# Efficient algorithms for ascertaining markers for controlling for population substructure
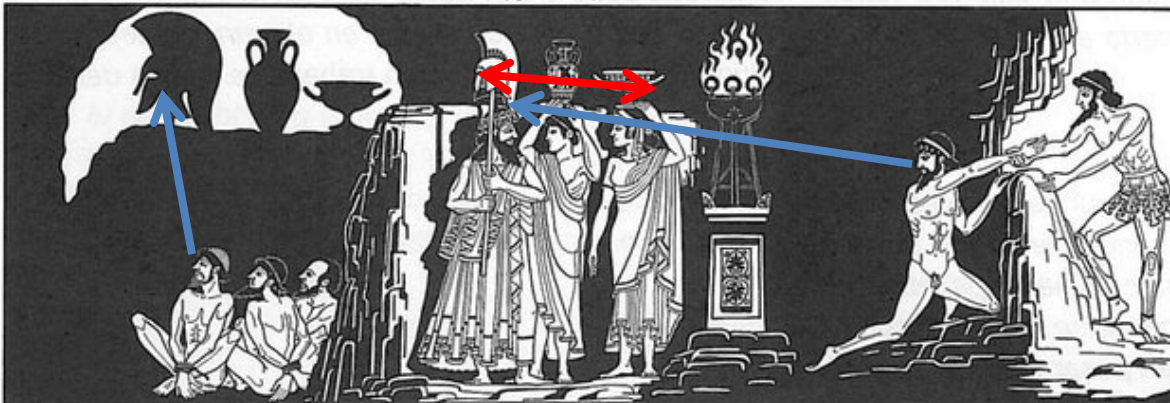
Oscar Lao
Department of Forensic Molecular Biology
Erasmus Medical Center (Rotterdam)
New Jersey 2009

# Workflow

1. Human population substructure

   - How to detect it?

   - How much?

   - Where does it come from?

2. Why does it matter?

3. Ancestry Sensitive Markers (ASMs) / Ancestry Informative Markers (AIMs)

   - Hypothesis driven. Particular individual clusters are preferred
     - ASMs
     - PhenoASMs

How much there is and how much can be detected. The two sides of the same coin

Plato's cave myth

# DETECTION

- STRUCTURE

- BAPS

- FRAPPE

- GENELAND

- PCA/MDS + K means

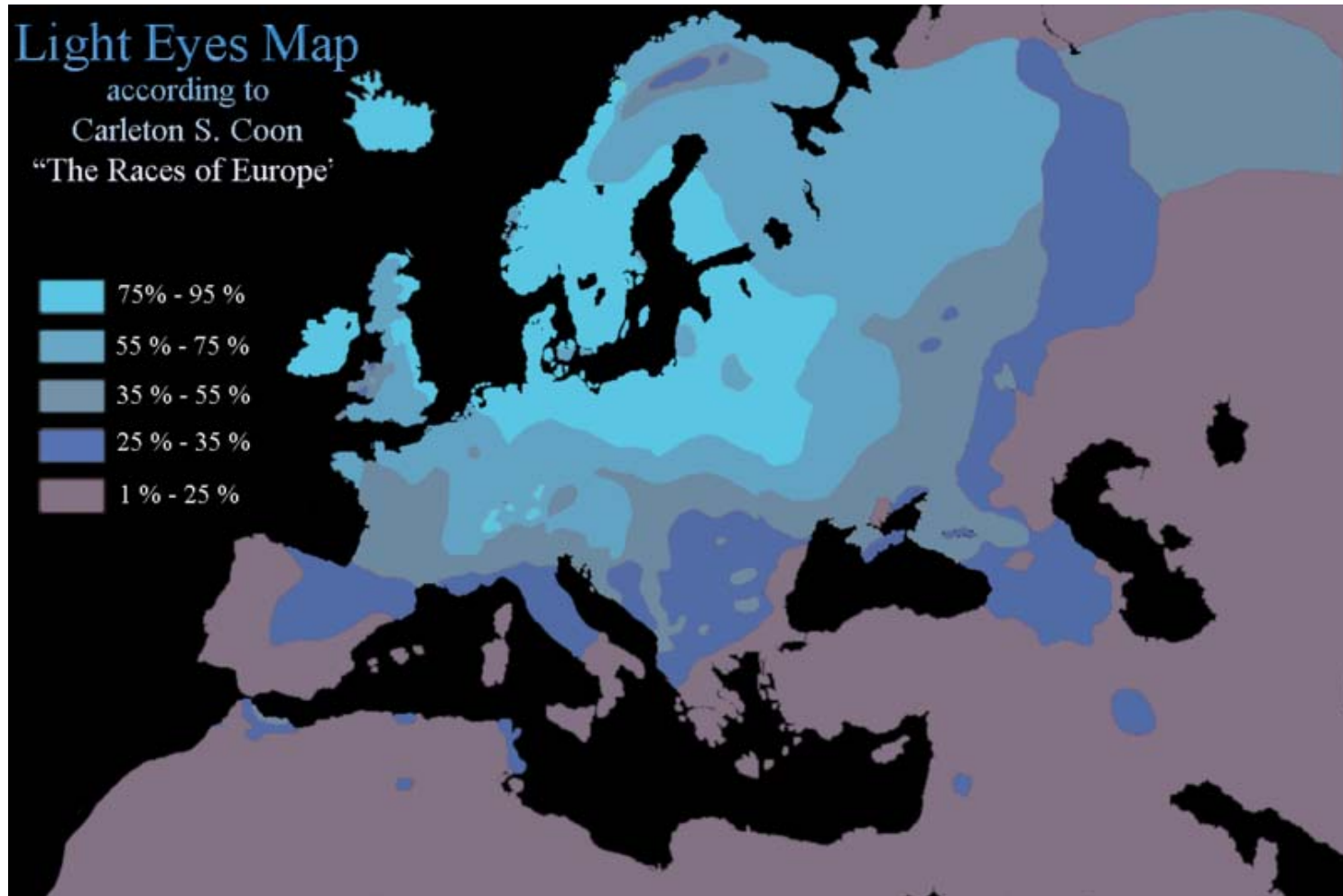- Neural Networks

- ...

