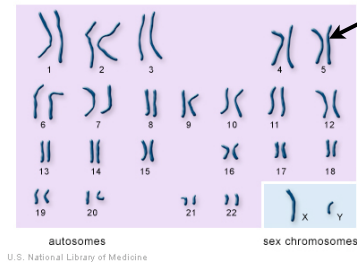


# HC70A & SAS70A Winter 2010 Genetic Engineering in Medicine, Agriculture, and Law

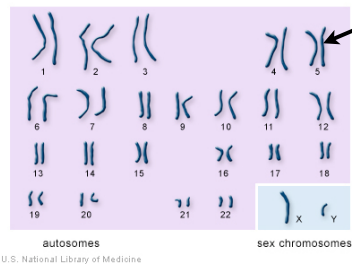
## Tracking Human Ancestry

Professor John Novembre

## A question



## Intuitive answers...



These are just some of the possible answers...

### Why might you be homozygous?

- 1) Your mother and father both gave you identical copies they received from a very recent common ancestor (identity-by-descent - inbreeding)
- 2) There is only one version of the DNA that is viable (purifying selection)

### Why might you be heterozygous?

- 1) Mutation rates at that locus are relatively high
- 2) Your parents are very distantly related (migration or large population size)
- 3) Being heterozygous for that locus is common because it is selectively favored (heterozygote advantage)

## Themes

- Global patterns of human genetic diversity
  - Tracing our ancient ancestry
  - Clines versus clusters debate
- Within-continent patterns
- Personalized genomic ancestry inference
  - What really is ancestry?
  - Admixture and Chromosome painting
- Natural selection and patterns of human genetic diversity
  - Salivary Amylase
  - Eye color (OCA2)

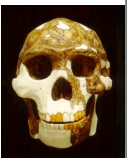
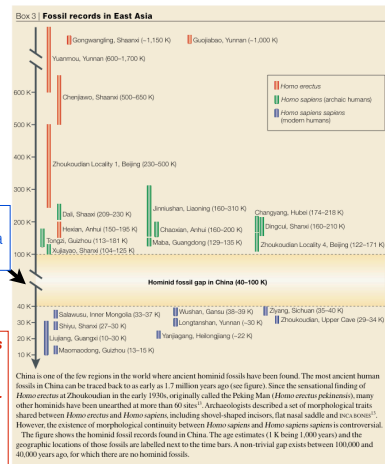
## Population genetics

- When and why should we expect to see diversity among copies of DNA sampled from a single population?
- What can we learn about a population from observing patterns of genetic diversity?
- Fundamental processes:
  - Genetic drift / inbreeding
  - Mutation
  - Natural selection
  - Migration
  - Recombination
- A field that integrates multiple scales of biology: molecular-level processes (mutation and recombination) and population-level (migration, genetic drift, natural selection) processes

## A puzzle of human ancestry

The fossil record in East Asia contains a 60,000 year gap

Did archaic humans outside-of-Africa die out, only to be replaced later on or is this just an incomplete fossil record?

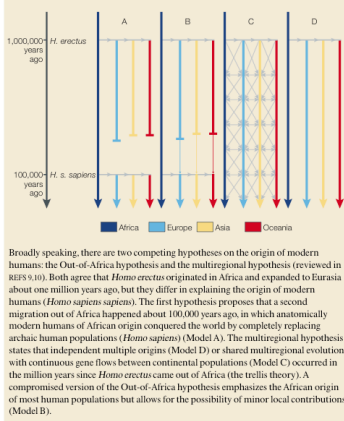


Homo erectus

Jin and Su, Nature Reviews Genetics (2000)

## Competing models of human origins

Box 1 | "Out-of-Africa" versus the multiregional hypothesis



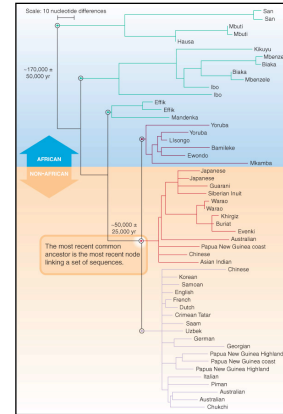
Which of these is the best supported hypothesis for origins of humans around the globe?

- (A) Recent Out-of-Africa origin
- (B) Recent Out-of-Africa origin with minor contribution from ancestral archaic humans
- (C) Shared multi-regional evolution
- (D) Independent multiple origins

Jin and Su, Nature Reviews Genetics (2000)

All human mitochondrial DNA sequences have a common ancestor "Eve" 120-220k years ago

Non-African sequences have a common ancestor at 25-75k years ago



Tree shape is consistent with Out-of-Africa origin

Mitochondrial "Eve" existed in the recent past!

