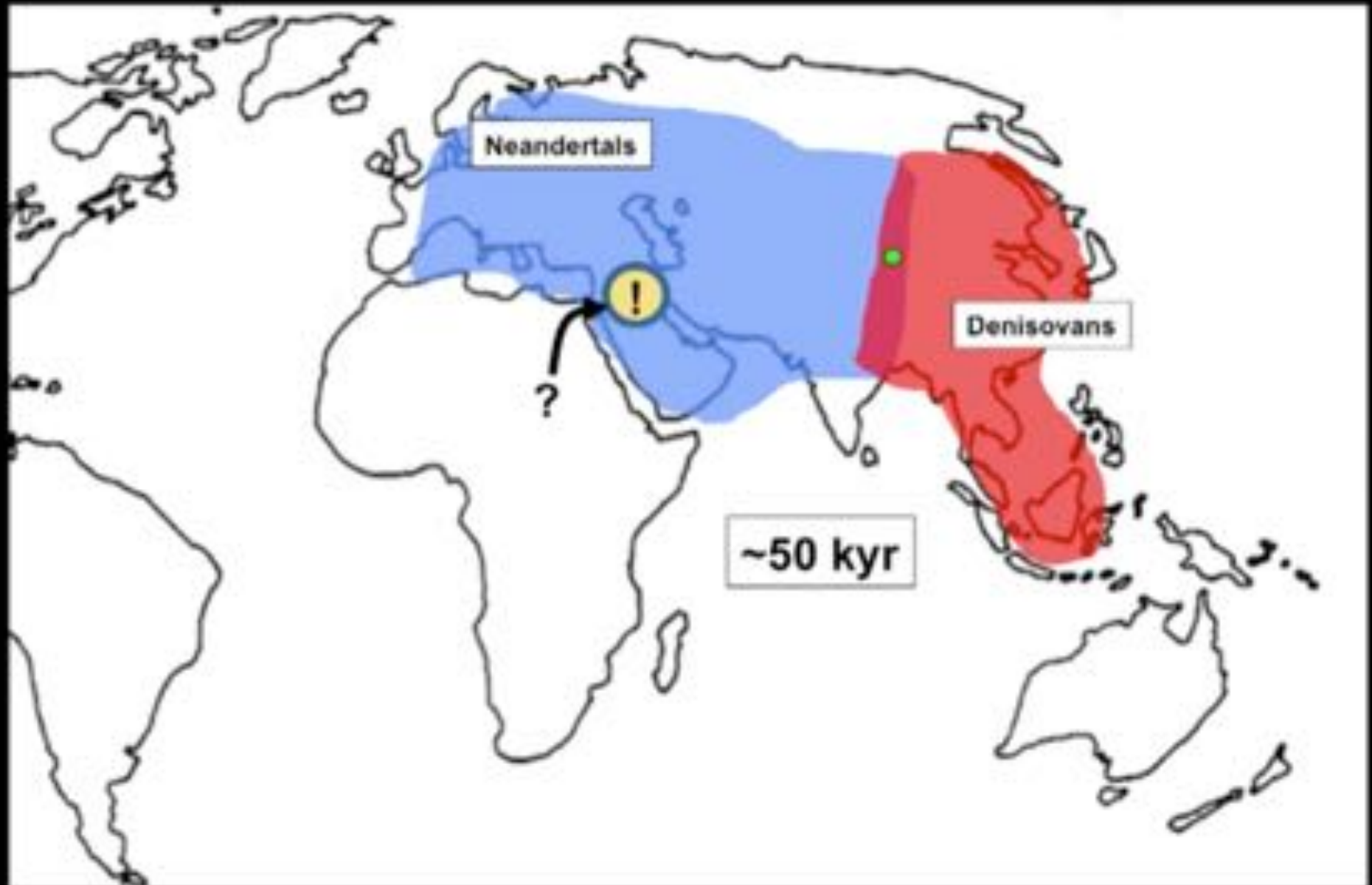


Map after Pickrell et al., 2009

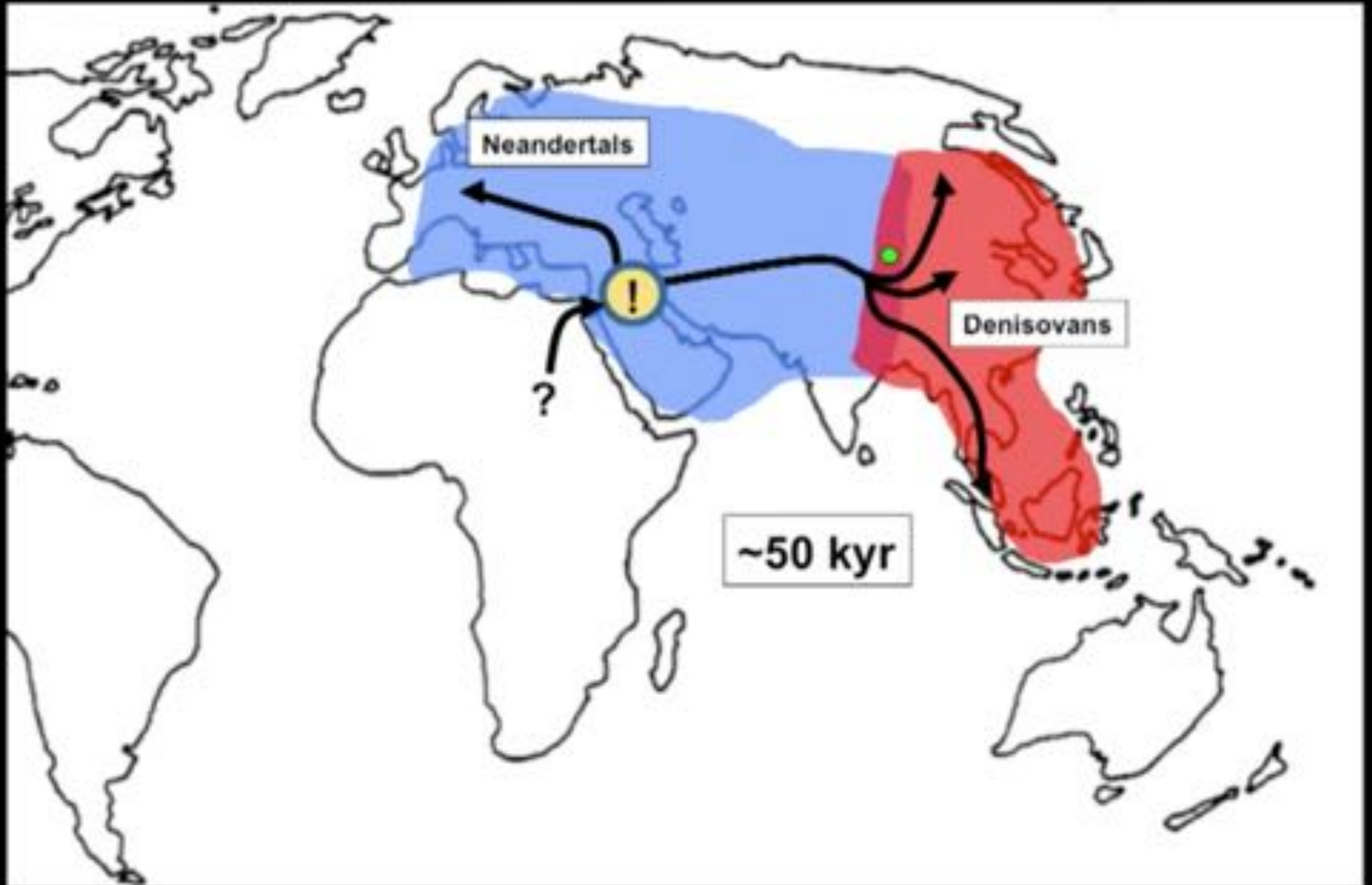
Timeline of ancient hominids



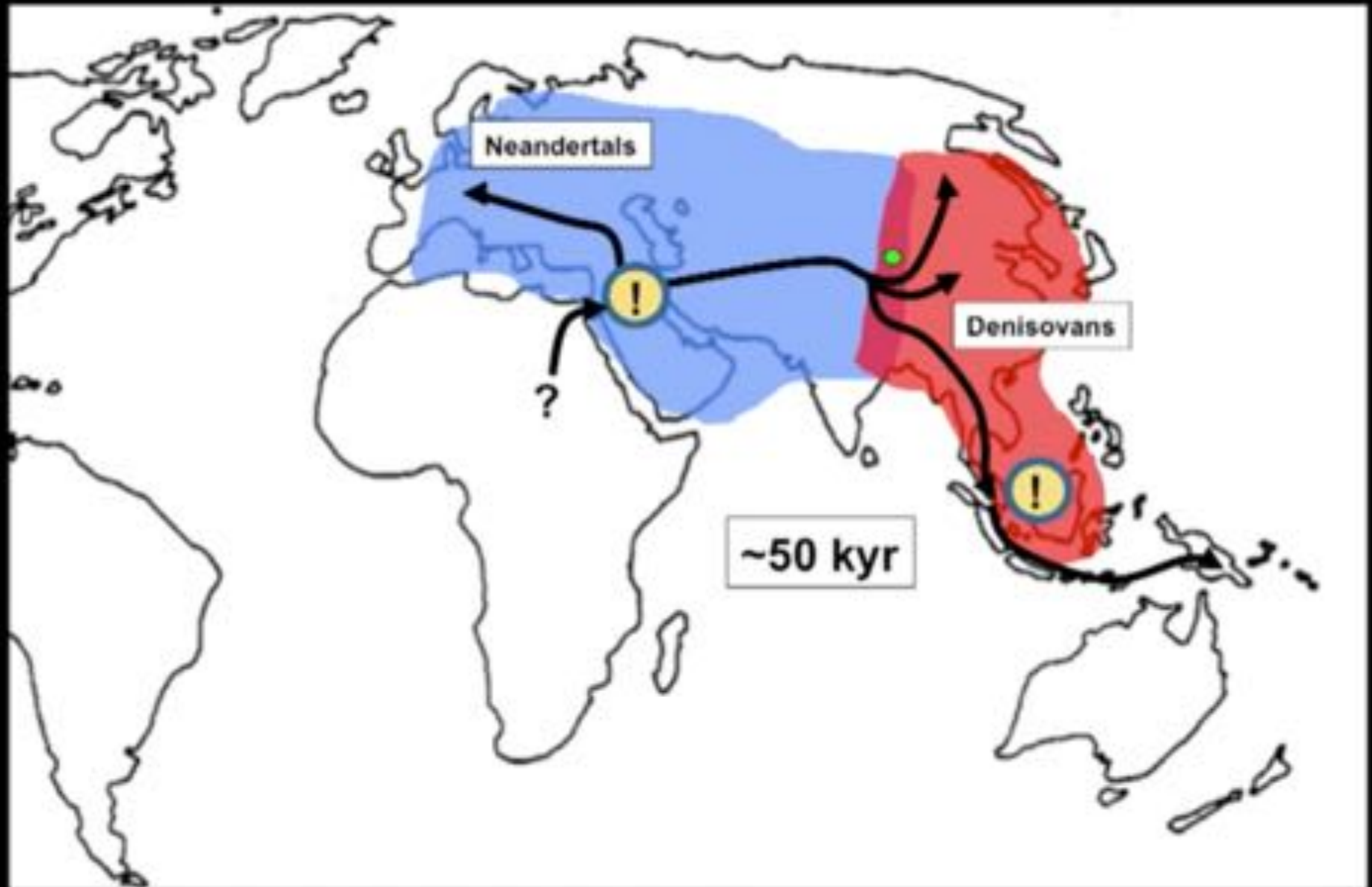
Timeline of ancient hominids



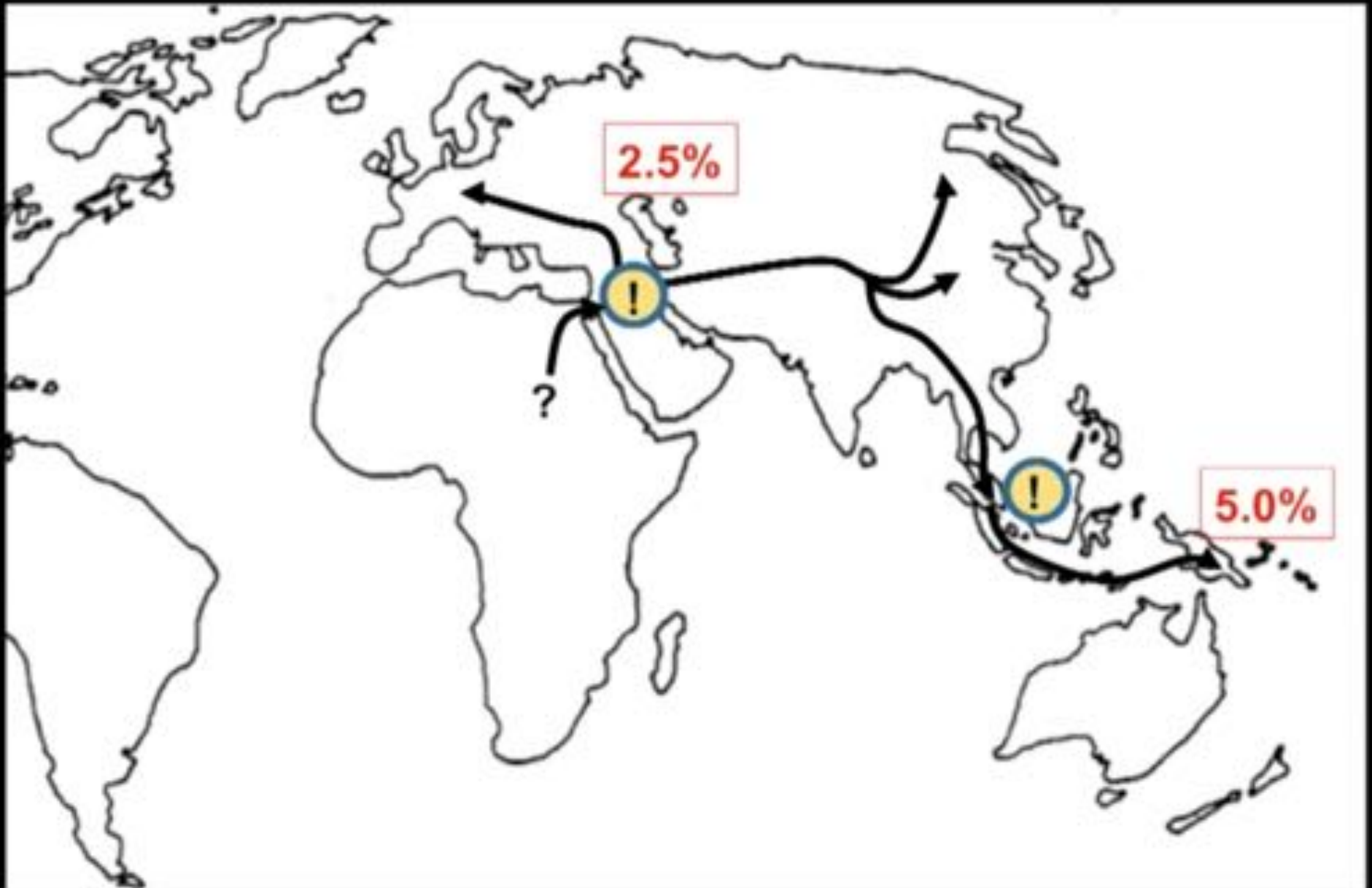