Sometimes results
are NOT reproducible
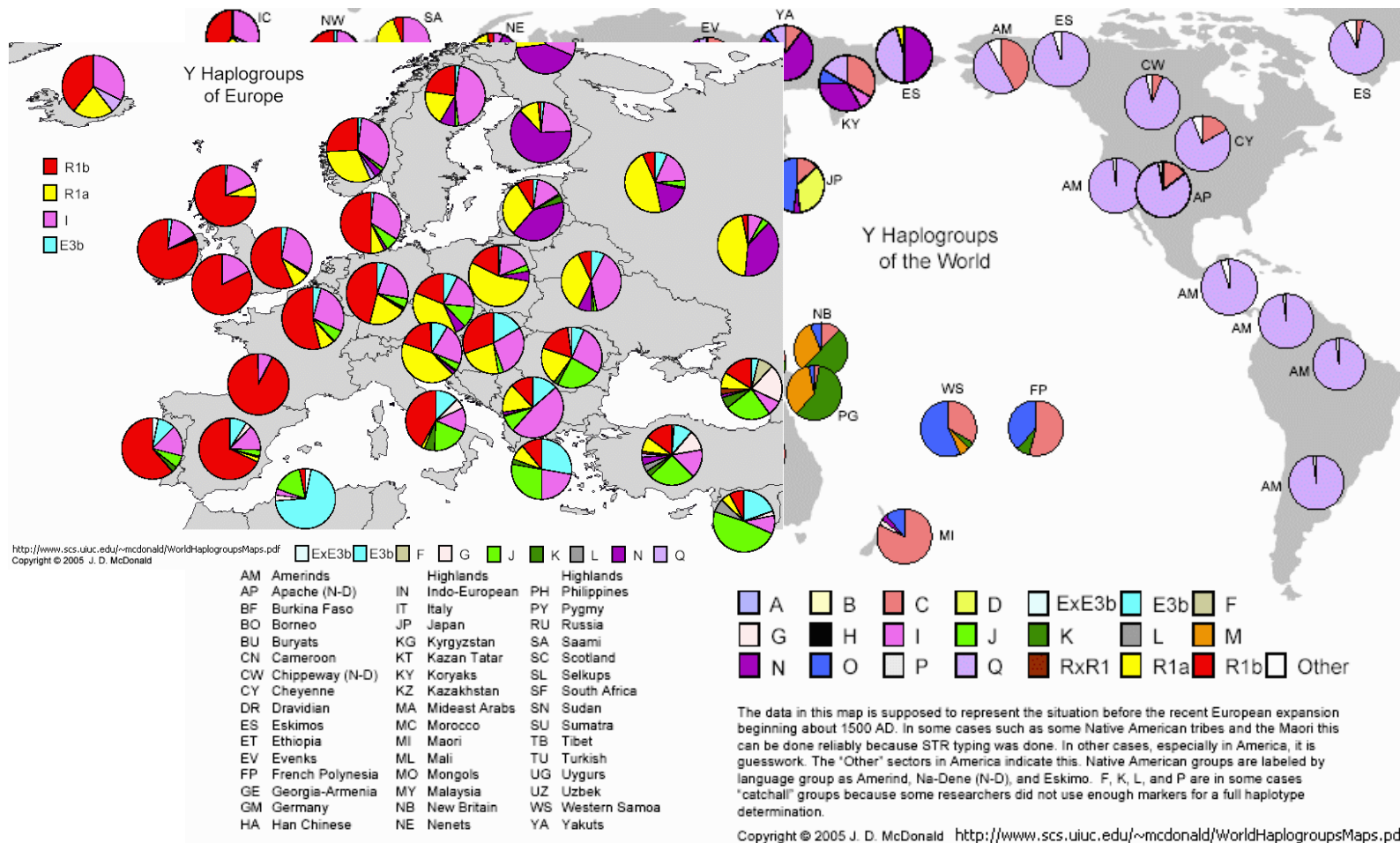
# HOW MUCH?

