

HC70A & SAS70A Winter 2009  
Genetic Engineering in Medicine,  
Agriculture, and Law

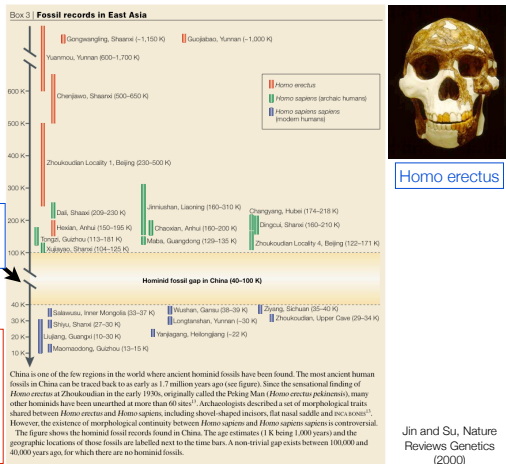
Tracking Human Ancestry

Professor John Novembre

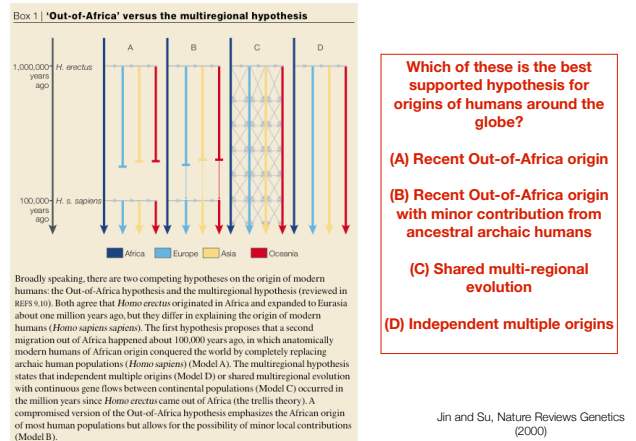
Themes

- 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)

A puzzle of human ancestry

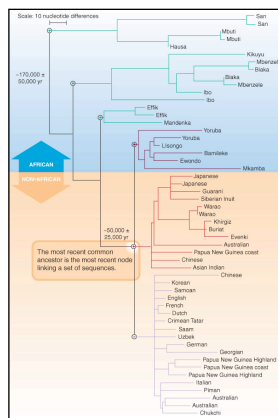


Competing models of human origins



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

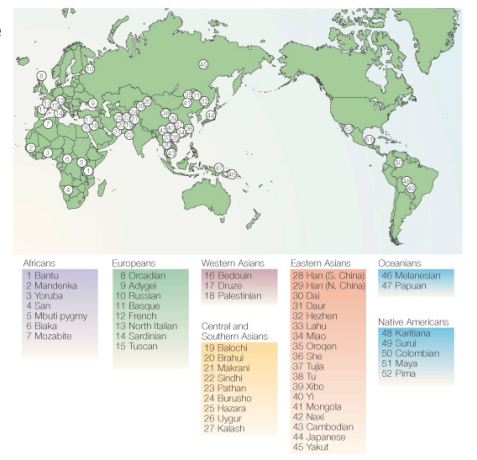
Non-African sequences have a common ancestor at 25-75k years ago



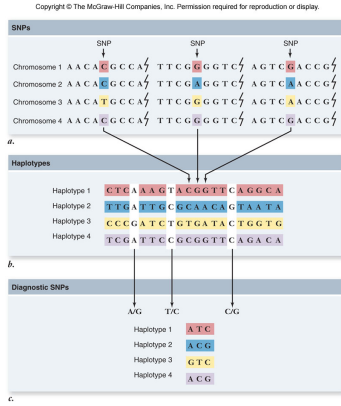
Human genome diversity panel

A global-scale sample of human genetic diversity

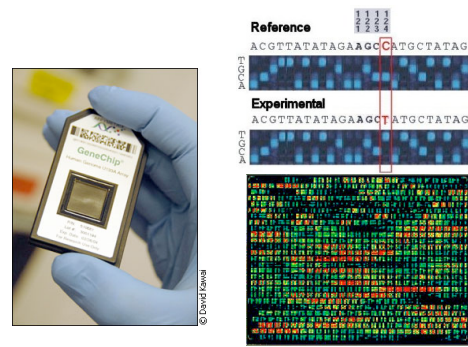
Not a perfect sampling, but the best studied yet...