Timeline of ancient hominids



Timeline of ancient hominids



Timeline of ancient hominids



We have always mixed!

Cite as: B. Vernot *et al.*, *Science*
10.1126/science.1254166 (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. Gitterman,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, N.J., 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.

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

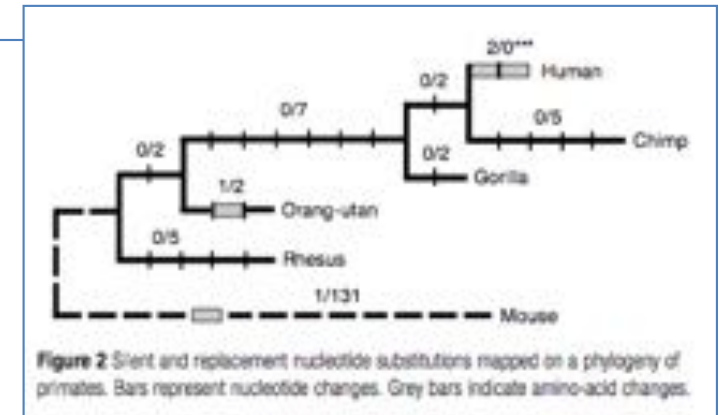
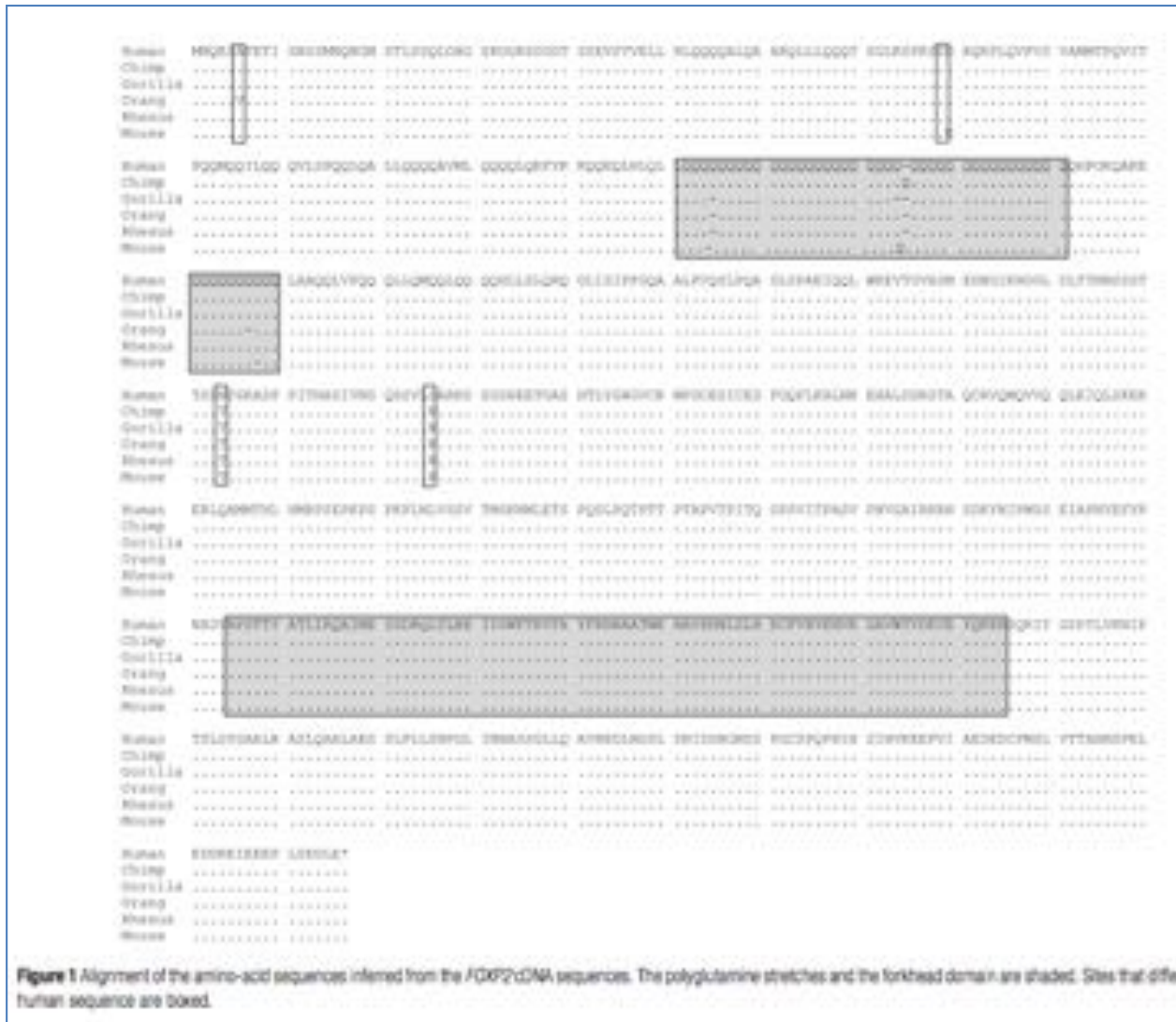
Enrichment analysis