## Which type?

- Phenotype

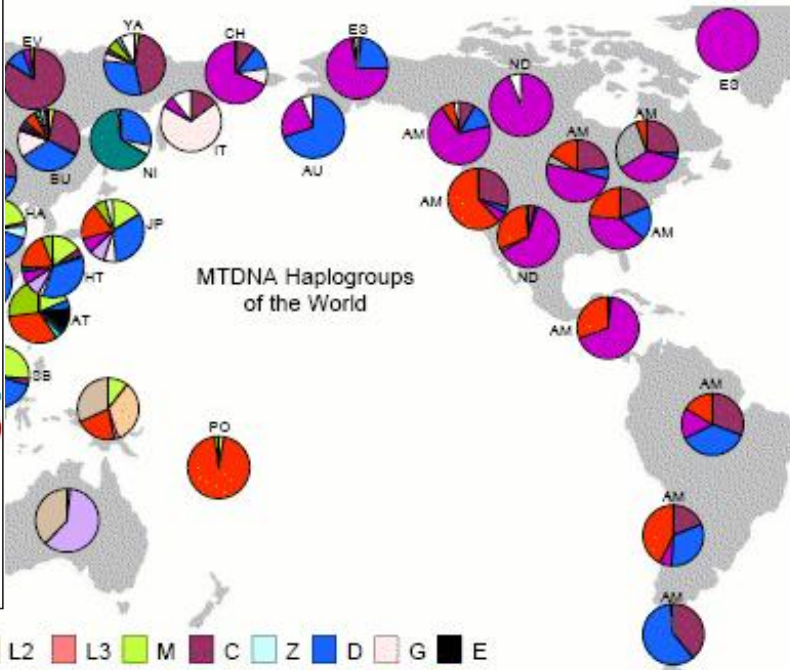- Genotype

    Y chromosome

    mtDNA

    Autosomal markers

## Where?

- Worldwide

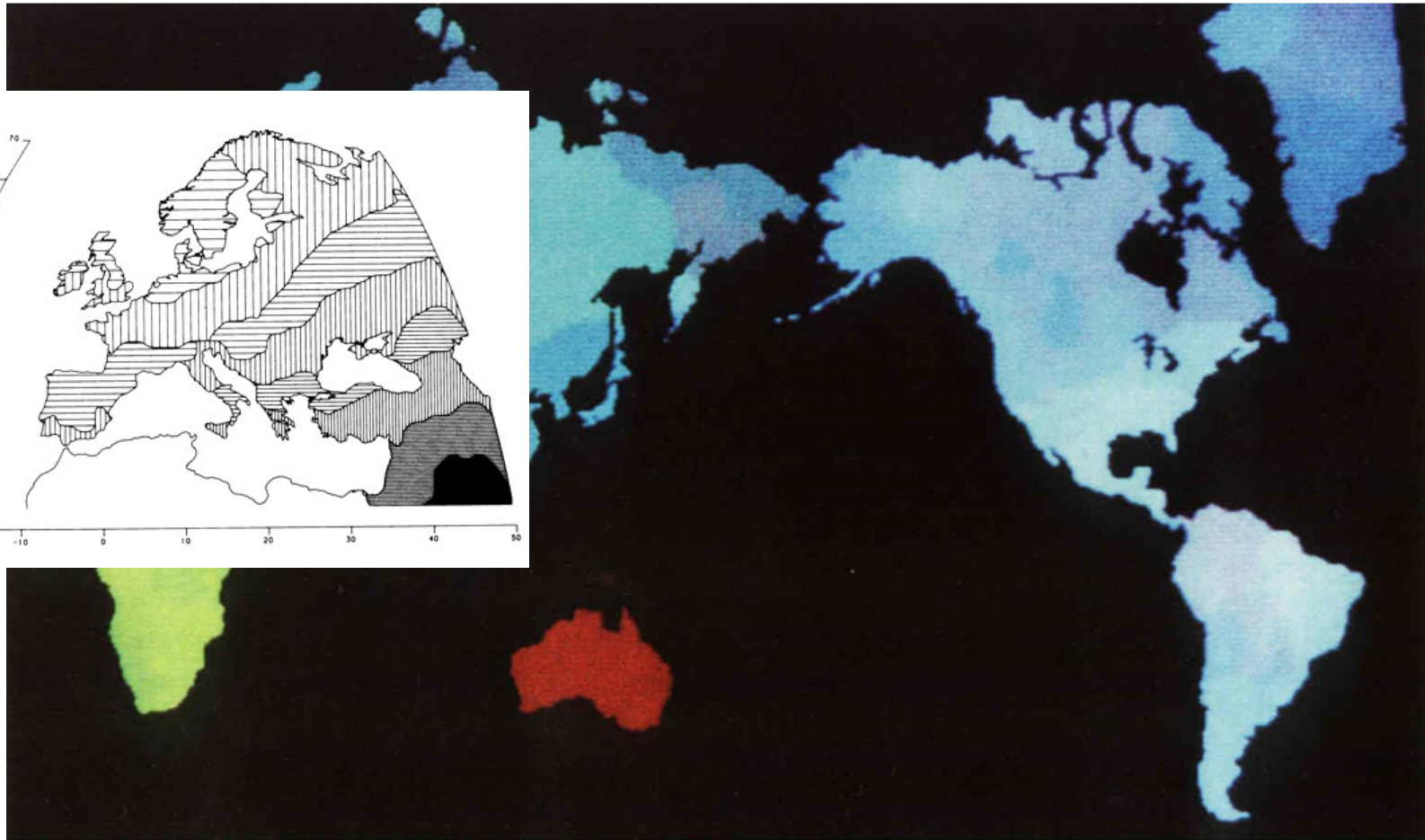- Regional (I will focus on Europe)

Light Eyes Map according to Carleton S. Coon "The Races of Europe"

MTDNA Haplogroups of the World

Specific tribes or locations are shown at left. Unlabelled pies are for general population in the area. African, American, and especially Polynesian areas are very large. The data in this chart is supposed to represent the situation before the recent European expansion beginning about 1500 AD. Assignments in Australia are somewhat iffy.
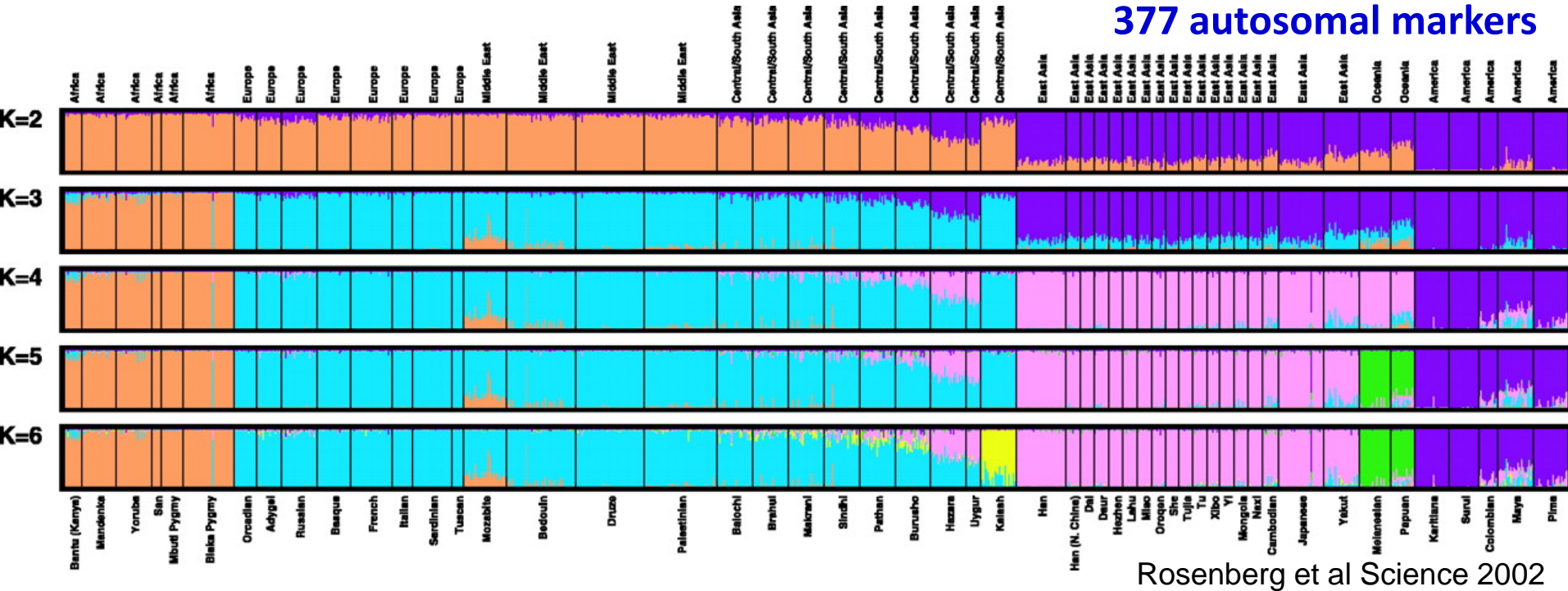
Copyright © 2005 J. D. McDonald

Cavalli-Sforza et al 1994

1064 samples
51 human populations of global distribution

Autosomal STRs

377 autosomal markers

Rosenberg et al Science 2002

993 autosomal markers

Africa | Europe | M. East | C/S-Asia | E-Asia | O | Ameri

Rosenberg et al Plos Genetics 2005

**A** 650,000 SNPs (FRAPPE)



550,000 SNPs (STRUCTURE)

Haplotypes

Li et al Science 2008





Jakobsson et al Nature 2008

23 populations

500 Affy Array

300,000 SNPs



Lao and Lu et al Current Biology 2008

Novembre et al Nature 2008

K = 2; Admixture



Correlation with latitude

**R² = 0.86**

Correlation with longitude

**R² = 0.01**

# World

# Europe



Anayet peak (2574 m), Pyrenees



Keukenhoof  garden(-2 m), Netherlands

Chr2. Comparison CEPH Europeans vs CHB Asians

EDAR (positive selection in Asians)

**Chromosome 2**

LCT (positive selection in North European populations)