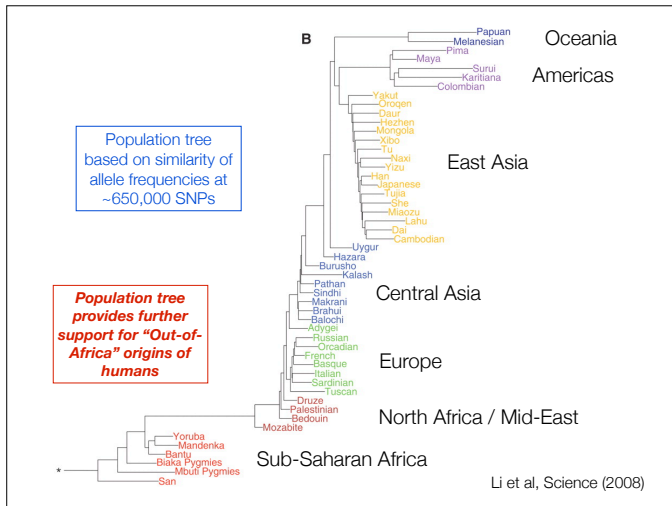
## Review: Single nucleotide polymorphisms (SNPs) and haplotypes



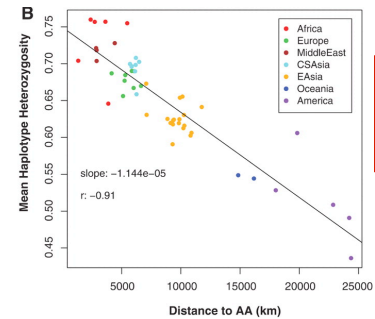
## DNA Chips Can Detect SNP Genotypes (Or Haplotypes) Across An Individual's Genome



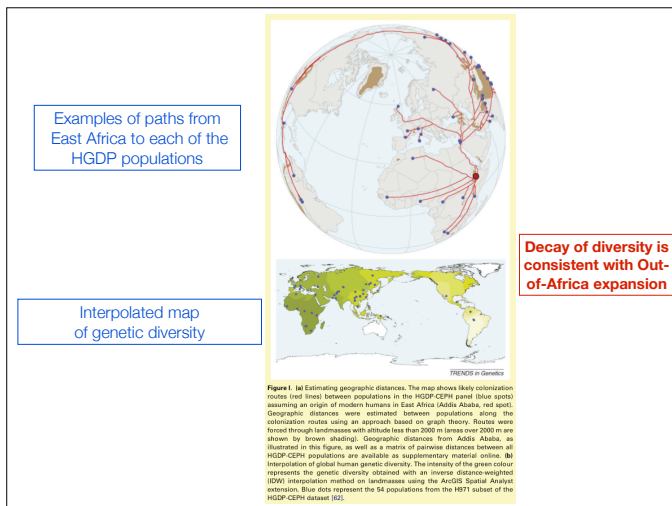
This Can Then Be Correlated With Diseases &/or Geographical Associations



## Global patterns of haplotype diversity



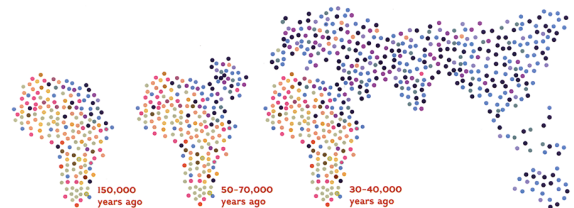
Decay of haplotype heterozygosity is consistent with a "serial bottlenecks" during Out-of-Africa expansion



## 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

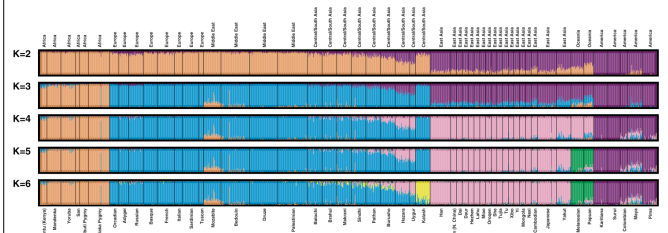


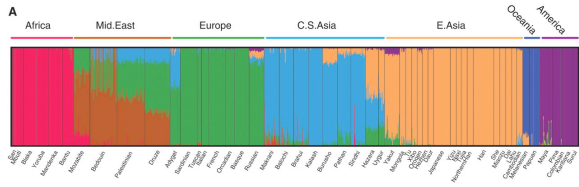
Fig. 1. Estimated population structure. Each individual is represented by a thin vertical line, which is partitioned into K colored segments that represent the individual's estimated membership fractions in K clusters. Black lines separate individuals of different populations. Populations are labeled 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 instead of Eurasia, 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?