## Human genome diversity panel

A global-scale sample of human genetic diversity



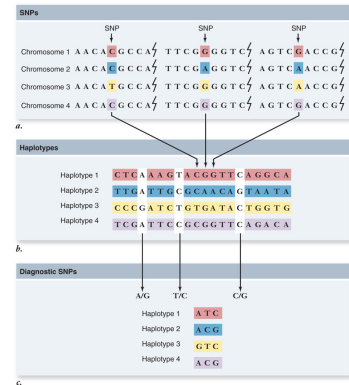
Not a perfect sampling, but the best studied yet...

Africans	Europeans	Western Asians	Eastern Asians	Oceanians
1 Yoruba	8 Orcadian	16 Bedouin	28 Han (C. China)	26 Malesian
2 Mende	9 Adygei	17 Druze	29 Han (N. China)	27 Papuan
3 Yoruba	10 Russian	18 Palestinian	30 Dai	
4 San	11 Basque	19 Daur	31 Daur	
5 Mbuti pygmy	12 French	20 Hazen	32 Hazen	
6 Baka	13 North Italian	21 Lahu	33 Lahu	
7 Mende	14 Tuscan	22 Sardinian	34 Miao	
		23 Pathan	35 Oroqen	
		24 Burusho	36 Shu	
		25 Hazara	37 Tu	
		26 Uyghur	38 Tu	
		27 Kalash	39 Yi	
			40 Mongol	
			41 Xibo	
			42 Naxi	
			43 Cambodian	
			44 Japanese	
			45 Yakut	

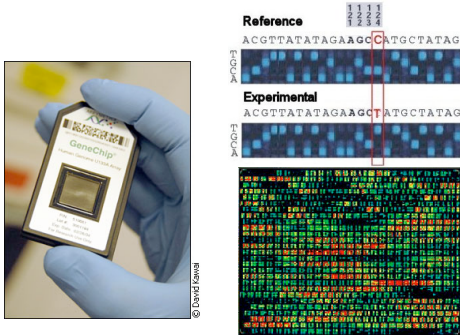
Nature Reviews | Genetics

## Review: Single nucleotide polymorphisms (SNPs) and haplotypes

Copyright © The McGraw-Hill Companies, Inc. Permission required for reproduction or display.



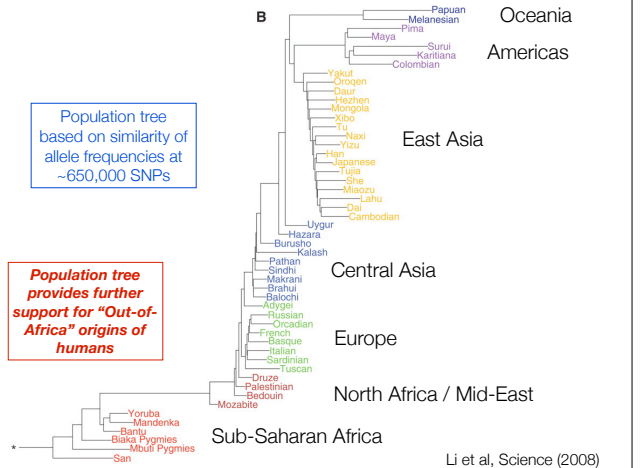
## DNA Chips Can Detect SNP Genotypes (Or Haplotypes) Across An Individual's Genome



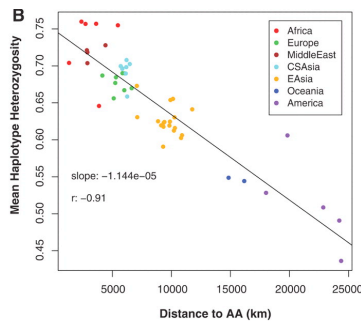
This Can Then Be Correlated With Diseases &/or Geographical Associations

Population tree based on similarity of allele frequencies at ~650,000 SNPs

Population tree provides further support for "Out-of-Africa" origins of humans



## Global patterns of haplotype diversity

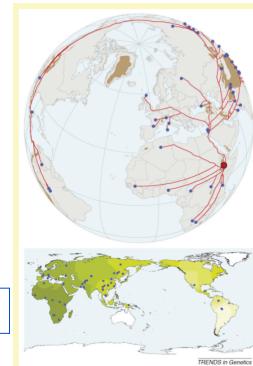


Decay of haplotype heterozygosity is consistent with a "serial bottlenecks" during Out-of-Africa expansion

