

Human variations and genetic diseases

Michael Schatz

October 30 – Lecture 17

EN.601.452 Computational Biomedical Research

AS.020.415 Advanced Biomedical Research



Assignment 2: Genome Assembly

Assignment Date: Wednesday, October 25, 2017

Due Date: Monday, November 6, 2017 @ 11:58pm

Assignment Overview

In this assignment, you will explore a few properties of the sequencing data. You can either submit your results in a jupyter notebook, or as a single PDF document. Feel free to sketch the figures by hand, and then include a photograph of your solution. I encourage you to discuss your solutions with other members of the class, but everyone should submit their own write up. You are allowed to use the notes from class, and notes found online to help you work through the problem.

Here are a few helpful resources:

- [Python 2 reference](#)
- [Jupyter notebooks](#)
- [Matplotlib and Gallery](#)
- [Numpy and Scipy](#)

Question 1. Read coverage (10pts)

1a. The cichlid fish genome is 1 Gbp. Approximately how many 100bp reads should we sequence so that we expect at least 99.99% of the genome will be sequenced at least 40 times? (hint: show your work)

1b. Sketch the expected coverage distribution for this number of reads; be sure to clearly label the mean coverage, and how 40 fold coverage relates to the mean. (hint: in a normal distribution, 68.2% of the data will be within 1 standard deviation, 95.4% within 2, 99.7% within 3, and 99.9% within 4)

Question 2. de Bruijn graph construction (10pts)

2a. Draw the de Bruijn graph for the following reads using $k=3$ (assume all reads are from the forward strand, no sequencing errors, complete coverage of the genome)