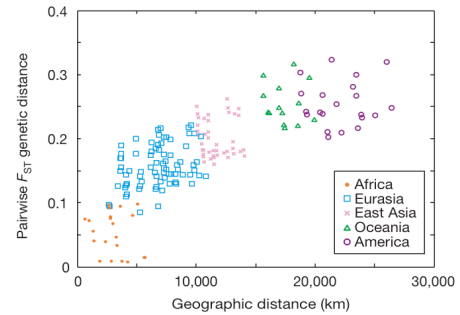
Results of **unsupervised** clustering algorithm on ~650,000 SNPs



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



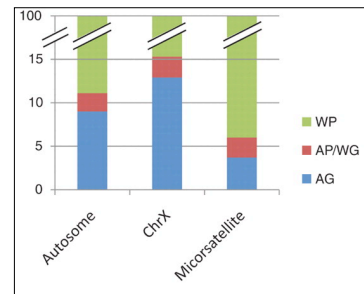
Africans	Europeans	Western Asians	Eastern Asians	Oceanians
1 Bantu	8 Orcadian	16 Bedouin	28 Han (S. China)	46 Melanesian
2 Mandinka	9 Argyle	17 Druze	29 Han (N. China)	47 Papuan
3 Yoruba	10 Russian	18 Palestinian	30 Dai	
4 San	11 Basque	31 Daur	32 Hui	
5 Mbuti pygmy	12 French	32 Hui	33 Lahu	
6 Biaka	13 North Italian	34 Mac	35 Oromo	
7 Mozabito	14 Sardinian	35 Oromo	36 She	
	15 Tuscan	36 She	37 Tuva	
		37 Tuva	38 Tu	
		38 Tu	39 Xibo	
		39 Xibo	40 Yi	
		40 Yi	41 Mongolian	
		41 Mongolian	42 Naei	
		42 Naei	43 Cambodian	
		43 Cambodian	44 Japanese	
		44 Japanese	45 Yakut	
		45 Yakut		

Perhaps clusters are due to impact of geographic barriers?