Distance *via waypoints* to:  
Addis Ababa, Rift Valley, Ethiopia  
(Putative origin of early modern humans)

Examples of paths from East Africa to each of the HGDP populations

Interpolated map of genetic diversity



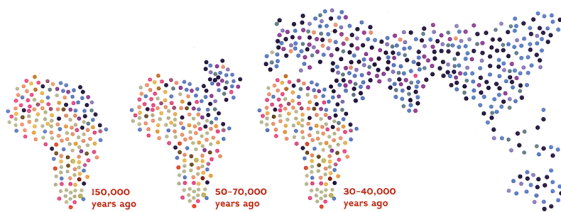
Decay of diversity is consistent with Out-of-Africa expansion

**Figure 1.** (a) Estimating geographic distances. The map shows likely colonization routes (red lines) between populations in the HGDP-CEPH panel (blue spots) assuming an origin of modern humans in East Africa (Addis Ababa, red spot). Geographic distances were estimated between populations along the colonization routes using an approach based on graph theory. Routes were found through landmasses with altitude less than 2000 m (lines over 2000 m are shown by brown shading). Geographic distances from Addis Ababa, as illustrated in this figure, as well as a matrix of pairwise distances between all HGDP-CEPH populations are available as supplementary material online. (b) Interpolation of global human genetic diversity. The intensity of the green color represents the genetic diversity obtained with an inverse distance weighted (IDW) interpolation method on landmasses using the ArcGIS Spatial Analyst extension. Blue dots represent the 54 populations from the HBT1 subset of the HGDP-CEPH dataset (35).

## Most Genetic Diversity Originated in the Founder Populations to Modern Humans!

### Diverse From the Start

The diversity of genetic markers is greatest in Africa (multicolored dots in map), indicating it was the earliest home of modern humans. Only a handful of people, carrying a few of the markers, walked out of Africa (center) and, over tens of thousands of years, seeded other lands (right). "The genetic makeup of the rest of the world is a subset of what's in Africa," says Yale geneticist Kenneth Kidd.

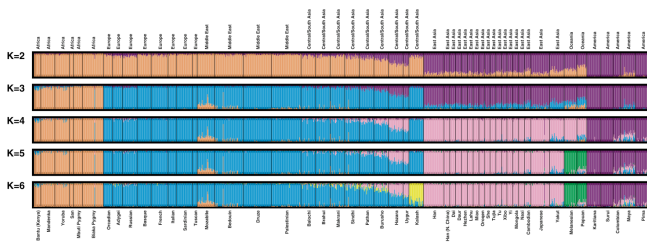


## Summary: Human origins

- At a global scale, genome-wide diversity patterns broadly consistent with a single, recent origin of modern humans in Africa
- Further support available from recent Y-chromosome and autosomal gene TMRCA dates
- Note:
  - A very small number of loci show very ancient TMRCA dates (e.g. 2 million years old).
  - **Open question:** Are these evidence for rare, ancient contributions from archaic humans?

## Continental-scale clusters of human variation?

Results of an *unsupervised* clustering algorithm on microsatellite diversity



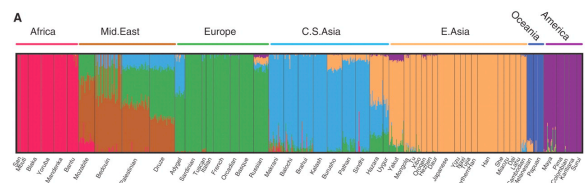
**Fig. 1.** Estimated population structure. Each individual is represented by a thin vertical line, which is partitioned into  $K$  colored segments that represent the individual's estimated membership fractions in  $K$  clusters. Black lines separate individuals of different populations. Populations are labeled below the figure, with their regional affiliations above it. Ten structure runs at each  $K$  produced nearly identical individual membership coefficients, having pairwise similarity coefficients above 0.97, with the exceptions of comparisons involving four runs at  $K = 3$  that separated East Asia from Kalash, and one run at  $K = 6$  that separated Karitiana instead of Kalash. The figure shown for a given  $K$  is based on the highest probability run at that  $K$ .