ATTG
ATTG
GATT
CTTA
GATT
TATT
TTAT

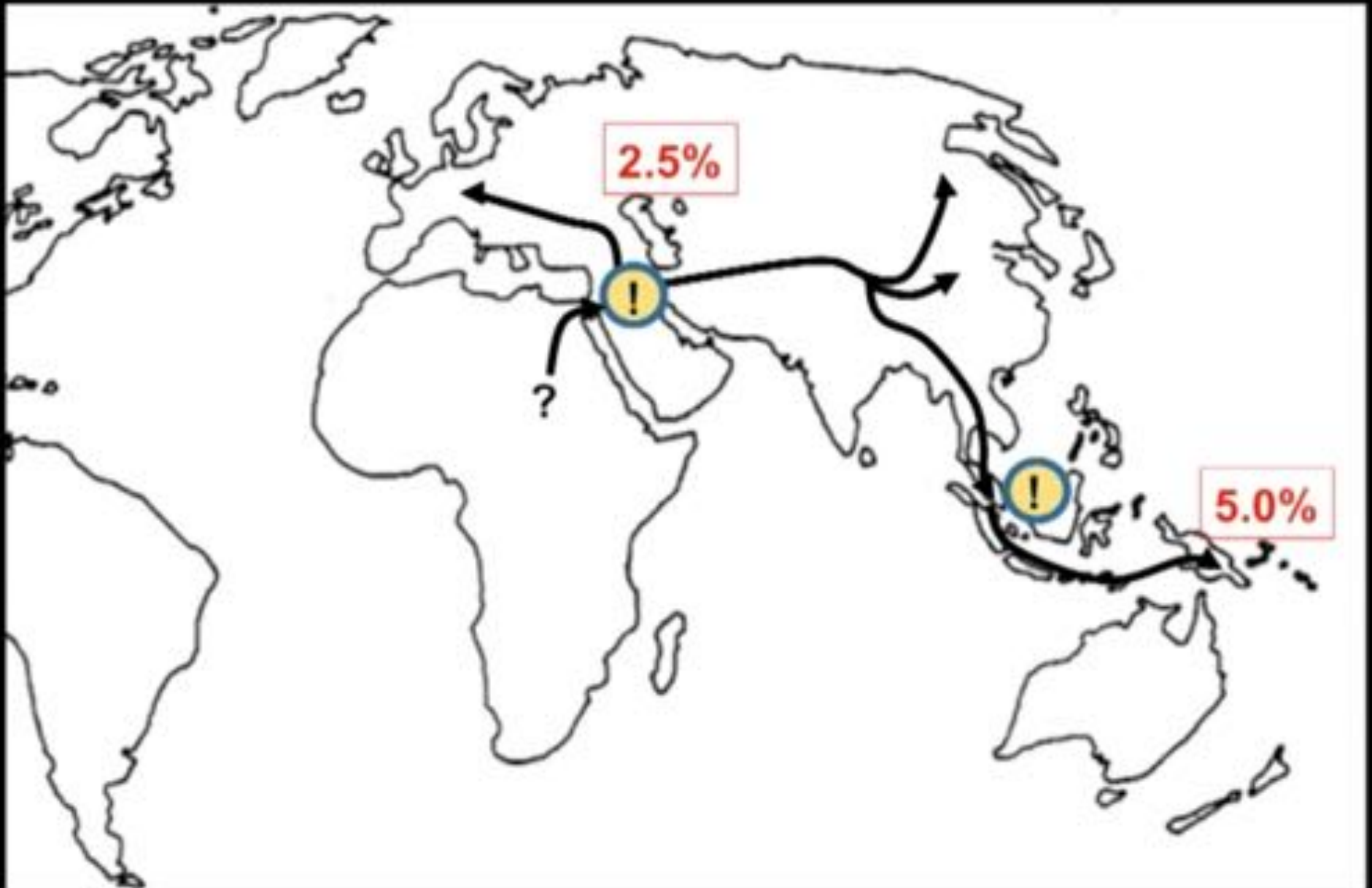
Sequencing ancient genomes

Janet Kelso

Max-Planck Institute



Migration of ancient hominids