Lao and Lu et al Current Biology 2008

Cavalli-Sforza & Feldman Nature Genetics 2003



Simoni et al AJHG 2000

- Selective pressures within the species (locus specific)



Lactose tolerance
Malaria resistence
Human pigmentation
…

- Population substructure & pigmentation (5 SNPs)



MDS

Europe
Africa
Middle East
Oceania
America
C/SAsia
EAsia
N Africa

Biasutti skin
pigmentation
units

Lao et al Ann Hum Genet 2007

Plato's cave myth



CHANGE THE ALGORITHM
FOR DETECTING
POPULATION
SUBSTRUCTURE

Plato's cave myth



INCREASE THE
RESOLUTION TO
SEE  THE OBJECTS

# AIMs/ASMs

- Markers that capture most of the genetic ancestry
  - Estimate ancestry
  - Reduce the number of markers to test for genetic homogeneity
    - Time cost (clustering algorithms can be extremely computational intensive)
    - Economical cost (i.e exclude individuals BEFORE doing the GWA)

- Based on the existing diversity between individuals (i.e Paschou et al 2008)

- Based on predefined groups of individuals
  - No phenotype linked
    - Large Genetic distances
    - Signals of positive selection
  - Phenotype linked
    - Covariates with the phenotype of interest

- Use a statistic to quantify the amount of differentiation between populations

- Compute the OVERAL non-redundant amount of In between set of SNPs

- Take the best combination of markers from all the possible combinations

- Repeat the process until the information of the set of markers is maximum

*informativeness for assignment*

$$I_n(Q;J) = \sum_{j=1}^{N} \left( -p_j \log p_j + \sum_{i=1}^{K} \frac{p_{ij}}{K} \log p_{ij} \right)$$



Am J Hum Genet. 2003 Dec;73(6):1402-22

- How much information a marker contains about the ancestry of one individual (measured in *nats*)

- Ranges from 0 to the natural logarithm of the number of clusters and it is proportional to the number of differentiated clusters

- Computes the **non-redundant** amount of information when considering more than one marker

- Requires computing the frequency of **ALL** the allelic combinations when considering more than 1 locus

- Problem: The number of combinations increases exponentially with the number of markers.

  – Number of allelic combinations considering 50 SNPs:

$$2^{50} = 1,125,899,906,842,624$$

$$I_n(Q;J) = \sum_{j=1}^{N}\left(-p_j \log p_j + \sum_{i=1}^{K}\frac{p_{ij}}{K}\log p_{ij}\right)$$

$$I_n(Q;J) = \sum_{j=1}^{N}\left(\overline{H_j} - \sum_{i=1}^{K}\frac{H_{ij}}{K}\right)$$
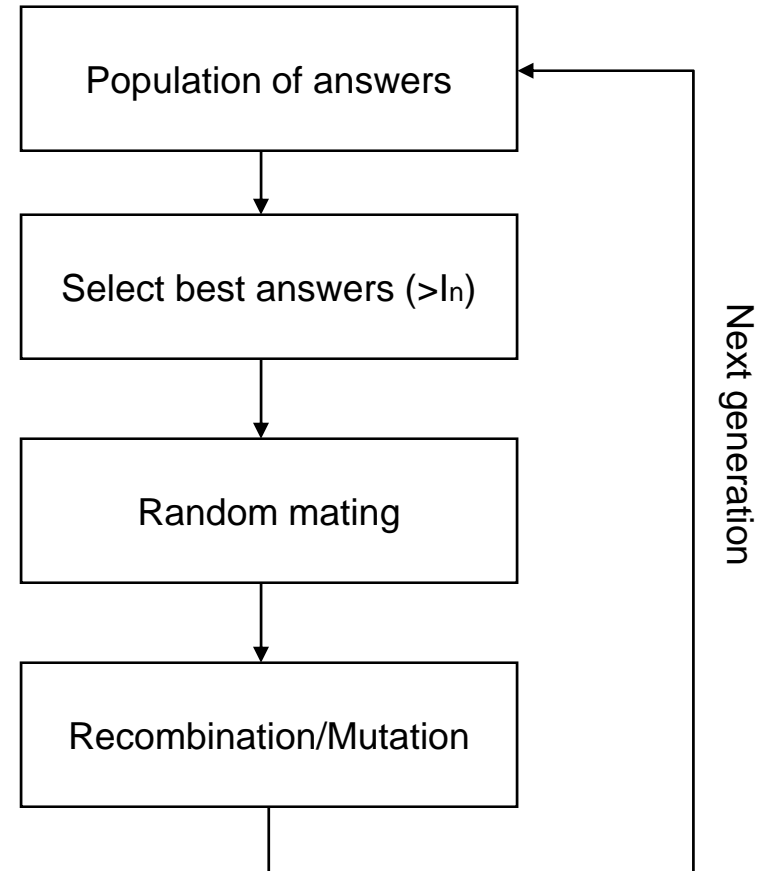
$$H \approx \frac{1}{N}\sum_{i=1}^{N}\ln(p)$$

By applying the Asymptotic Equipartition Property of Entropy

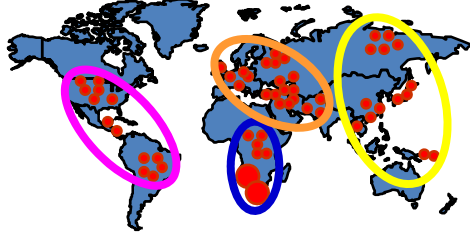- Problem: Considering 8,000 markers, ascertaining the best set of 50 markers requires computing :

$$N_{combinations} = \frac{8,000!}{50!(8,000-50)!} \approx 4 \times 10^{130}$$

Population of answers

Select best answers (>$I_n$)

Random mating

Recombination/Mutation

Next generation

## CEPH-HGDP panel

1064 samples
51 human populations of global distribution

## YCC-panel

76 human individuals
21 sampling localities

**Reproducibility of geographic structure in a different dataset**

SNP ascertainment
(10 SNPs)

(10 SNPs)