Genetic "clusters" approximate continental-scale regions

Rosenberg et al (2002) Science

## Continental-scale clusters of human variation?

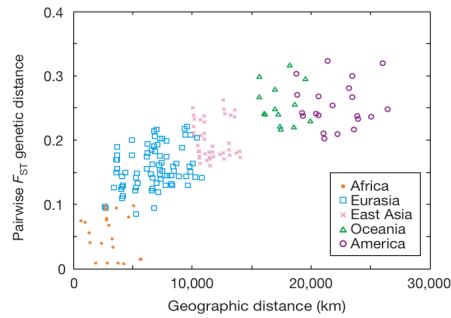
Results of *unsupervised* clustering algorithm on ~650,000 SNPs



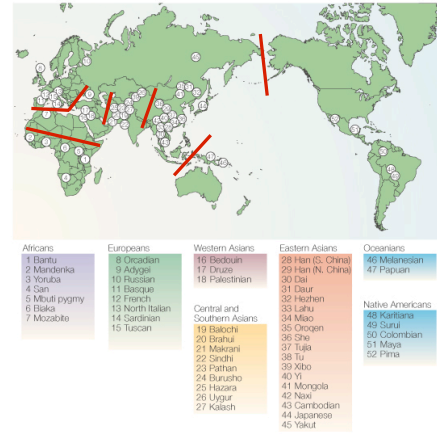
Clusters are now more detailed  
Europe, Middle East, and Central South Asia  
distinguishable

Li et al (2008) Science

Or does differentiation increase smoothly with geographic distance ("clines")?



## Clines vs. clusters debate



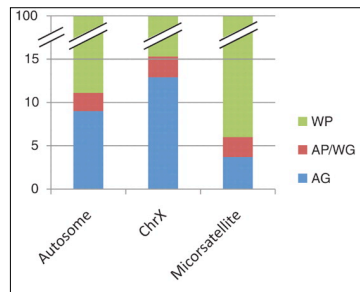
Perhaps clusters are due to impact of geographic barriers?

Sahara  
Tibetan Plateau  
Bering Strait  
Malay archipelago

Or perhaps clusters are an artifact due to gaps in sampling?

Nature Reviews | Genetics

Either way: Variation is mostly found within rather than between groups

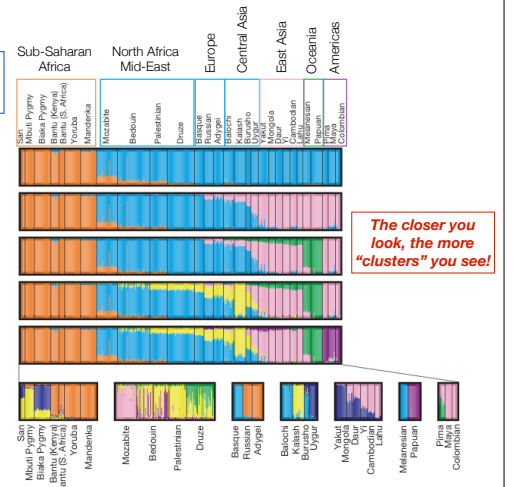


Within populations  
Among population / within groups  
Among groups

Li et al (2008) Science

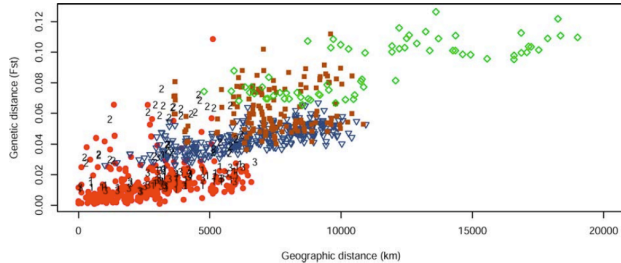
Admixture proportions assuming  $K$  unknown populations

$K = 2$   
 $K = 3$   
 $K = 4$   
 $K = 5$   
 $K = 6$



The closer you look, the more "clusters" you see!

Jakobsson et al (2008) Nature



**Figure 6.** Genetic and Geographic Distance for Pairs of Populations  
Red circles indicate comparisons between pairs of populations with majority representation in the same cluster in the  $K = 5$  plot of Figure 2; blue triangles indicate pairs with one population from Eurasia and one from East Asia; brown squares indicate pairs with one population from Africa and the other from Eurasia; and green diamonds indicate pairs with one population from East Asia and the other from either Oceania or America. Comparisons involving one of Hazara, Kalash, and Uyghur and other populations from Eurasia or East Asia are marked 1, 2, and 3, respectively. No comparisons are shown between any of these three groups and any African population.  
DOI: 10.1371/journal.pgen.0010070.g006