Part I: Clustering Refresher

Clustering Refresher

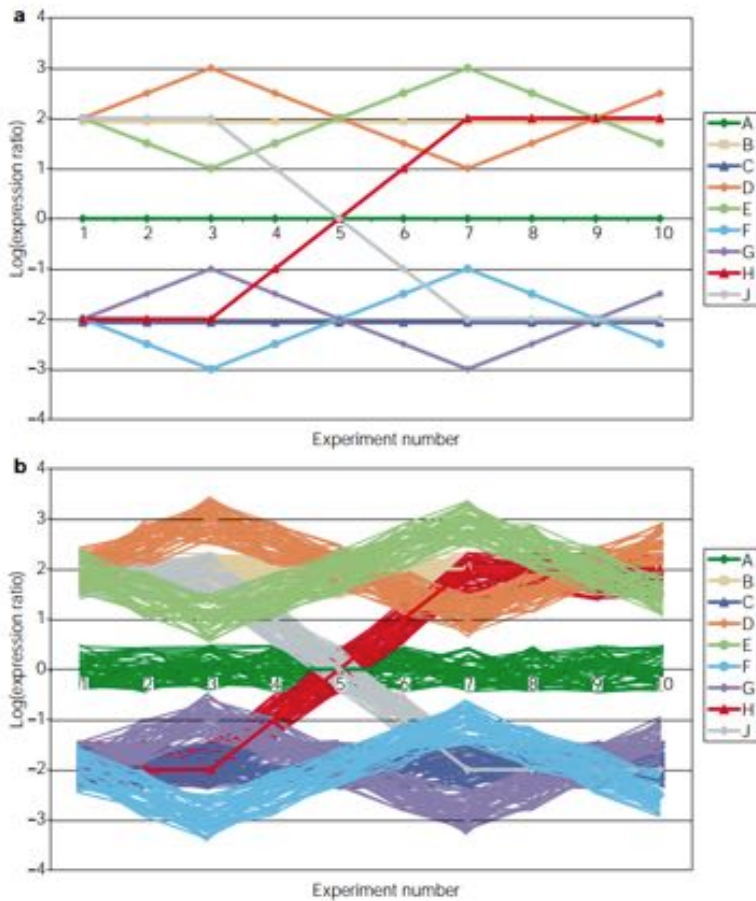
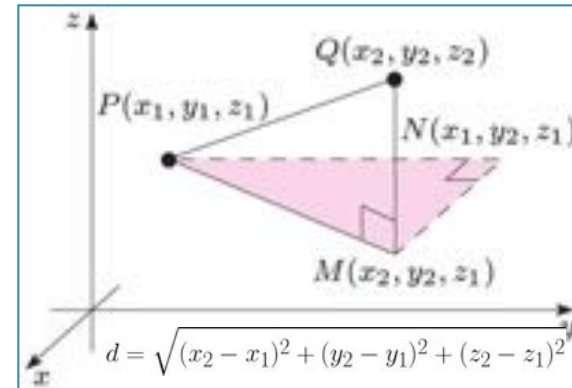
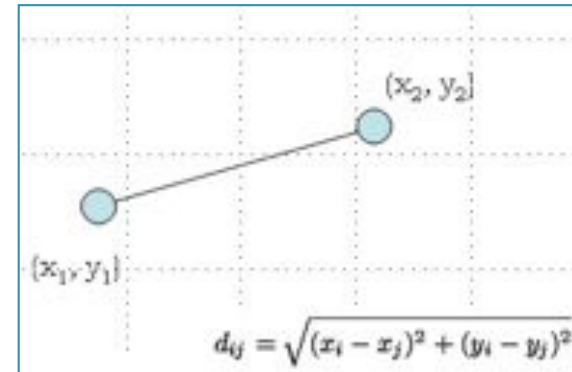