Sahara  
Tibetan Plateau  
Bering Strait  
Malay archipelago

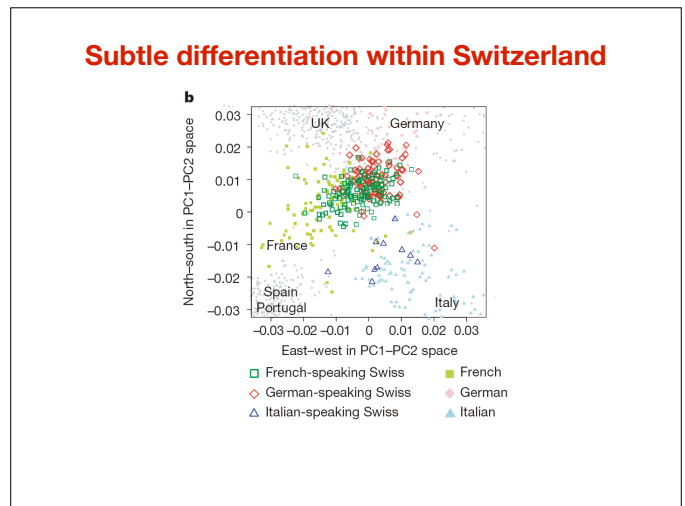
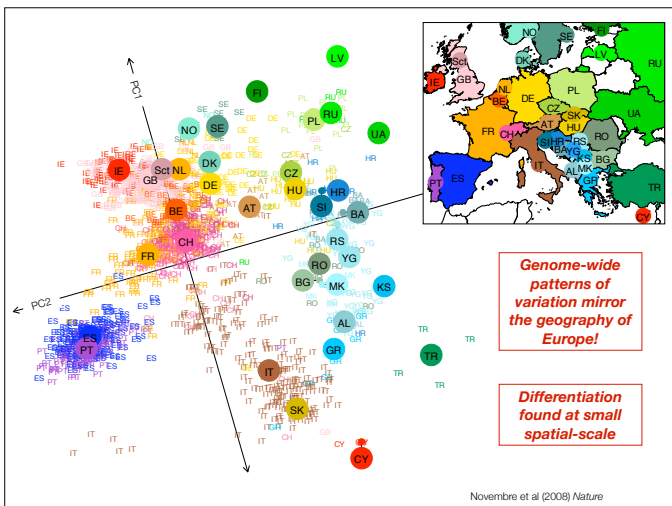
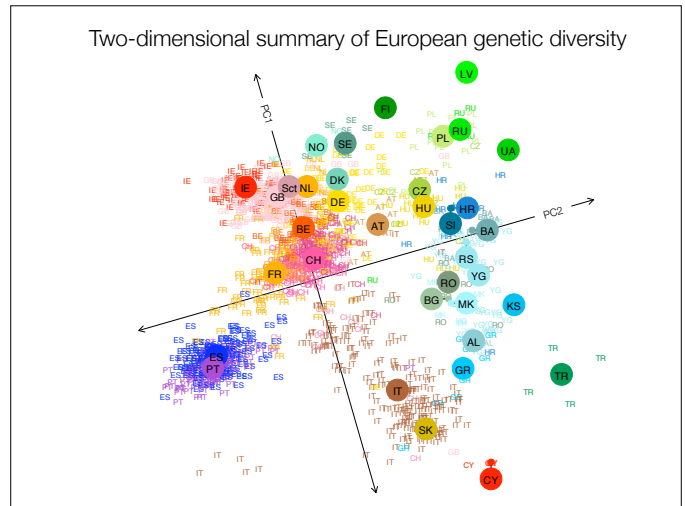
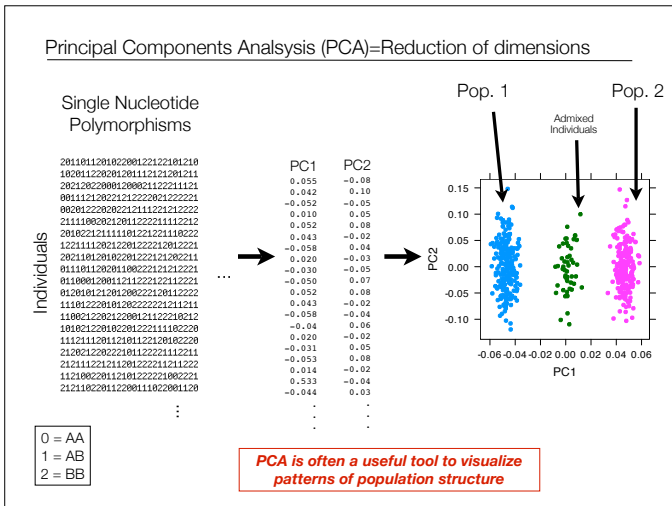
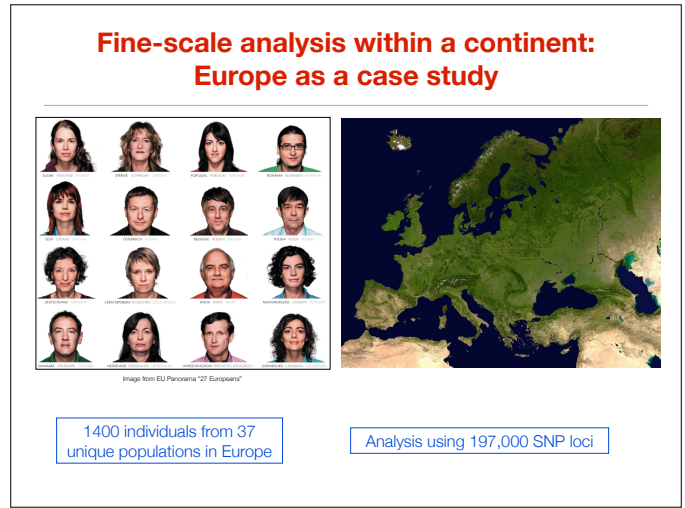
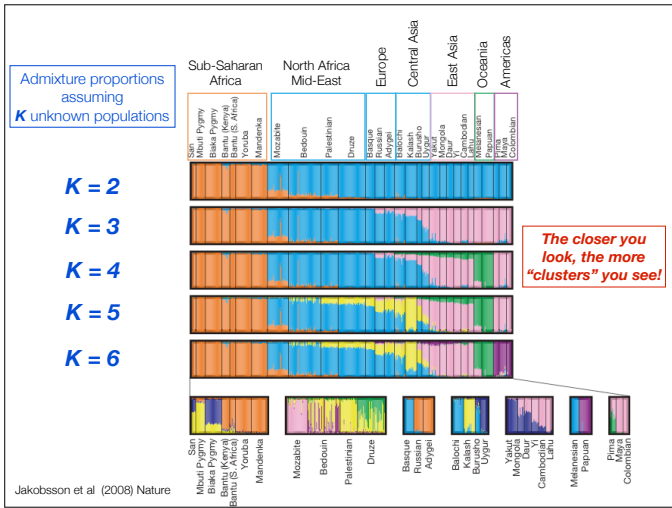
Or perhaps clusters are an artifact due to gaps in sampling?