Patterns of genetic distance as a function of geographic distance reveal cluster and cline patterns

## Fine-scale analysis within a continent: Europe as a case study



1400 individuals from 37 unique populations in Europe



Analysis using 197,000 SNP loci

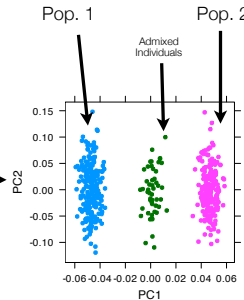
## Principal Components Analysis (PCA)=Reduction of dimensions

### Single Nucleotide Polymorphisms

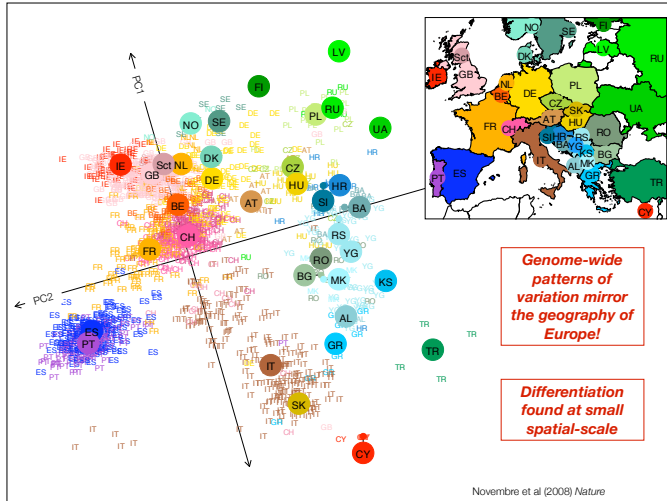
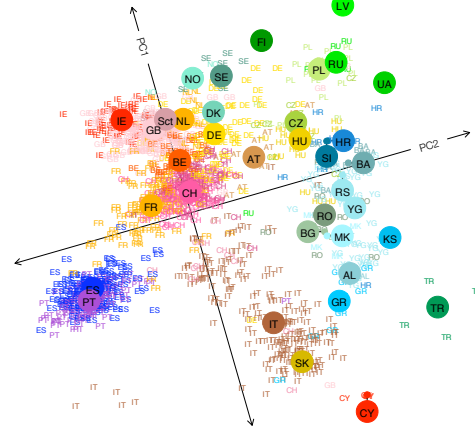
Individuals	PC1	PC2
20110112010220112212010210	0.055	-0.08
1020112202012011121201211	0.042	0.10
2021202200120021122211121	-0.052	-0.05
00111212022121222202122221	0.010	0.05
00301222020221211221212222	0.052	0.08
21111002012011222211112212	0.043	-0.02
20102212111101221221110222	-0.058	0.04
1221112021201222212012221	0.020	-0.03
20211012010220122212120221	-0.030	-0.05
0110112020100222121212221	-0.050	0.07
0110001200112122212212221	0.052	0.08
01201012120120022120112222	0.043	-0.02
111012201012022222121211	-0.058	-0.04
110021220212200121122210212	-0.04	0.06
101021220102201222111102220	-0.030	-0.02
1112112011201122120102220	-0.031	0.05
21202122022101122221121211	-0.053	0.08
2121122121201222212112222	0.014	-0.02
11210022011210122221002221	0.533	-0.04
21210220112200111022001120	-0.044	0.03