Figure 2 | **A synthetic gene-expression data set.** This data set provides an opportunity to evaluate how various clustering algorithms reveal different features of the data. **a** | Nine distinct gene-expression patterns were created with log₂(ratio) expression measures defined for ten experiments. **b** | For each expression pattern, 50 additional genes were generated, representing variations on the basic patterns.

Euclidean Distance

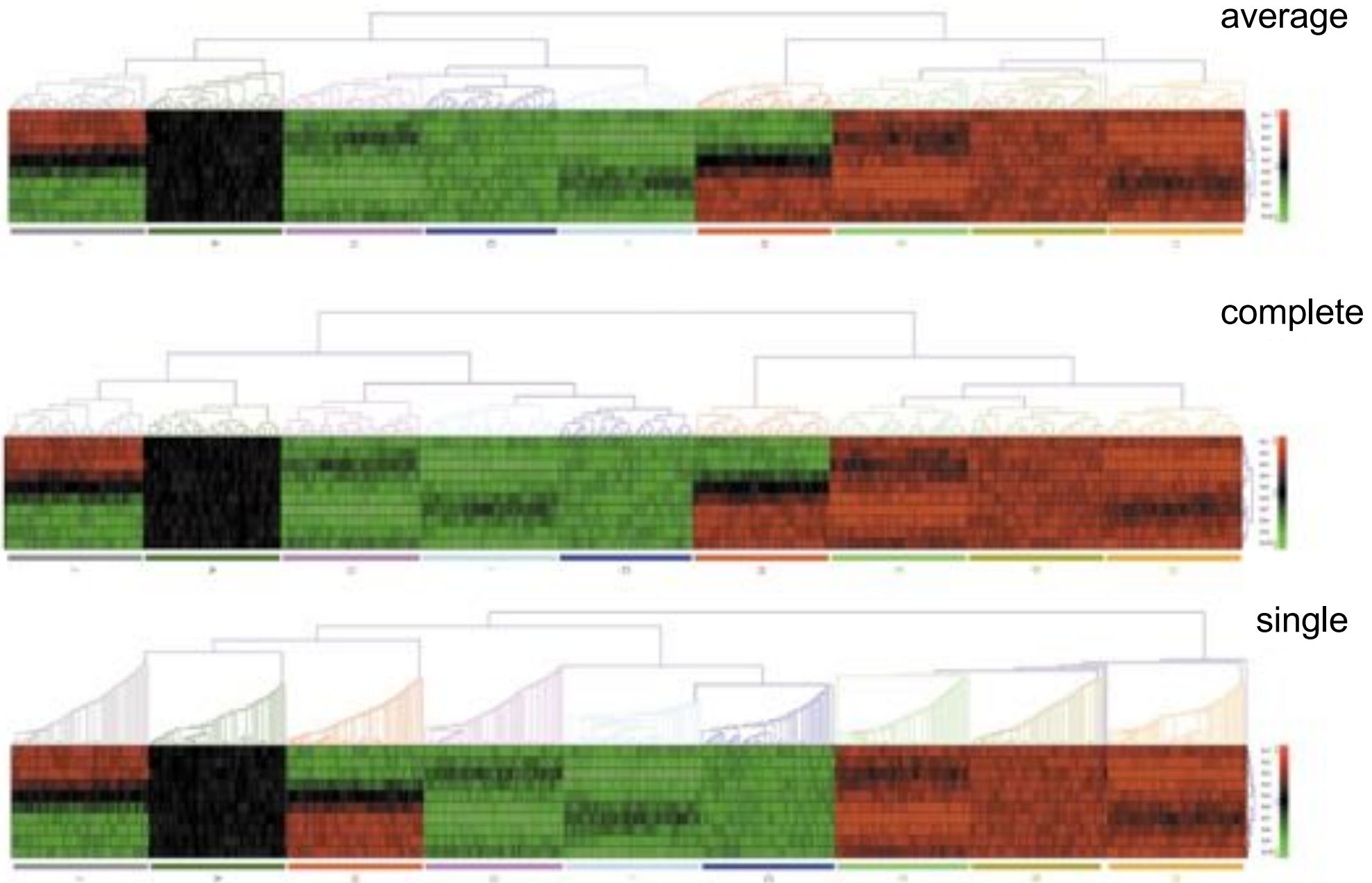


$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}.$$

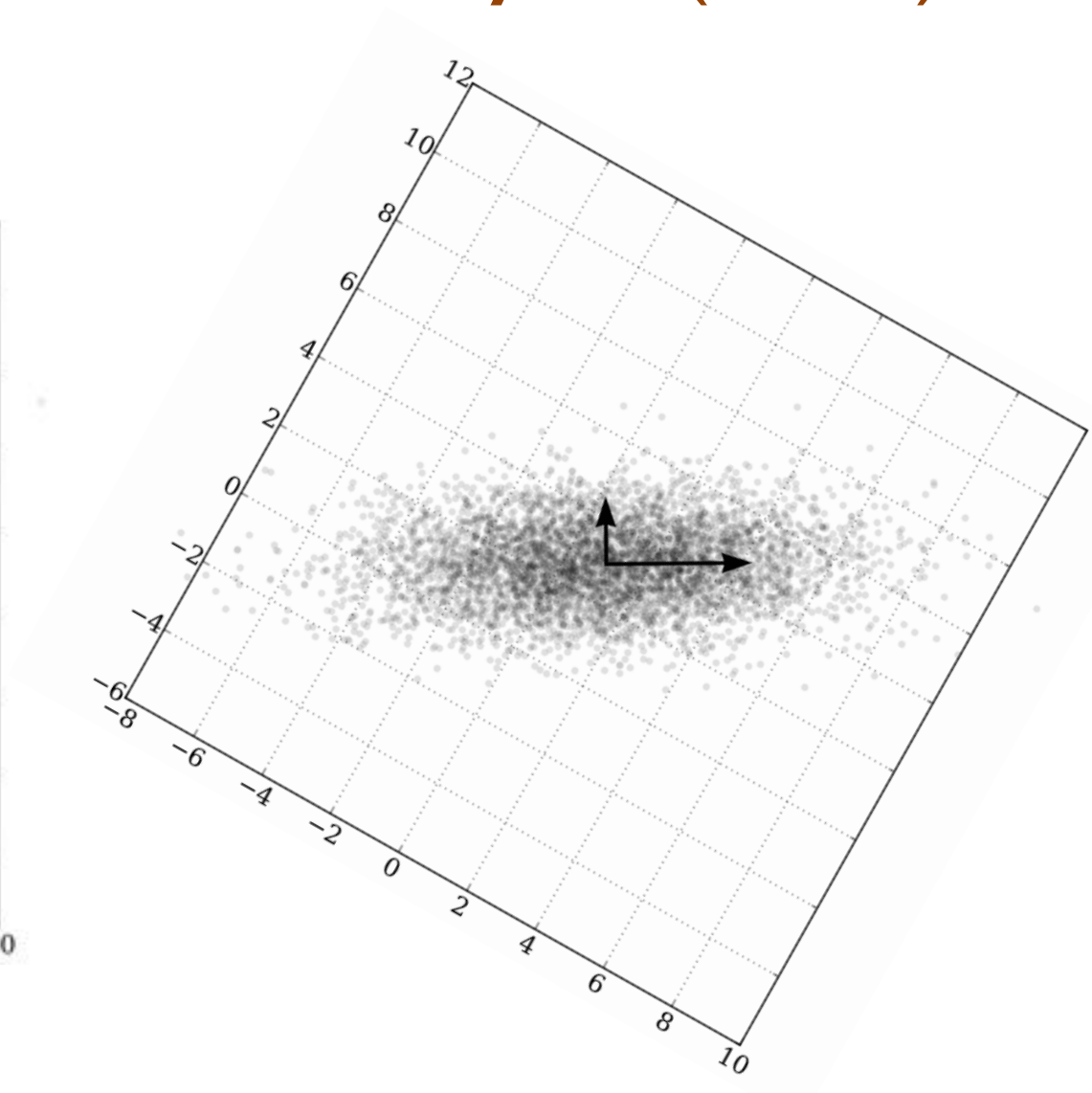
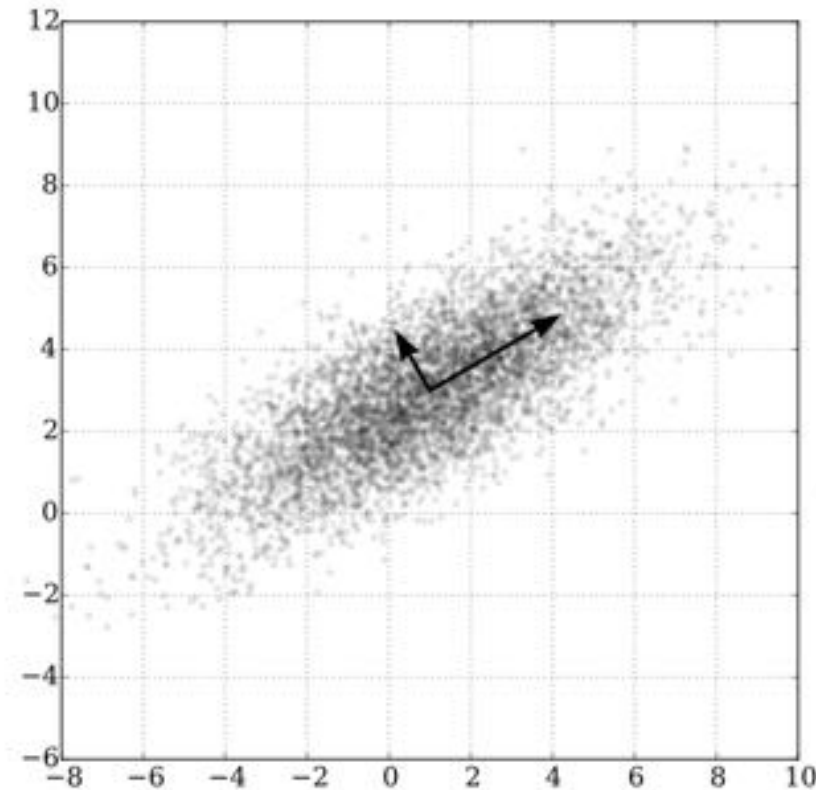
Computational genetics: Computational analysis of microarray data