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

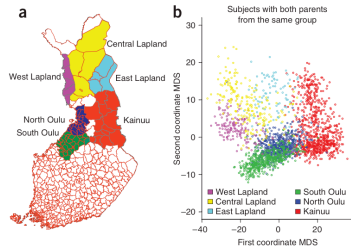


Within populations  
Among population / within groups  
Among groups

Li et al (2008) Science



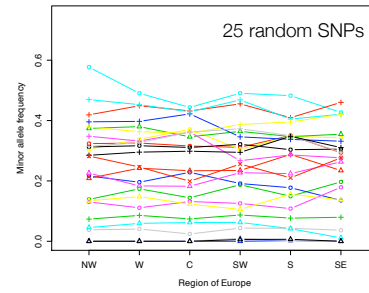
## Differentiation within Finland



**Figure 1** Linguistic/geographic groups of Northern Finland and their genetic signature. (a) Map of Finland with county boundaries. The subjects in NFC1965 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

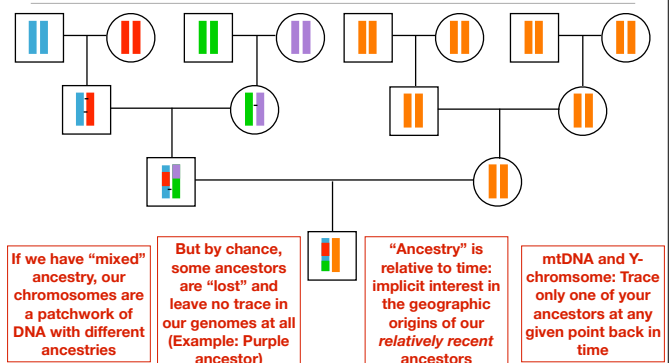
- At a global scale, support for clusters and clines can be observed
  - **Unresolved:** Do geographic barriers lead to global-scale “cluster” results or is it uneven sampling?
- With large numbers of SNP markers, patterns of differentiation are detectable even at small-scales within continents (although 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

## 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 is 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 was likely Native American.
  - (d) Not enough information to tell.

## What do we mean by ancestry?





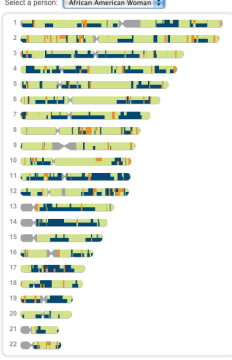
### ancestry painting

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 15**




**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.

- Africa 65%
- Europe 29%
- Asia 7%
- Not Genotyped

**Worldwide Examples**

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

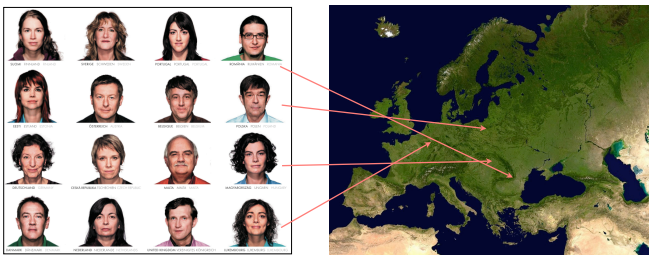


Tell Me About...