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)
		<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;
		<ul style="list-style-type: none"> - FWER=0.538; FDR=0.0887928) - 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)
		<ul style="list-style-type: none"> - morphologies of the larynx and the epiglottis

skeletal morphologies (limb length, digit development)

morphologies of the larynx and the epiglottis

FOXP2 Analysis



- 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.

Molecular evolution of FOXP2, a gene involved in speech and language

Enard et al (2002) *Nature*. doi:10.1038/nature01025



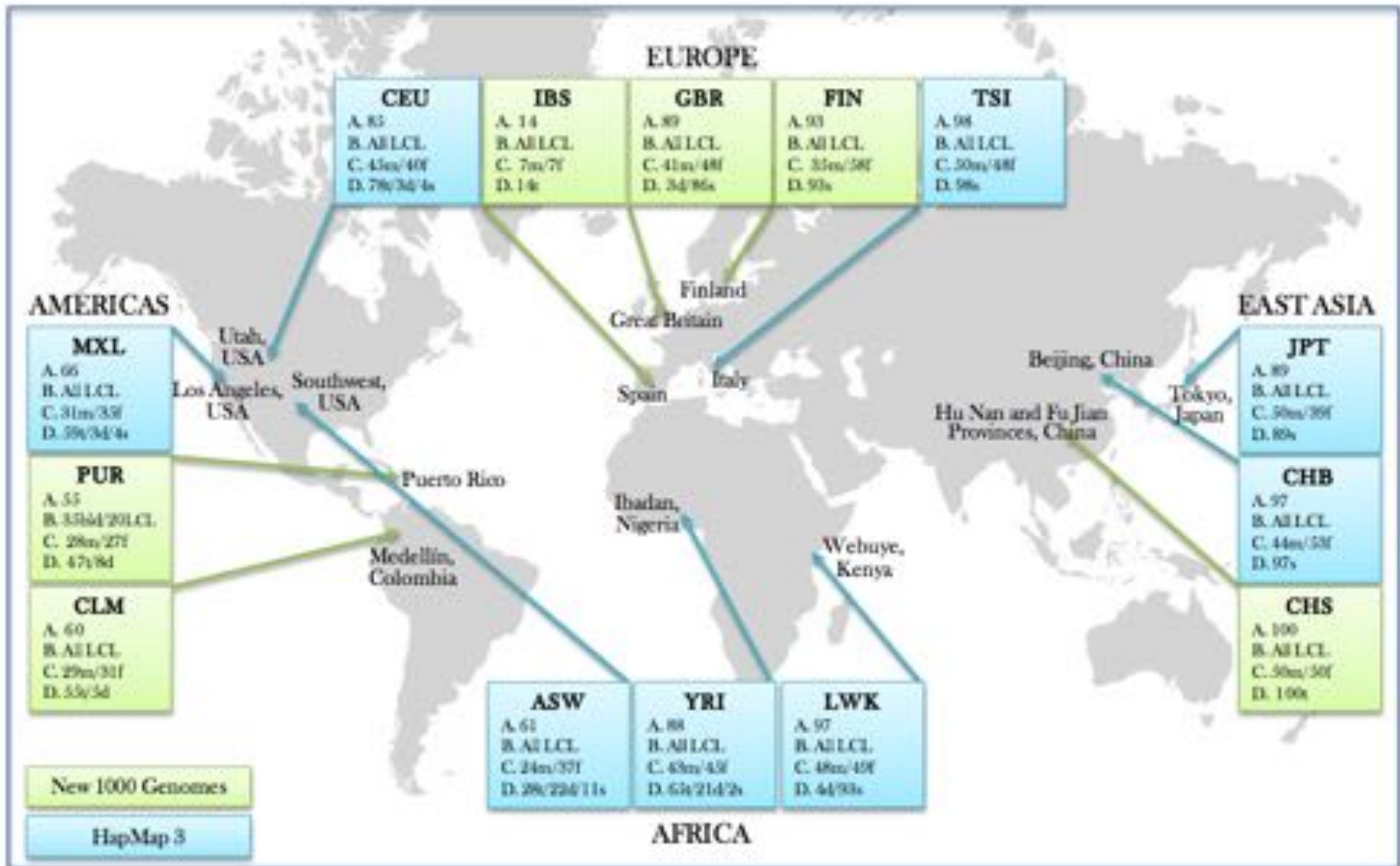
Part II: Modern Humans

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



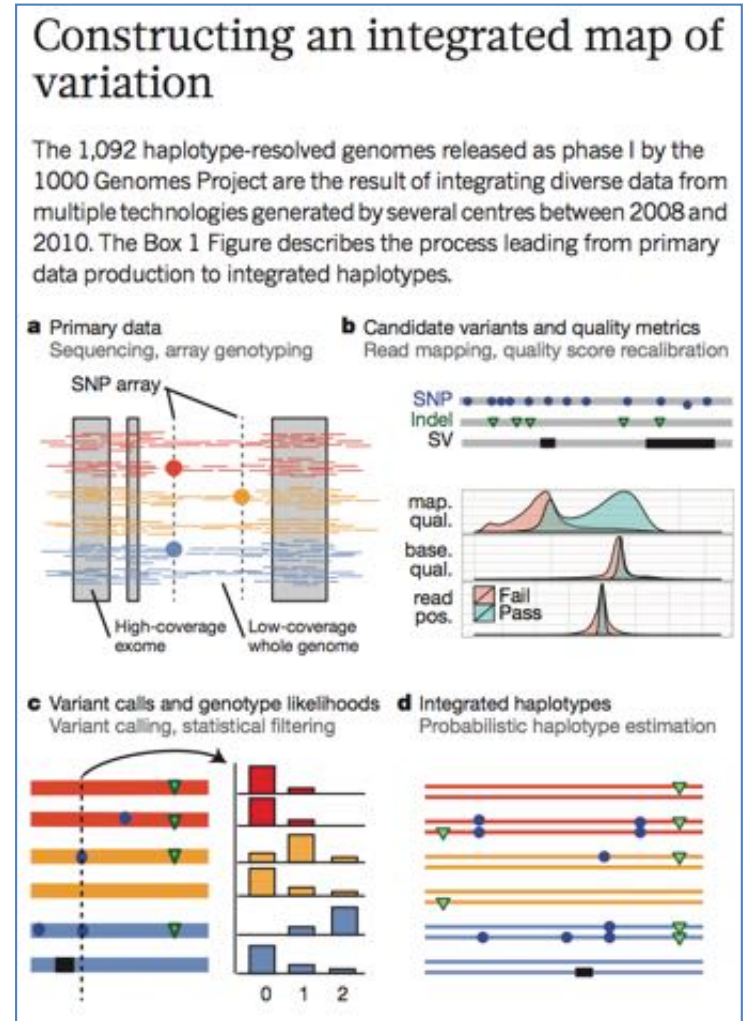
1000 Genomes Populations

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 (GHI)	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	62	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	609	601	691
British in England and Scotland (GBR)	no	yes	0	89	92	90	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	508	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	87	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	392	347	395