0 = AA  
1 = AB  
2 = BB

PCA is often a useful tool to visualize patterns of population structure

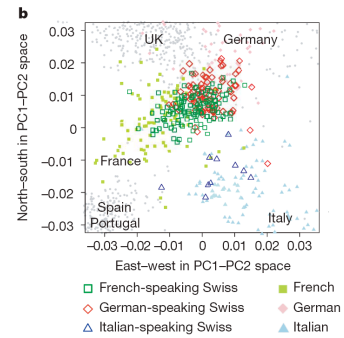


## Two-dimensional summary of European genetic diversity

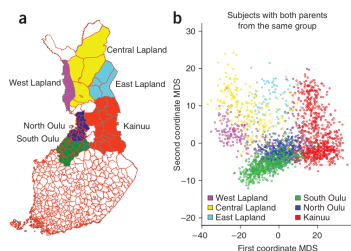


Novembre et al (2008) Nature

## Subtle differentiation within Switzerland



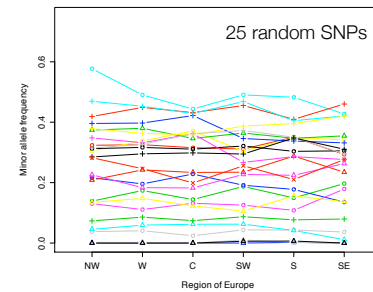
## Differentiation within Finland



**Figure 1** Linguistic/geographic groups of Northern Finland and their genetic signature. (a) Map of Finland with county boundaries. The subjects in NFB1965 were all born in the two northern provinces. Counties in Northern Finland are color coded to correspond to the six linguistic/geographical groups that can be identified. (b) Scatterplot of the two first components identified by MDS on the matrix of genetic similarity between individuals. Only subjects with both parents born in the same population group are plotted, and they are color coded according to the group of origin.

VOLUME 41 | NUMBER 1 | JANUARY 2009 NATURE GENETICS

## Per locus information content extremely low



At any given SNP, there is little variation among populations

PCA methods pool weak signals across many, many loci to reveal differentiation

• Across loci mean  $F_{ST}$ =0.004!



## Summary: Clines versus Clusters

- At a global scale, support for clusters and clines can be observed
- With large numbers of SNP markers, patterns of differentiation are detectable even at small-scales within continents (although often more “clinal” than “clustered”)
  - Note: Variation is still predominately within vs. between groups
- **Future directions:** With whole genomes - we may detect even more subtle patterns of differentiation

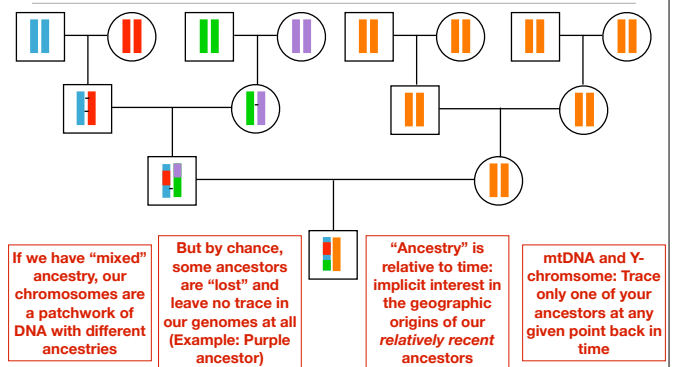
## Applications: Ancestry Inference in Personalized Genomics

The screenshot shows the 23andMe website. At the top, it says "See your genes in a whole new light. TIME Magazine's 2008 Invention of the Year, now \$399." Below this is a navigation bar with "How it works", "Buy US \$399", and "Try a demo". On the right, there's a login section with fields for "Login name", "Password", and "Forgot password?". Below the login section, there's a section titled "deCODE your ancestry" with a sub-header "learn about your ancestral story". The text describes how genetic data is used to make inferences about people's ancestors and geographical relationships. It mentions that all humans descend from a common ancestor group that originated in Africa about 200,000 years ago. Below this, there's a section titled "The ancestry results" which includes a map of kinship, genetic data, ancestral origins, and a map of the world. A link to "OUR COMPLETE STORY" is also present.

## What do we mean by ancestry?

- My mtDNA test comes back and says I have a Native American mitochondrial haplotype. Does this mean:
  - (a) I am likely to be 100% Native American
  - (b) My mother's side is likely to be 100% Native American
  - (c) My mother's mother's mother's mother's mother's mother's mother was likely Native American.
  - (d) Not enough information to tell.

## What do we mean by ancestry?



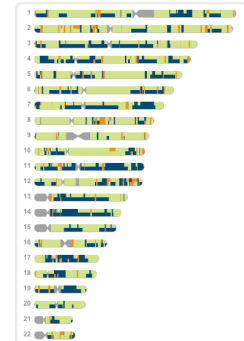
## ancestry painting

Trace the ancestry of your chromosomes, one segment at a time. Last updated April 23th, 2008.

### Chromosome View

Solid segments indicate that both chromosomes come from the same geographic region. See a Cambodian Woman's painting. Dual-colored segments indicate chromosomes from different geographic regions. See an African American Man's painting.

Select a person: **African American Woman**



### African American Woman

Most African Americans today trace a large part of their ancestry to sub-Saharan Africa as a result of the slave trade. Over the generations since, both Europeans and Native Americans have intermarried with African Americans and contributed ancestry, as seen in the ancestry painting of this woman, who identified herself as African American.



### Worldwide Examples

Click on the icons in the map below to see example paintings of individuals from across the globe.