**Test for signatures of positive selection (EHH test)**

## Perlegen Database

3 Human populations
~1,500,000 SNPs
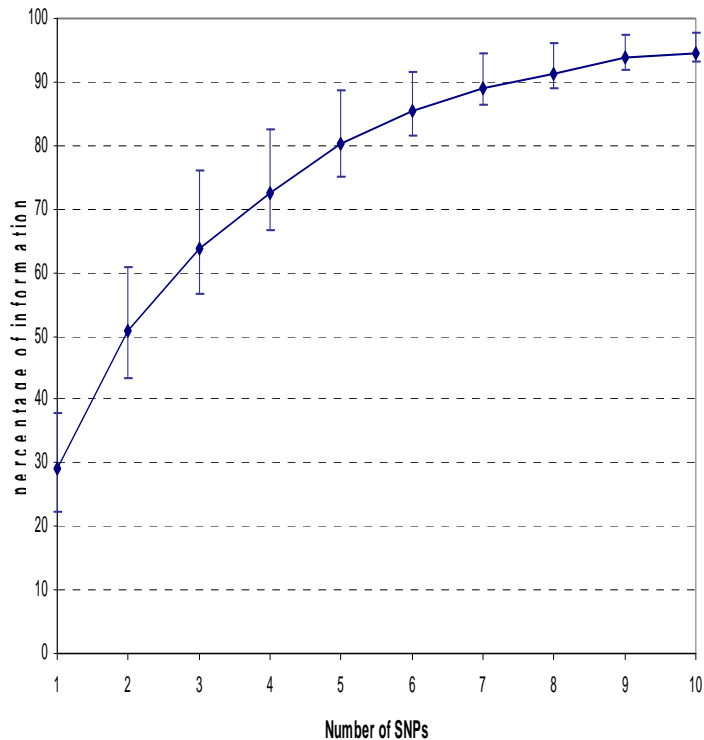(most informative 5 SNPs)

10k Affymetrix Array
(~9000 SNPs after excluding X-SNPs & missing SNPs)

Lao et al. Am J Hum Genet. 2006 Apr;78(4):680-90

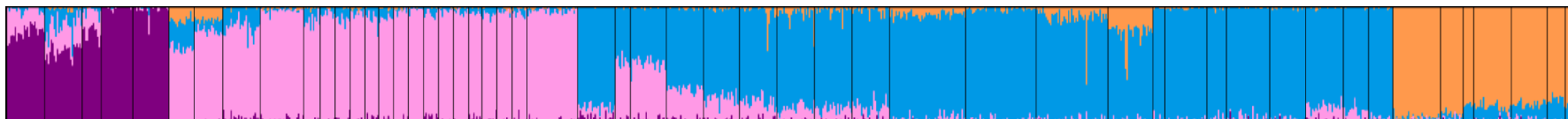**The genetic algorithm was applied increasing every time the number of selected SNPs**

**Selected SNPs in the final 10 SNPs run**



| Marker name | Chromosome | Gene name | $I_N$ (%) from 4 groups YCC panel | $I_N$ (%) from 7 groups CEPH-HGDP |
|---|---|---|---|---|
| rs722869 | 14 | VRK1 | 29.066 | 7.960 |
| rs1858465 | 17 | | 25.637 | 9.228 |
| rs1876482 | 2 | LOC442008 | 24.589 | 10.290 |
| rs1344870 | 3 | | 22.810 | 11.074 |
| rs1363448 | 5 | PCDHGB1 | 19.418 | 4.552 |
| rs952718 | 2 | ABCA12 | 18.739 | 9.472 |
| rs2352476 | 7 | | 18.317 | 5.603 |
| rs714857 | 11 | | 18.083 | 6.157 |
| rs1823718 | 15 | | 17.845 | 5.451 |
| rs735612 | 15 | RYR3 | 14.315 | 5.530 |

Lao et al. Am J Hum Genet. 2006 Apr;78(4):680-90

**993 autosomal markers**



**10 SNPs**  **No admixture**



**Admixture**
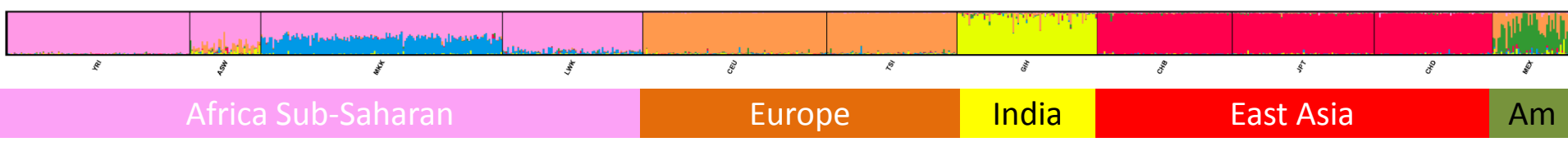


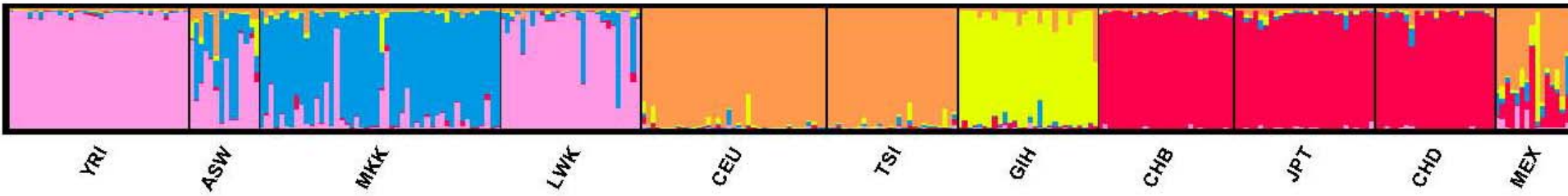| Ameri | O | E-Asia | C/S-Asia | Europe | M. East | Africa |

Lao et al. Am J Hum Genet. 2006 Apr;78(4):680-90

K = 6 (1000 (randomly ascertained) markers, Admixture, 10,000 burning, 10,000 retained simulations)



| Africa Sub-Saharan | Europe | India | East Asia | Am |

K = 5 (50 markers, Admixture,  500,000 burning, 500,000 retained simulations)



YRI  ASW  MKK  LWK  CEU  TSI  GIH  CHB  JPT  CHD  MEX

K = 6 (100 markers, Admixture,  100,000 burning, 100,000 retained simulations)



YRI  ASW  MKK  LWK  CEU  TSI  GIH  CHB  JPT  CHD  MEX

25 ascertained markers. PCA

- # CEPH



550,000 SNPs

K = 5 (50 ascertained markers, Admixture, 500,000 burning, 500,000 retained simulations)