Total			343	1092	2530	2504	2877

26 populations from 5 major population groups

1000 Genomes: Human Mutation Rate

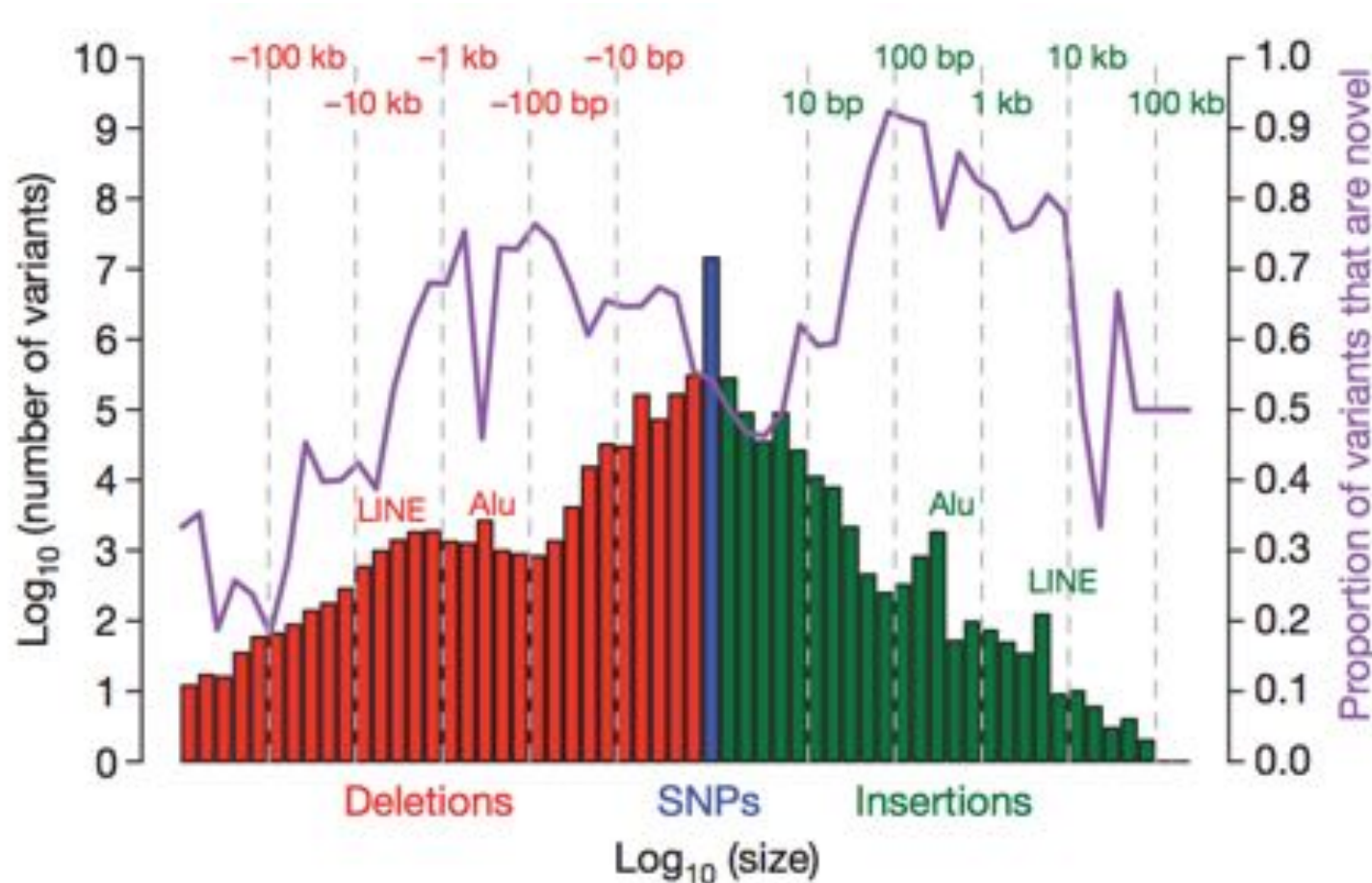
- Phase I Release
 - 1092 individuals from 14 populations
 - Combination of low coverage WGS, deep coverage WES, and SNP genotype data
- Overall SNP rate between any two people is ~1/1200bp to ~1/1300
 - ~3M SNPs between me and you (.1%)
 - ~30M SNPs between human to Chimpanzees (1%)
- De novo mutation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



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

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 Zhao⁵, Asif Rasheed², Sumeet A. Khetarpal⁶, Hong-Hee Won⁷, Konrad J. Karczewski^{4,5}, Anne H. O'Donnell-Luria^{4,5,8}, Kaitlin E. Samocha^{4,5}, Benjamin Weisburd^{4,5}, Namrata Gupta⁴, Moazzam Zaidi², Maria Samuel², Atif Imran², Shahid Abbas⁹, Faisal Majeed², Madiha 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^{4,5}, Stacey Gabriel⁴, Eric S. Lander⁴, Mark J. Daly^{4,5}, Philippe Froggert^{19,20}, John Danesh^{19,20}, Daniel J. Rader^{4,20} & Sekar Kathiresan^{1,2,4}

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 *SLC9A3R1*, 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