Quackenbush (2001) *Nature Reviews Genetics*. doi:10.1038/35076576

Hierarchical Clustering



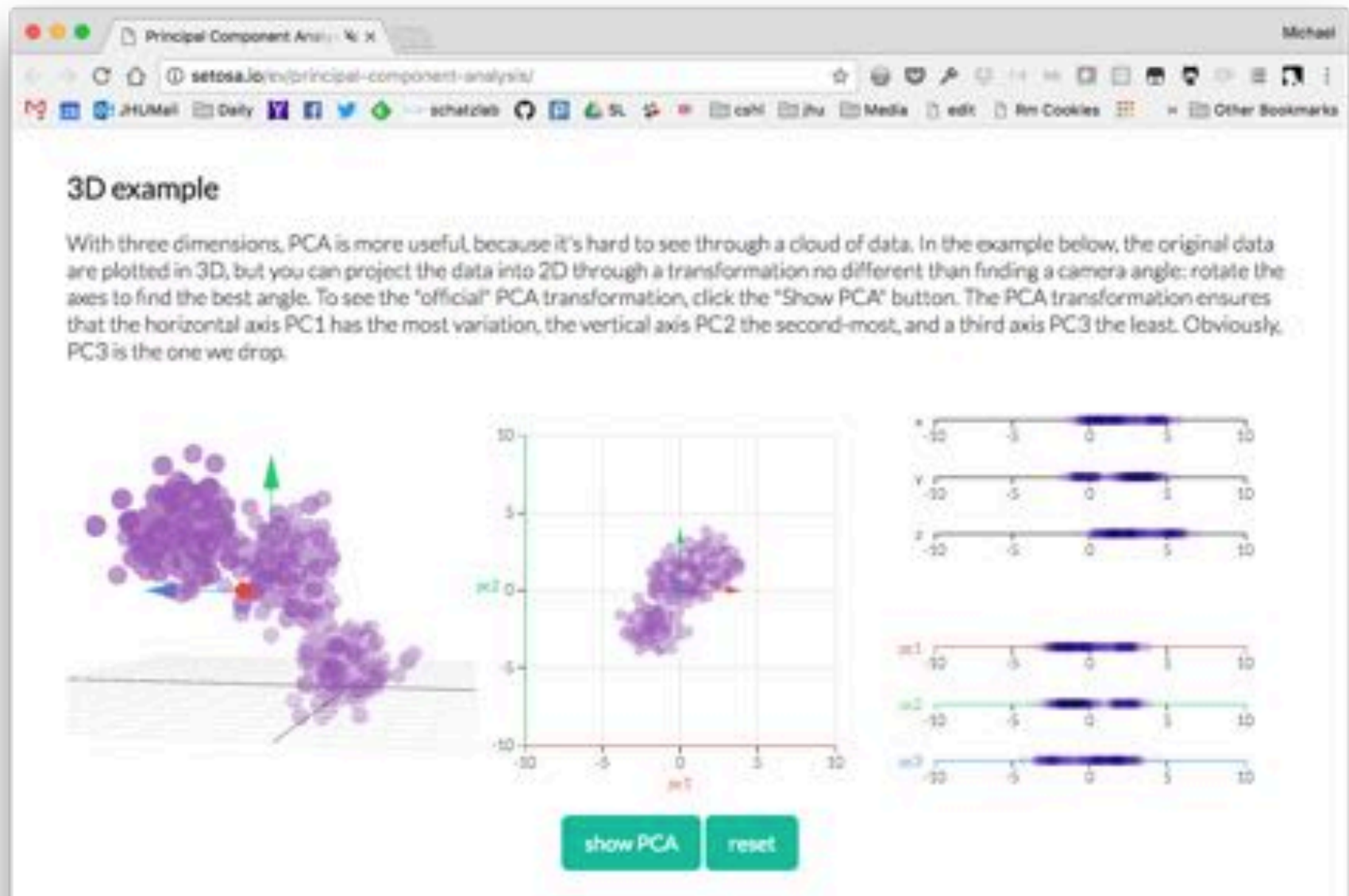
Principle Components Analysis (PCA)



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

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

PCA in Higher Dimensions



<http://setosa.io/ev/principal-component-analysis/>

PC1: "New X"- The dimension with the most variability

PC2: "New Y"- The dimension with the second most variability

Principle Components Analysis (PCA)