Real geographic location
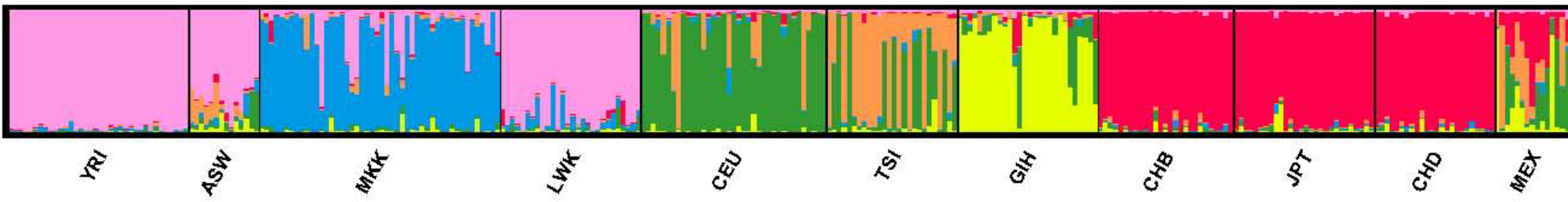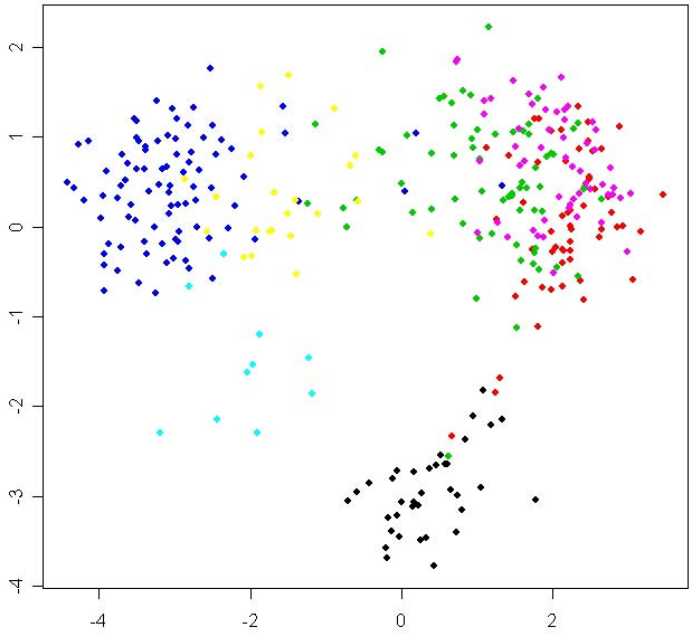
K = 2 (5000 random markers, Admixture,  10,000 burning, 10,000 retained simulations)



K = 3 (500 ascertained markers, Admixture,  10,000 burning, 10,000 retained simulations)
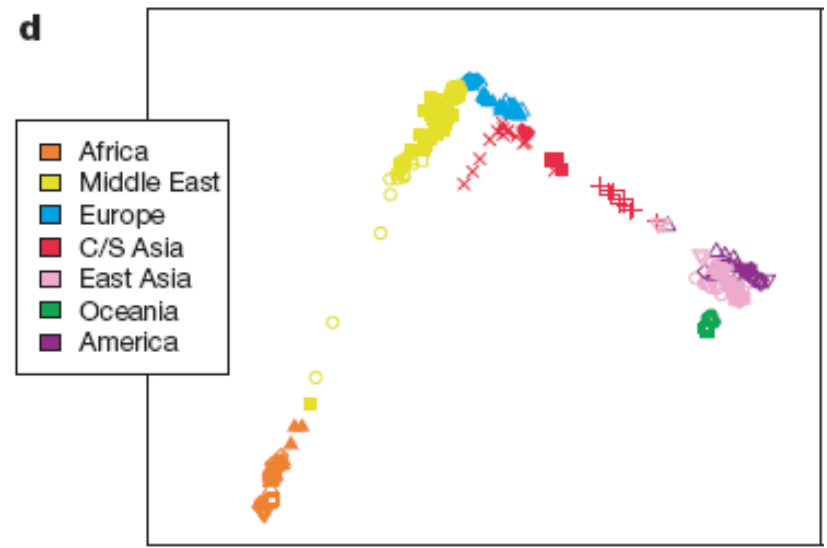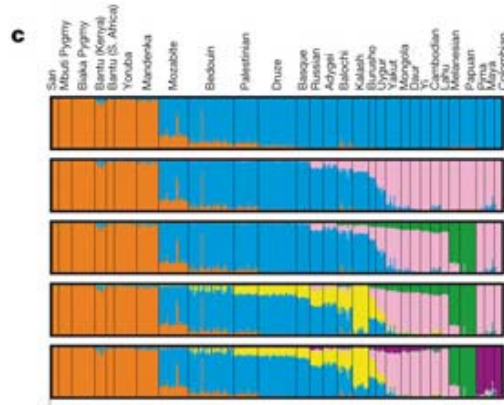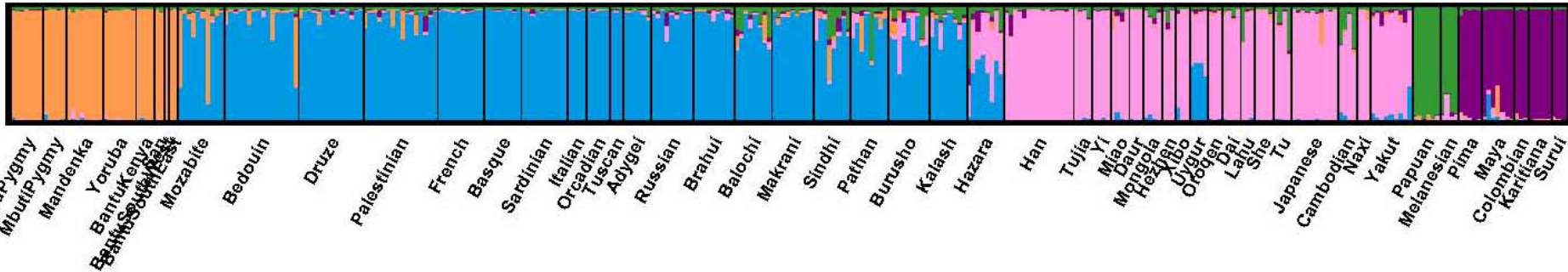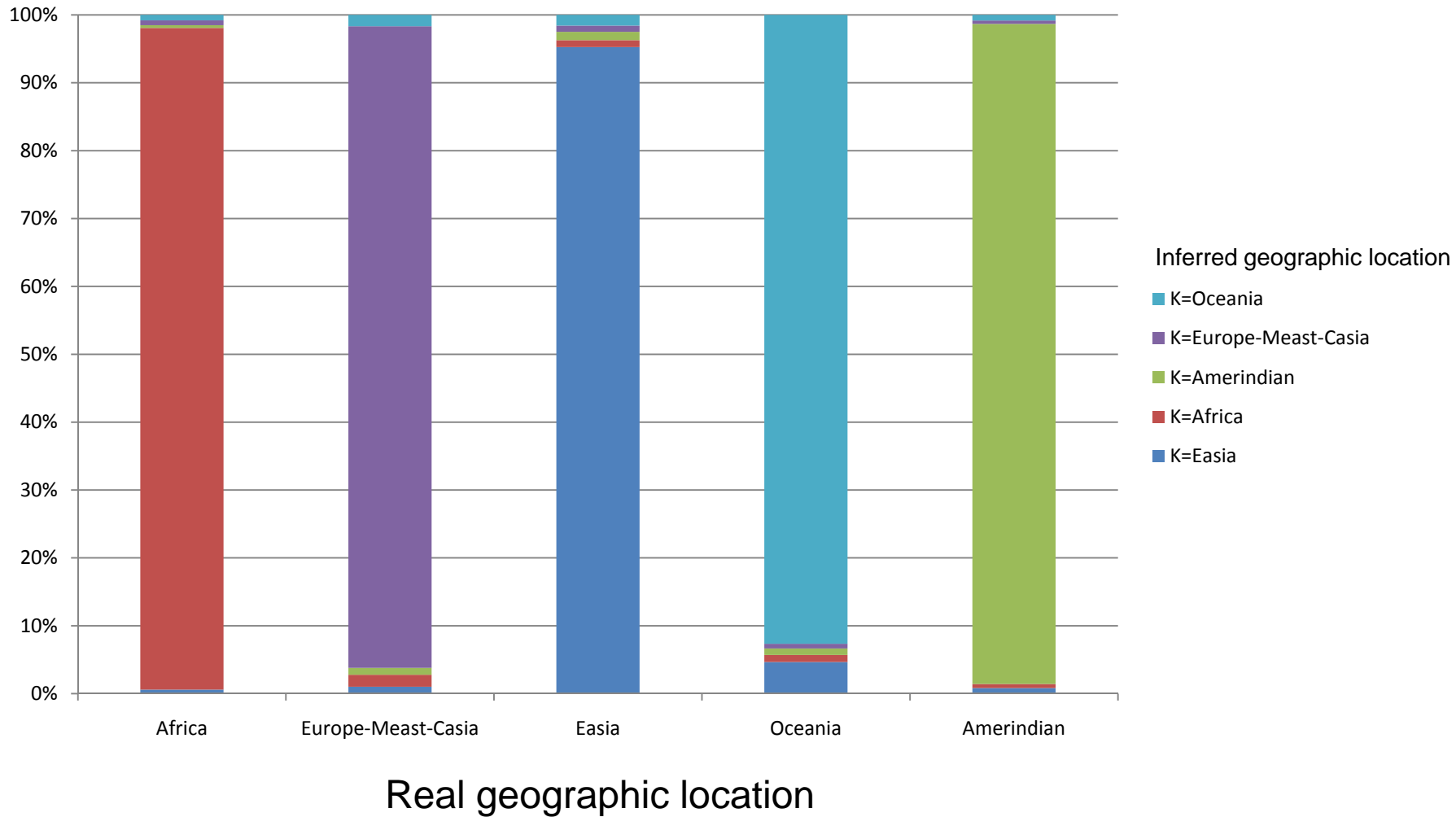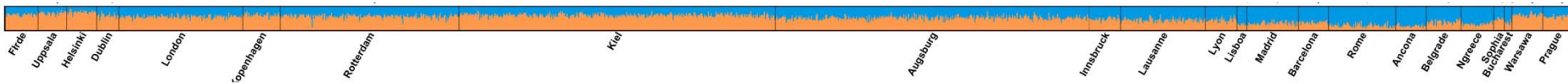
## Association between OCA_HERC2 region and iris color adjusted for ancestry sensitive markers