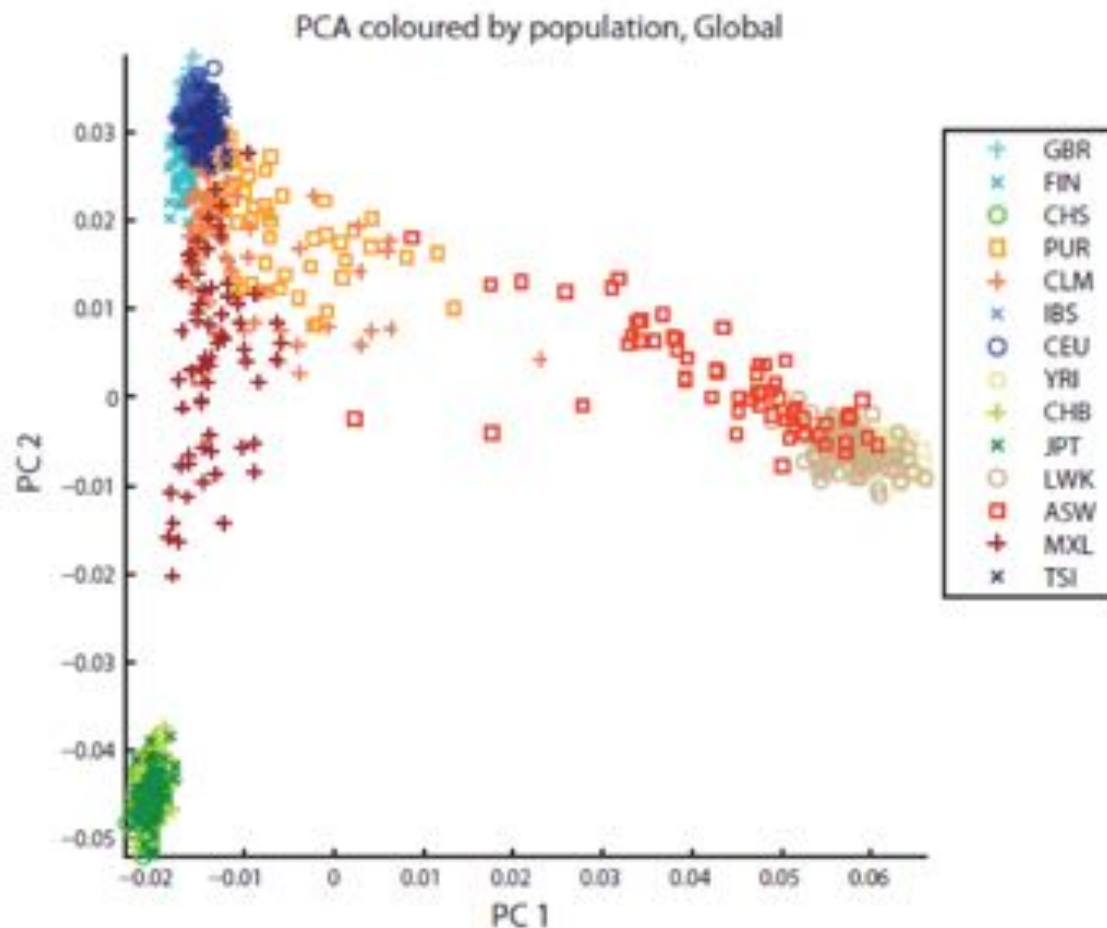
across 14,345 autosomal genes were annotated as pLoF mutations (that is, nonsense, frameshift, or canonical splice-site mutations predicted to inactivate a gene). To increase the probability that mutations are correctly annotated as pLoF by automated algorithms, we removed nonsense and frameshift mutations 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



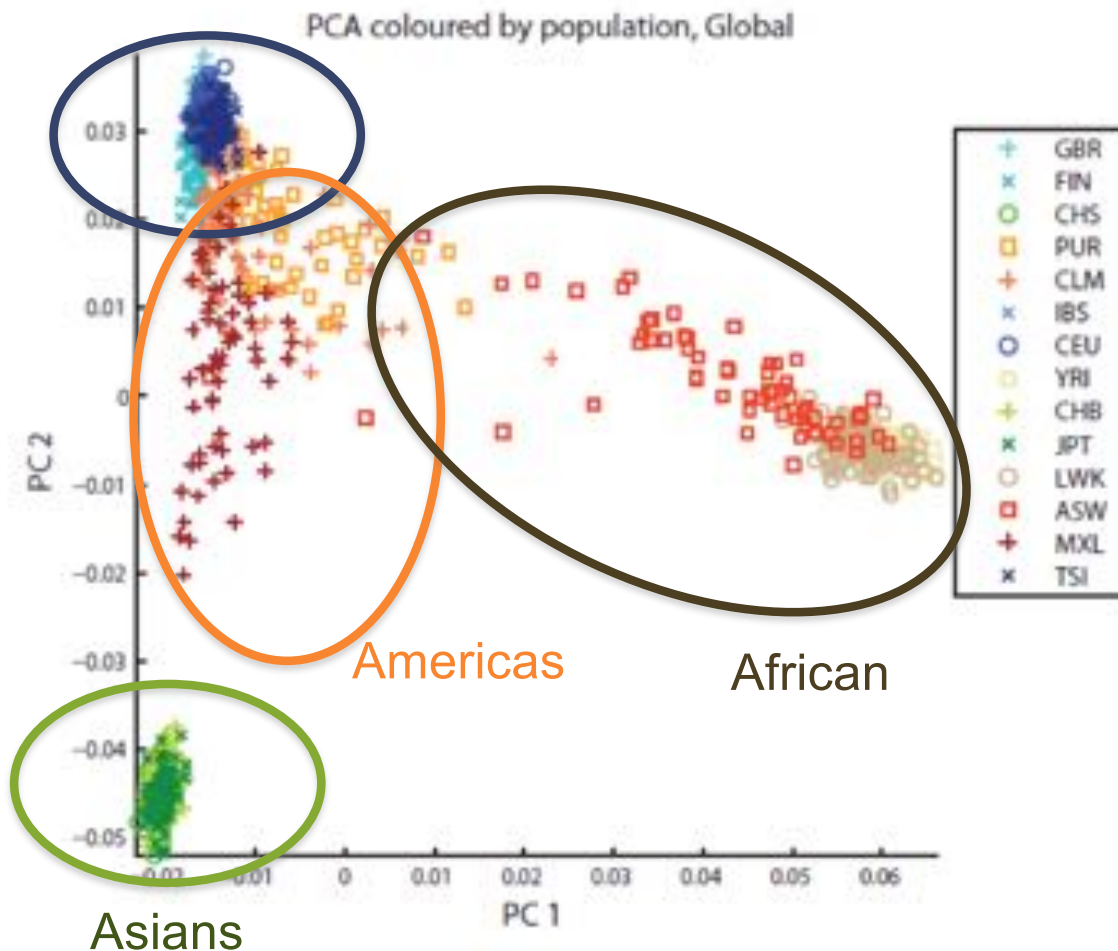
LEVEL	POP_PAIR	# of Highly differentiated SNPs	% in transcribed regions*
AFR	ASW-LWK	258	46.8
AFR	LWK-YRI	251	50.2
AFR	ASW-YRI	213	45.8
ASN	CHS-JPT	275	48.1
ASN	CHB-JPT	176	43.7
ASN	CHB-CHS	79	38.7
EUR	FIN-TSI	343	42.6
EUR	CEU-FIN	201	40.7
EUR	FIN-GBR	197	43.2
EUR	GBR-TSI	100	38.9
EUR	CEU-TSI	57	53.8
EUR	CEU-GBR	17	14.3
CON	AFR-EUR	348	52.2
CON	AFR-ASN	317	52.6
CON	ASN-EUR	190	53.4