**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**

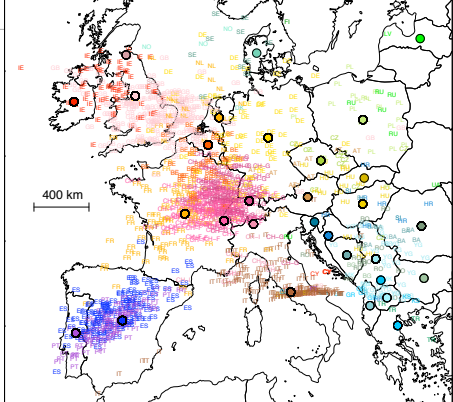
### Pushing the limits: Ancestry inference within continents



### 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

global similarity

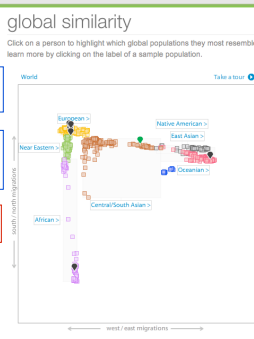
Click on a person to highlight which global populations they most resemble. Zoom in or learn more by clicking on the label of a sample population.

PCA of HGDP populations

Customer with mixed East Asian / European ancestry show in green.

Customer appears to be Central Asian!

Note: 23andme FAQ warns customer of this issue



### 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?

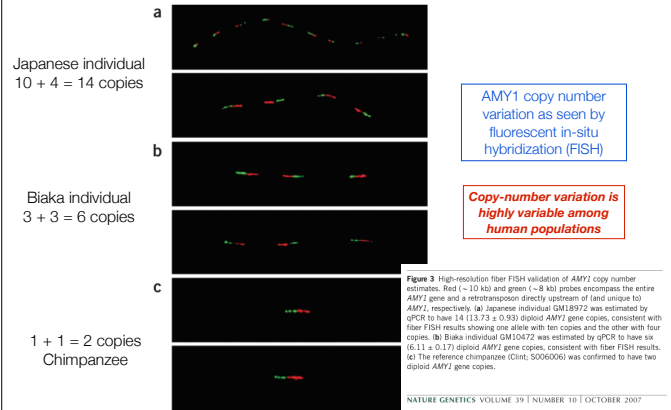
### What's your opinion:

- 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.

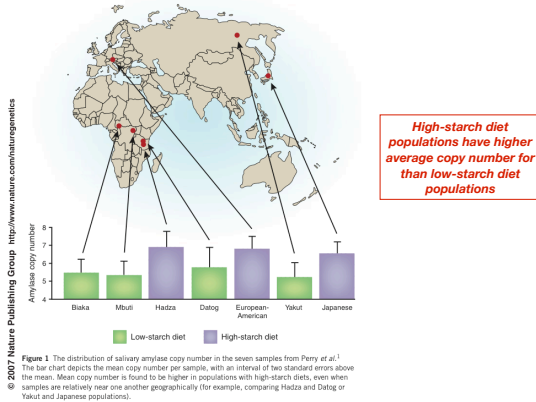
## Question:

- Are humans still evolving?
  - (a) Yes
  - (b) No

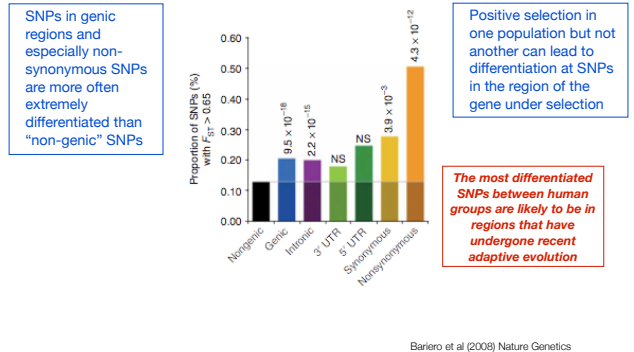
## Salivary amylase copy number variation



## Diet and Variation in salivary amylase copy number



## The impact of natural selection on population differences



Categories of genes that appear to have been under selection judging by extreme population differentiation