- Recall
  - Population substructure is only a problem when PHENOTIPIC and GENOTYPIC variation covariates
  - Why not ascertaining markers that are associated to the particular spatial pattern of the phenotype?

$$I_n(Q;P \mid J) = I_n(Q;P;J) - I_n(Q;J)$$

*"Amount of information of the phenotype (P) conditional on the genotype (J): How well could we correctly classify one individual given that we know his phenotype if we already know his genotype in a particular locus"*

Lactose intolerance in Europe

# Replication of signals from recent studies of Crohn's disease identifies previously unknown disease loci for ulcerative colitis

Andre Franke[1,5], Tobias Balschun[1,5], Tom H Karlsen[2], Jürgen Hedderich[3], Sandra May[1], Tim Lu[3], Dörthe Schuldt[1,4], Susanna Nikolaus[4], Philip Rosenstiel[1], Michael Krawczak[3] & Stefan Schreiber[1,4]
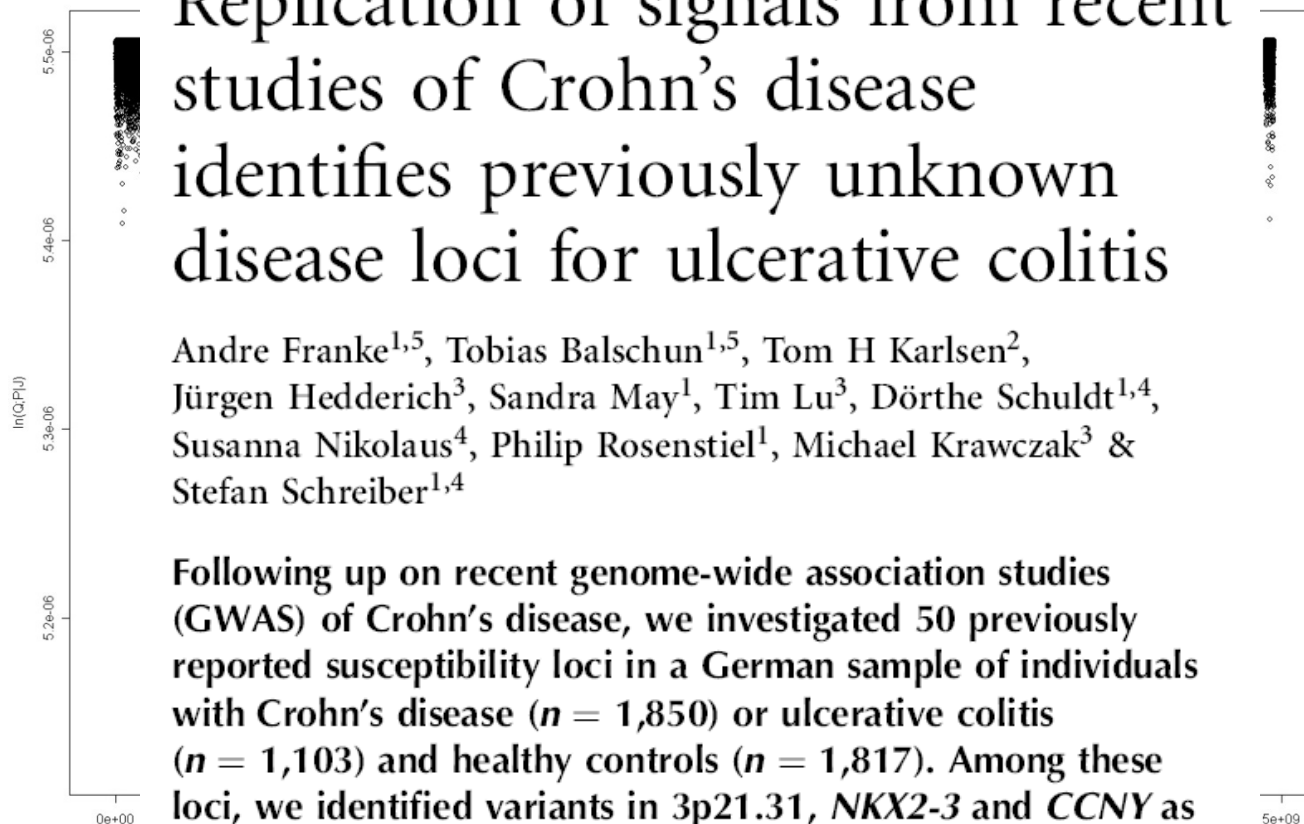
Following up on recent genome-wide association studies (GWAS) of Crohn's disease, we investigated 50 previously reported susceptibility loci in a German sample of individuals with Crohn's disease ($n = 1,850$) or ulcerative colitis ($n = 1,103$) and healthy controls ($n = 1,817$). Among these loci, we identified variants in 3p21.31, *NKX2-3* and *CCNY* as susceptibility factors for both diseases, whereas variants in *PTPN2*, *HERC2* and *STAT3* were associated only with ulcerative colitis in our sample collection.