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

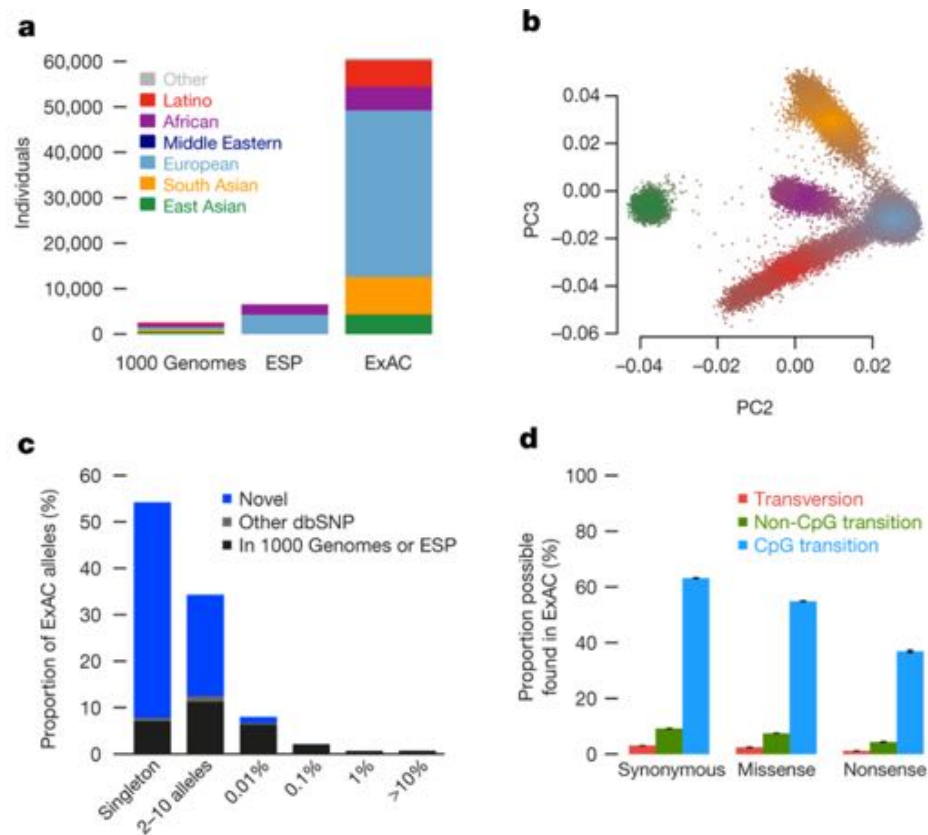


LEVEL	POP_PAIR	# of Highly differentiated SNPs	% in transcribed regions*
AFR	ASW-LWK	258	46.8
AFR	LWK-YRI	251	50.2
AFR	ASW-YRI	213	45.8
ASN	CHS-JPT	275	48.1
ASN	CHB-JPT	176	43.7
ASN	CHB-CHS	79	38.7
EUR	FIN-TSI	343	42.6
EUR	CEU-FIN	201	40.7
EUR	FIN-GBR	197	43.2
EUR	GBR-TSI	100	38.9
EUR	CEU-TSI	57	53.8
EUR	CEU-GBR	17	14.3
CON	AFR-EUR	348	52.2
CON	AFR-ASN	317	52.6
CON	ASN-EUR	190	53.4