**Table 1** Genes showing the strongest signatures of positive selection

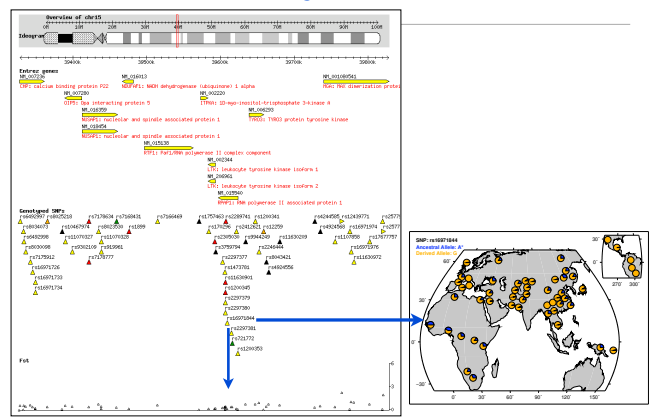
Phenotype category	Genes
Morphological traits (for example, skin pigmentation and hair development)	<i>ABCC11</i> , <i>EDAR</i> , <i>SLC45A2</i> , <i>PKP1</i> , <i>PLEKH44</i> , <i>SLC24A5</i>
Immune response to pathogens	<i>CEACAM1</i> , <i>CRI1</i> , <i>DUOX2</i> , <i>VHW2</i>
DNA repair and replication	<i>MFG</i> , <i>POLG2</i> , <i>TDP1</i>
Sensory functions (for example, olfaction and eye development)	<i>COL18A1</i> , <i>OR52K2</i> , <i>RP1L1</i>
Insulin regulation, metabolic syndrome (obesity, diabetes, hypertension)	<i>ALMS1</i> , <i>CEACAM1</i> , <i>ENPP1</i>
Various metabolic pathways (for example, ethanol, intestinal zinc and citrulline)	<i>ADH1B</i> , <i>ASS1</i> , <i>SLC39A4</i>
Miscellaneous	<i>FBXO31</i> , <i>RTTN</i> , <i>SPAG6</i>
Unknown	<i>ABCC12</i> , <i>ADAT1</i> , <i>AK127117</i> , <i>C17orf46</i> , <i>CBorf14</i> , <i>COLECD1</i> , <i>CPSF3L</i> , <i>DNAJC5B</i> , <i>DINH1</i> , <i>EIF2H</i> , <i>EXOC5</i> , <i>FAIM</i> , <i>CCDC142</i> , <i>FLJ37454</i> , <i>FXR1</i> , <i>GHSL2</i> , <i>KIAA0984</i> , <i>LAMB4</i> , <i>LOC64851</i> , <i>LIMCH1</i> , <i>PCDF1</i> , <i>PLEKHG4</i> , <i>POL33P</i> , <i>RNF135</i> , <i>SLC30A9</i> , <i>SYTL3</i> , <i>TEX15</i> , <i>TYC3P</i> , <i>VPS33B</i> , <i>ZNF646</i>

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 genetic SNPs with  $F_{ST} > 0.65$  is provided in **Supplementar 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>.

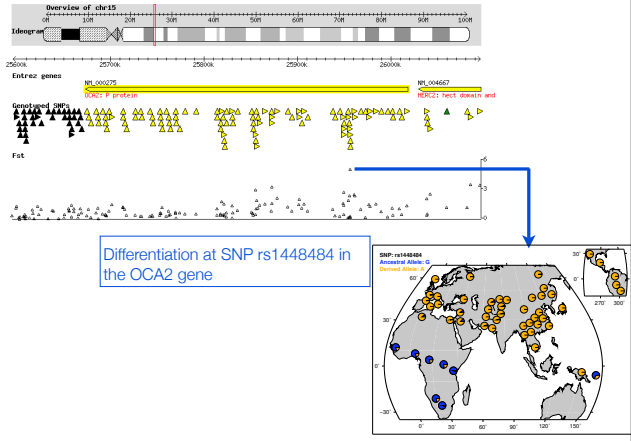
<sup>2</sup>These genes have not yet been attributed a HUGO-approved symbol. <sup>3</sup>These three genes are located in a linkage-disequilibrium block in chromosome 2. <sup>4</sup>These two genes are located in a linkage-disequilibrium block in chromosome 16.