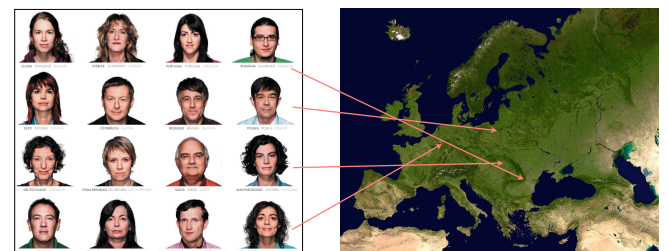


Tell Me About...

Sensitivity to reference populations: "Asian" ancestry might actually be Native American

Painting is the result of statistical inference, but often difficult to communicate the uncertainty

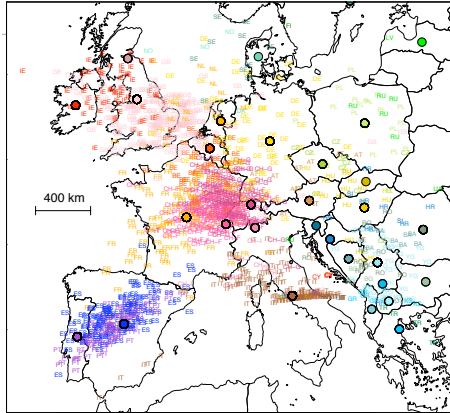
## Pushing the limits: Ancestry inference within continents



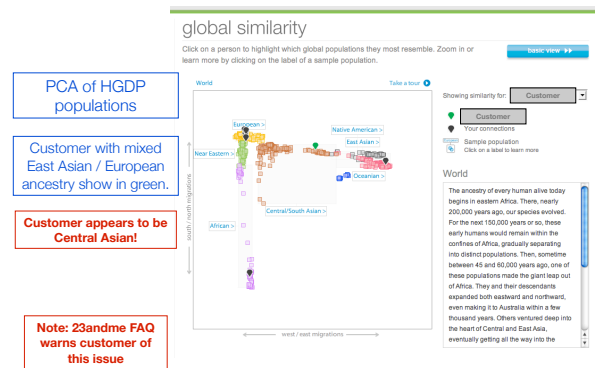
## Performance: Spatial prediction

- Regression-based approach using PC scores
- "Leave-one-out" cross-validation performance:
- For countries with  $n > 15$ :  
50% : < 240 km  
90% : < 630 km
- Overall:  
50% : < 540 km  
90% : < 840 km

**Caveat: With PCA, mixed ancestry leads to an intermediate positioning**



## PCA and mixed ancestry



## Summary: Personal ancestry inference

- Immense potential, but numerous challenges:
  - Obtaining the appropriate reference samples
  - Communicating to customer:
    - What we really mean by ancestry
    - Statistical uncertainty in the inferences
    - Limitations of particular analyses
  - Danger of customers forming notions of group membership / race?

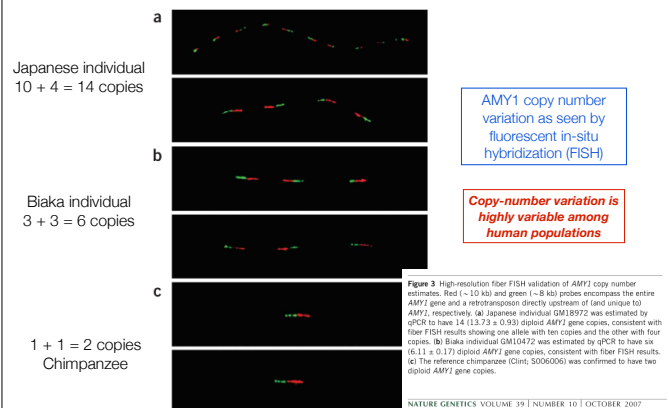
## What's your opinion:

- Personal ancestry inference is:
  - (a) A waste of money.
  - (b) A harmless hobby if that's what you're in to.
  - (c) Fascinating - sign me up!
  - (d) The first steps towards genetic elitism.