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

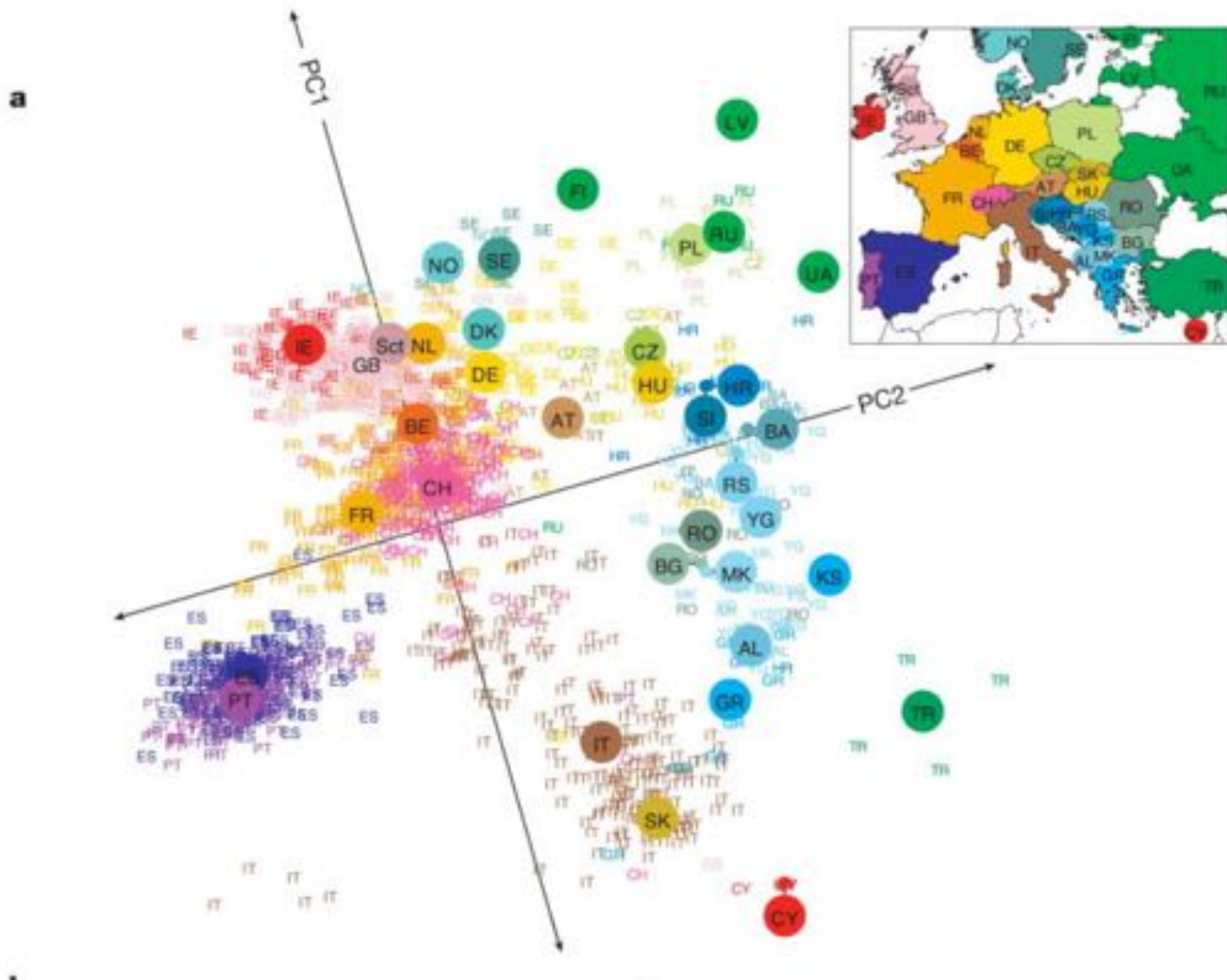
ExAC: Exome Aggregation Consortium



- 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

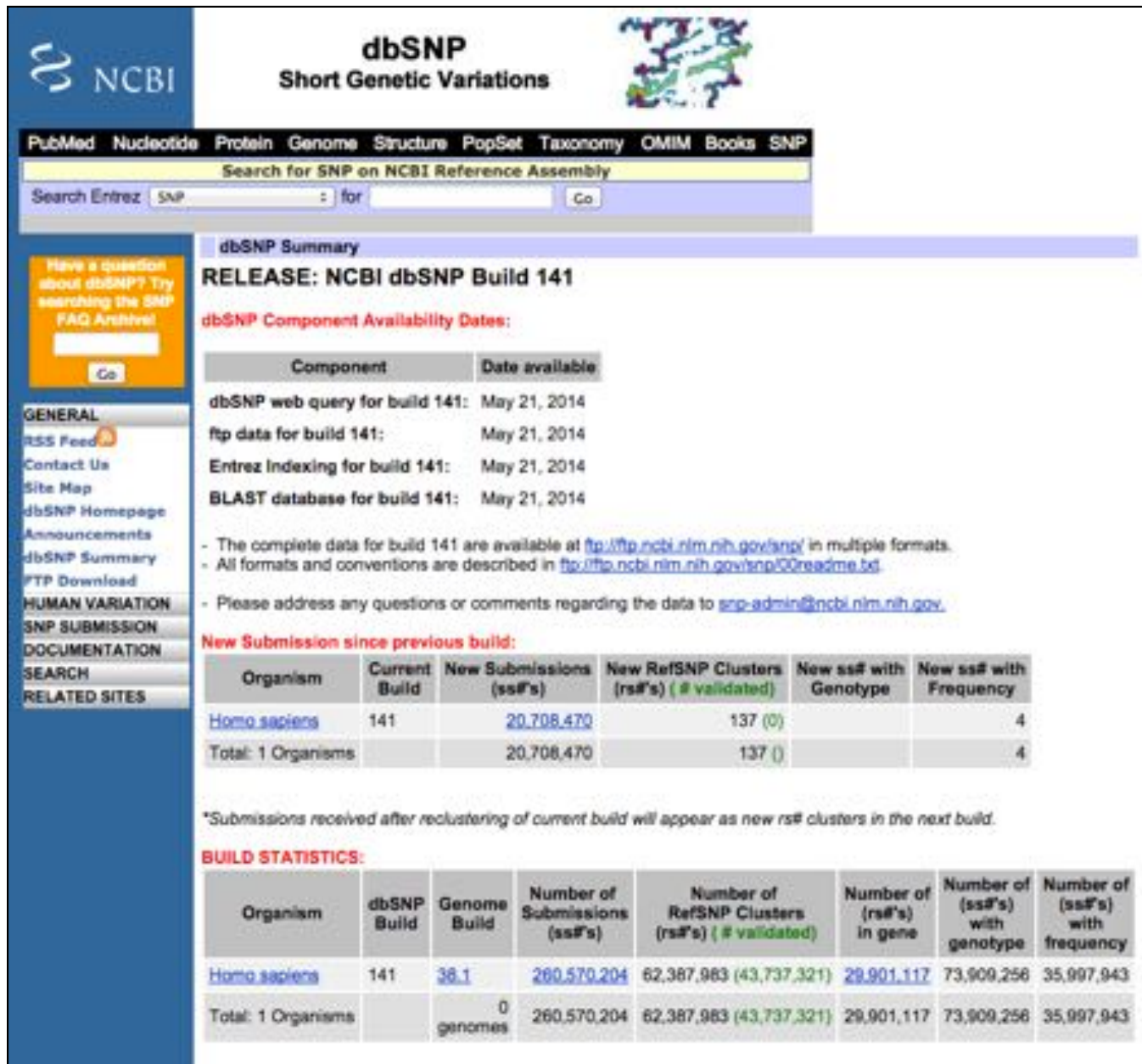
Lek et al (2016) Nature. doi:10.1038/nature19057



Genes mirror geography within Europe

Novembre et al (2008) Nature. doi: 10.1038/nature07331

dbSNP



The screenshot shows the NCBI dbSNP website interface. At the top, the NCBI logo is on the left, and the dbSNP logo with the text "Short Genetic Variations" is in the center. A navigation bar includes links to PubMed, Nucleotide, Protein, Genome, Structure, PopSet, Taxonomy, OMIM, Books, and SNP. Below this is a search bar with the text "Search for SNP on NCBI Reference Assembly". A sidebar on the left contains a "Have a question about dbSNP? Try searching the SNP FAQ Archive!" section with a "Go" button, and a "GENERAL" section with links to RSS Feed, Contact Us, Site Map, dbSNP Homepage, Announcements, dbSNP Summary, FTP Download, HUMAN VARIATION, SNP SUBMISSION, DOCUMENTATION, SEARCH, and RELATED SITES. The main content area features a "dbSNP Summary" section with the heading "RELEASE: NCBI dbSNP Build 141". Below this, a "dbSNP Component Availability Dates:" section lists the availability of various components for build 141. A "New Submission since previous build:" section provides a table of submission statistics for Homo sapiens. At the bottom, a "BUILD STATISTICS:" section provides a detailed table of build statistics for Homo sapiens.

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.