Bariero et al (2008) Nature Genetics

## Differentiation at a random region of Chromosome 15



## Differentiation at around the OCA2 gene



Differentiation at SNP rs1448484 in the OCA2 gene

## OCA2 came up in Lecture 5...

- Was it a gene that is known to affect:
  - (a) Type II diabetes susceptibility
  - (b) Eye-color
  - (c) Intelligence
  - (d) Cystic-fibrosis

## SNPs in Human P Protein Gene Lead To Different Eye Colors

**OCA2**  
From SNPEdit

OCA2, the melanocortin 1-receptor gene (also known as the human P protein gene, or P170), is a gene associated with albinism and certain pigmentation effects in general with eye color, skin color, and hair color.

A large (3,000 individuals study) of Caucasians indicates that the following OCA2 variants, all located in the first intron of the gene, are preferentially linked to blue eye color inheritance: together, they form haplotypes that in some cases at least predict eye color with greater than 90% odds. (PMID: 17256393; OMIM: 203300;0001)

<http://www.ncbi.nlm.nih.gov/omim/gene/gene?term=203300;AlleleVariant0111>

The haplotypes are defined in order as listed above for these 3 SNPs, i.e. for example, the H1E haplotype refers to C/S/S/T/C/G/A/C/S/S/S/S/S/S/S/S/S/S. The correspondence between haplotypes (the two haplotypes in one individual) and the % of individuals with blue/green/grey/brown and brown eye color, respectively, was reported as follows for the most common diploypes (PMID: 17256393).

Haplotype	Eye Color	% of Individuals
H1E/H1E	Blue	28.0, 9.2
H1E/H1E'	Blue	2.1, 2.3, 3.9
H1E/H1E''	Blue	3.9, 3.1, 4.9
H1E/H1E'''	Blue	2.0, 2.1, 3.0
H1E'/H1E'	Blue	3.1, 3.6, 2.1
H1E'/H1E''	Blue	3.1, 3.6, 2.1
H1E'/H1E'''	Blue	3.1, 3.6, 2.1
H1E''/H1E''	Blue	3.1, 3.6, 2.1
H1E''/H1E'''	Blue	3.1, 3.6, 2.1
H1E'''/H1E'''	Blue	3.1, 3.6, 2.1

The haplotypes shown in **bold** letters represent the ones reported by the authors of this study to be most associated with brown eye color. Furthermore, the haplotypes shown above are as published, and the associated SNPs - which have since changed if as well - are not in the orientation shown in dbSNP.

More recently, a study of a large Danish family led to associations with 2 SNPs in a different region of OCA2, also linked to blue or brown eye color:

- rs1293322
- rs1293323

Further studies found different regions of the OCA2 gene to also be predictive of eye color:

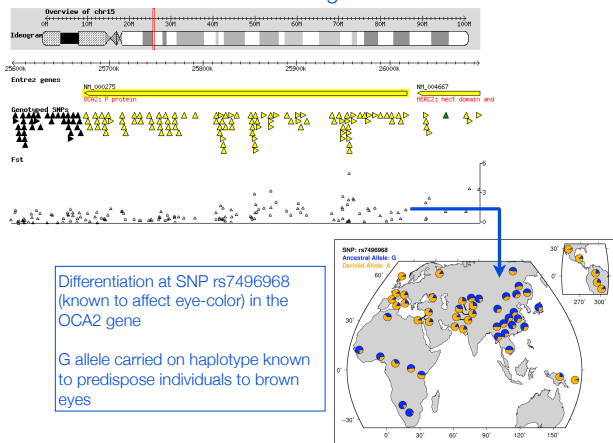
- HCA2 SNP rs180801 helps predict brown eye color. (PMID: 12163354; PMID: 1589196; OMIM: 203300.0011; <http://www.ncbi.nlm.nih.gov/omim/gene/gene?term=203300;AlleleVariant0111>)
- OCA2 SNP rs180807 may be associated with green/brown eye color in some populations, but not others. (PMID: 12163354; PMID: 1589196; OMIM: 203300.0012; <http://www.ncbi.nlm.nih.gov/omim/gene/gene?term=203300;AlleleVariant0121>)



Human Eye Color

**NOTE: OCA2 also affects skin color and hair color**

## Differentiation at around the OCA2 gene



Differentiation at SNP rs749698 (known to affect eye-color) in the OCA2 gene

G allele carried on haplotype known to predispose individuals to brown eyes

## Conclusions:

- 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