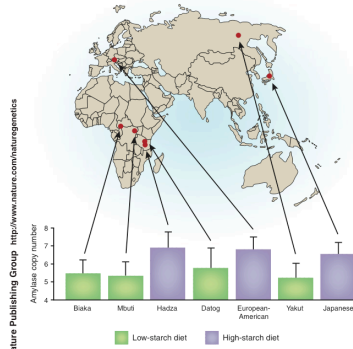
## Question:

- Have cultural changes had an effect on human evolution?
  - (a) Yes
  - (b) No

## Salivary amylase copy number variation



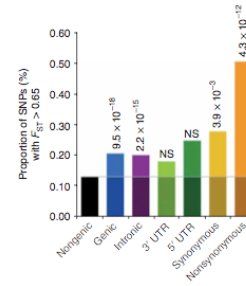
## Diet and Variation in salivary amylase copy number



High-starch diet populations have higher average copy number for than low-starch diet populations

## The impact of natural selection on population differences

SNPs in genic regions and especially non-synonymous SNPs are more often extremely differentiated than "non-genic" SNPs



Positive selection in one population but not another can lead to differentiation at SNPs in the region of the gene under selection

The most differentiated SNPs between human groups are likely to be in regions that have undergone recent adaptive evolution

Barreiro et al (2008) Nature Genetics

Table 1 Genes showing the strongest signatures of positive selection

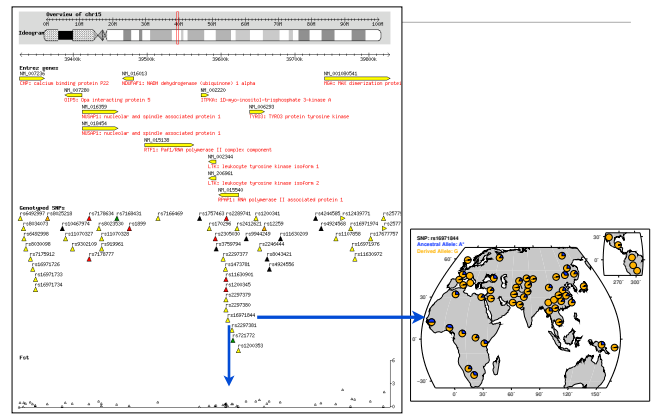
Phenotype category	Genes
Morphological traits (for example, skin pigmentation and hair development)	ABCC11, EDAR, SLC45A2, PKP1, PLEKHA4, SLC24A5
Immune response to pathogens	CEACAM1, CR1, DUOX2, VHW2
DNA repair and replication	MPG, POLG2, TDP1
Sensory functions (for example, olfaction and eye development)	COL18A1, OR52K2, RP1L1
Insulin regulation, metabolic syndrome (obesity, diabetes, hypertension)	ALMS1, CEACAM1, ENPP1
Various metabolic pathways (for example, ethanol, intestinal zinc and citrulline)	ADH1B, ASS1, SLC39A4
Miscellaneous	FBXO31, RTTN, SPAG6
Unknown	ABCC12, ADAT1, AK127117, C17orf46, C8orf14, COLEC11, CPSF3L, DNAJC5B, DNHD1, ETVF1, EXOC5, FAM, CDC142, FLJ37464, FOXL1, GOSL2, KIAA0984, LAMB4, LOC64851, LIMCH1, PCGF1, PLEKHG4, POLS3P, RNF135, SLC30A9, SYTL3, TEX15, TTC3P, VPS33B, ZNF646

What we see as surface indicators of "race" (skin, hair, eye-color) are some of the most unique regions of the genome

These genes contain at least one nonsynonymous or 5'-UTR mutation with  $F_{ST} > 0.65$ . An exhaustive list of 582 genes containing other classes of genic SNPs with  $F_{ST} > 0.65$  is provided in Supplementary Table 1. Genes in bold correspond to those also presenting significant long-range haplotypes, as measured by the iHS statistic<sup>2</sup>, or defined as top candidates for recent selective sweeps<sup>3</sup>. \*These genes have not yet been attributed a HUGO-approved symbol. †These three genes are located in a linkage-disequilibrium block in chromosome 2. ‡These two genes are located in a linkage-disequilibrium block in chromosome 16.

Barreiro et al (2008) Nature Genetics

## Differentiation at a random region of Chromosome 15





- Humans have been recently, and continue to be, evolving
- Patterns of genetic variation in humans point towards recent common ancestry in Africa
- With modern large-scale data sets, we can identify the personal ancestry of individuals to a fine spatial scale
- Many of the most differentiated regions of the genome seem to be the results of selection related to:
  - Novel diets
  - External morphology in response to different climates
  - Immune system / disease evasion
- Beyond these few differentiated regions (some of which are relevant to medicine), most variation is found globally (and most of the genome doesn't even vary!)
- Intelligent views about our common genetic heritage and diversity will be crucial in our post-genomic world