| | AA | AB | BB | Marginal phenotype |
|---|---|---|---|---|
| C | P(AA)P(C\|AA) | P(AB)P(C\|AB) | P(BB)P(C\|BB) | ∑P(g)P(C\|g) |
| D | P(AA)P(D\|AA) | P(AB)P(D\|AB) | P(BB)P(D\|BB) | ∑P(g)P(D\|g) |

- Update $\theta$ with a Metropolis algorithm
- Update the covariance matrix of the proposal distribution by means of a *"quasi-perfect adaptive MCMC"* (Andrieu and Atchade)
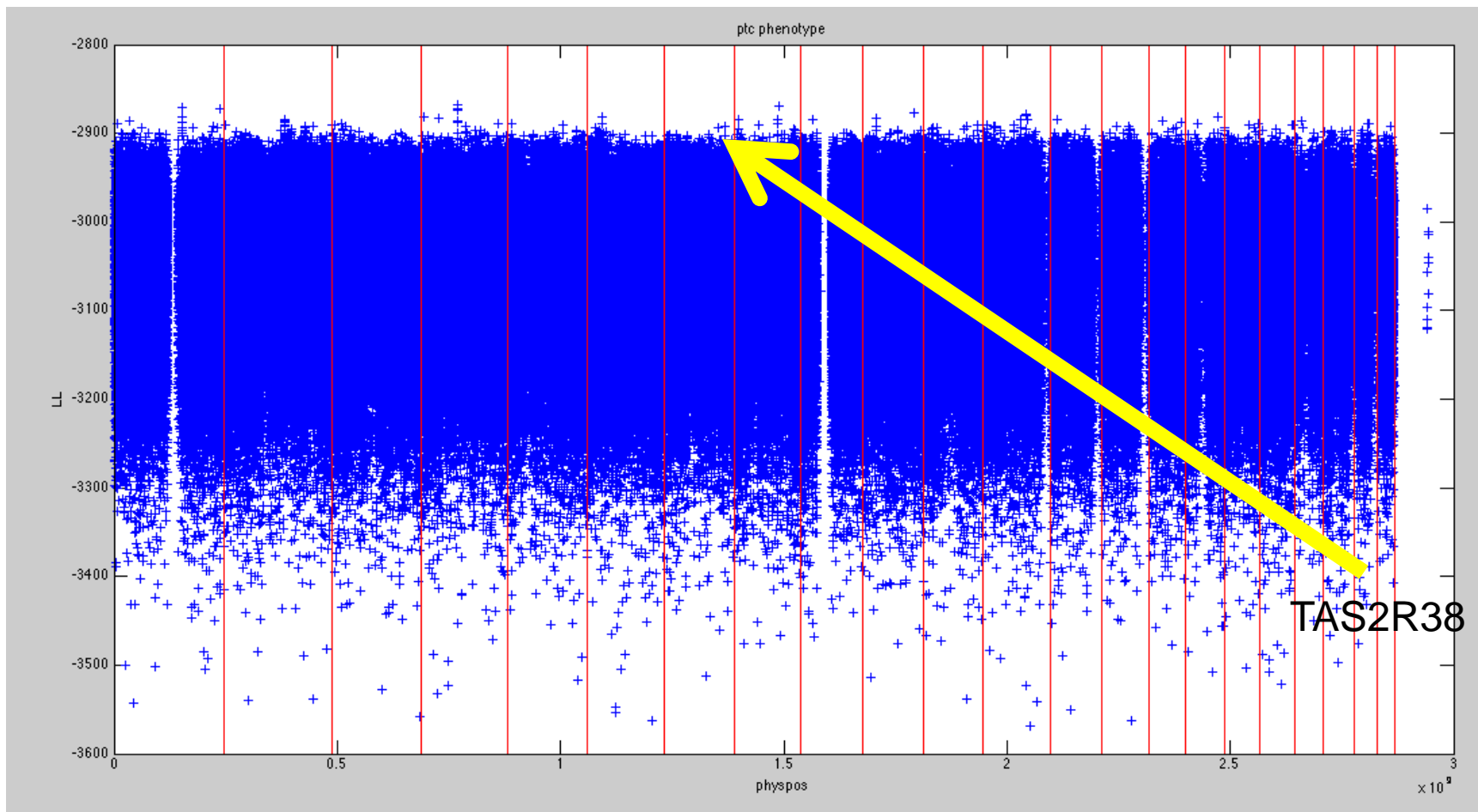- Compute the harmonic mean of the likelihood in order to obtain a rough estimate of P(M|D)

eye color phenotype

TAS2R38

TAS2R38

# Conclusions

- Low to moderate human population differentiation

- Mainly associated to geography

- No sharp discontinuities, except in particular genomic regions (selection?)

- Results depend on the clustering algorithm

- ASMs can improve the detection of population substructure

- $I_n$ is a good statistic for ascertaining markers to differentiate predefined populations

- If a prior definition of a population is used, ASMs will tend to differentiate such population, independently of the biological meaning

- PhenoASMs as the next level of ASMs?

# In collaboration with

M. Balascakova, C. Becker, J. Bertranpetit, L.A. Bindoff, D. Comas, U. Gether, C. Gieger, G. Holmlund, A. Kouvatski, M. Macek, I. Mollet, M. Nelson, P. Nuernberg, W. Parson, R. Ploski, A. Ruether, A. Sajantila, S. Schreiber, A. Tagliabracci, A. Uiterlinden, T. Werge, and E. Wichmann.

# Acknowledgements



Tim Lu



Manfred Kayser



Michael Krawczak

Andreas Wollstein

Petros Drineas

Peristeia Paschou

# *Thank you very much!*