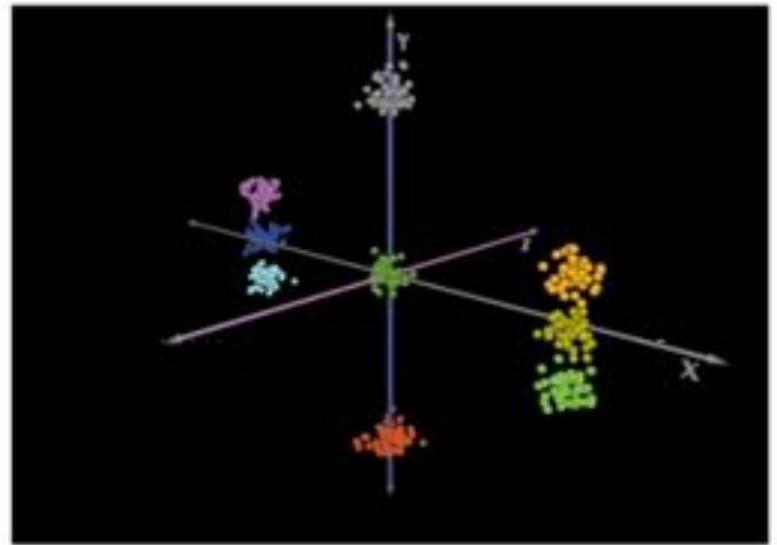
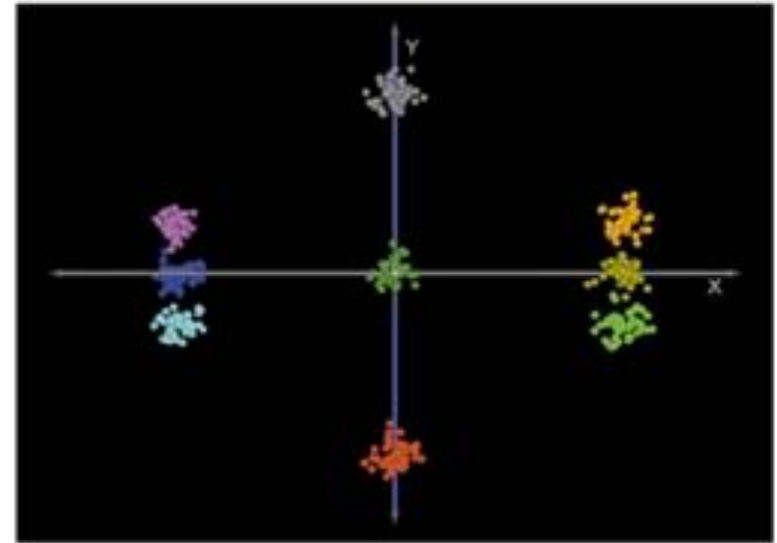
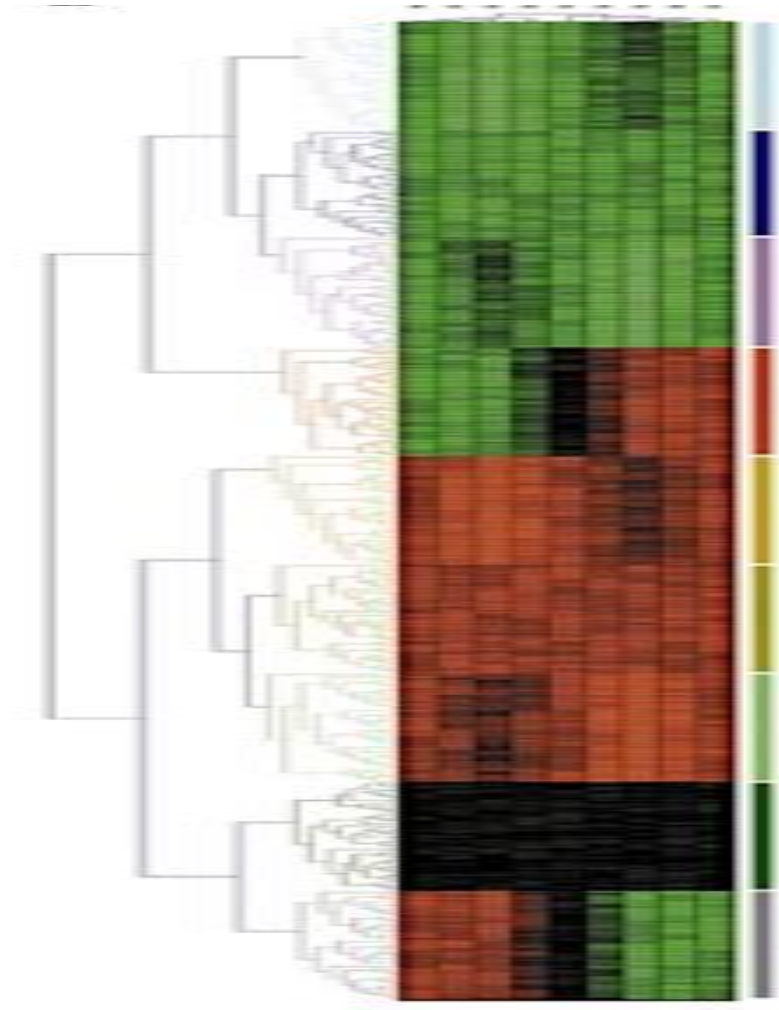
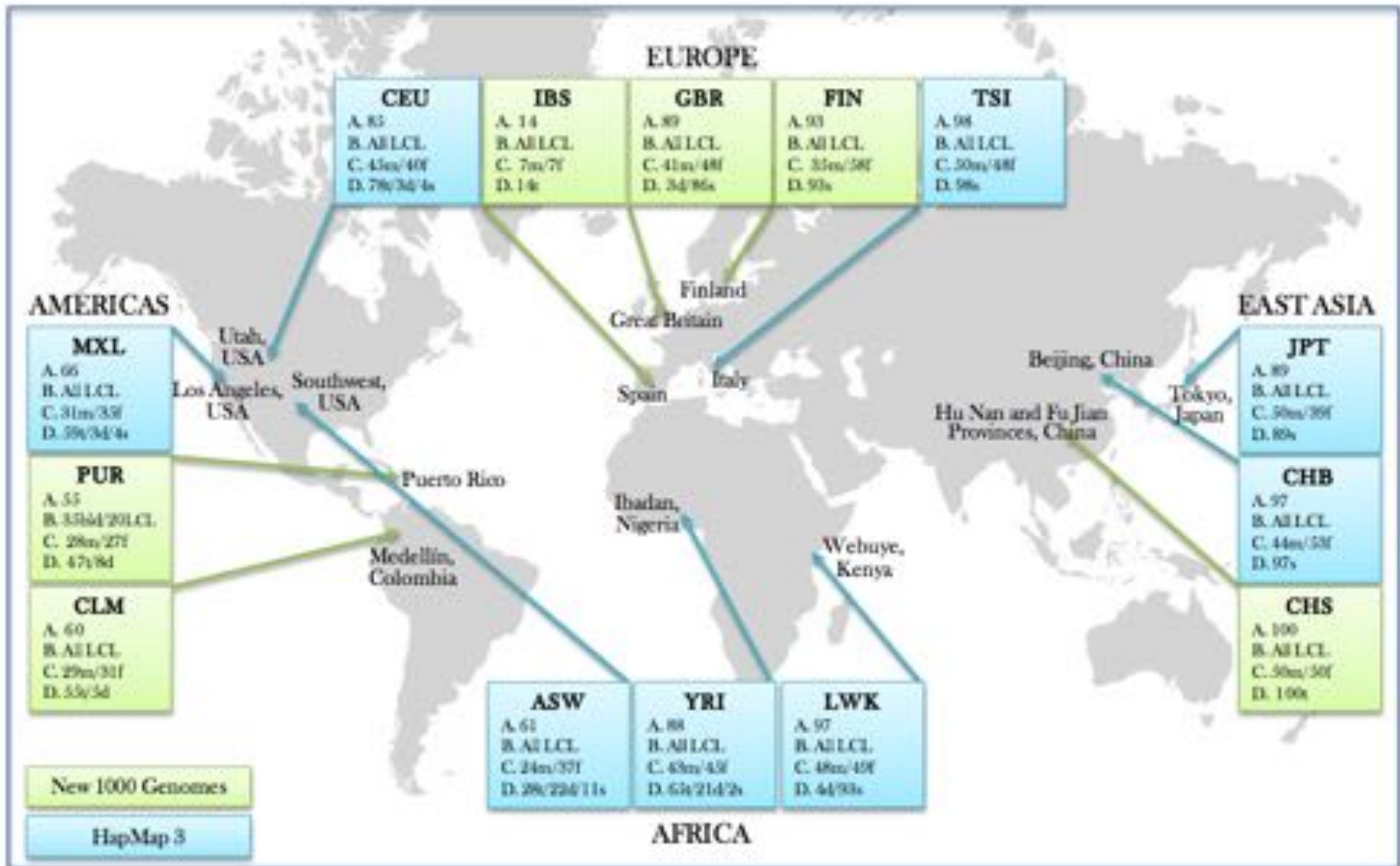


Figure 4 | **Principal component analysis.** The same demonstration data set was analysed using **a** | hierarchical (average-linkage) clustering and **b** | principal component analysis using Euclidean distance, to show how each treats the data, with genes colour coded on the basis of hierarchical clustering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



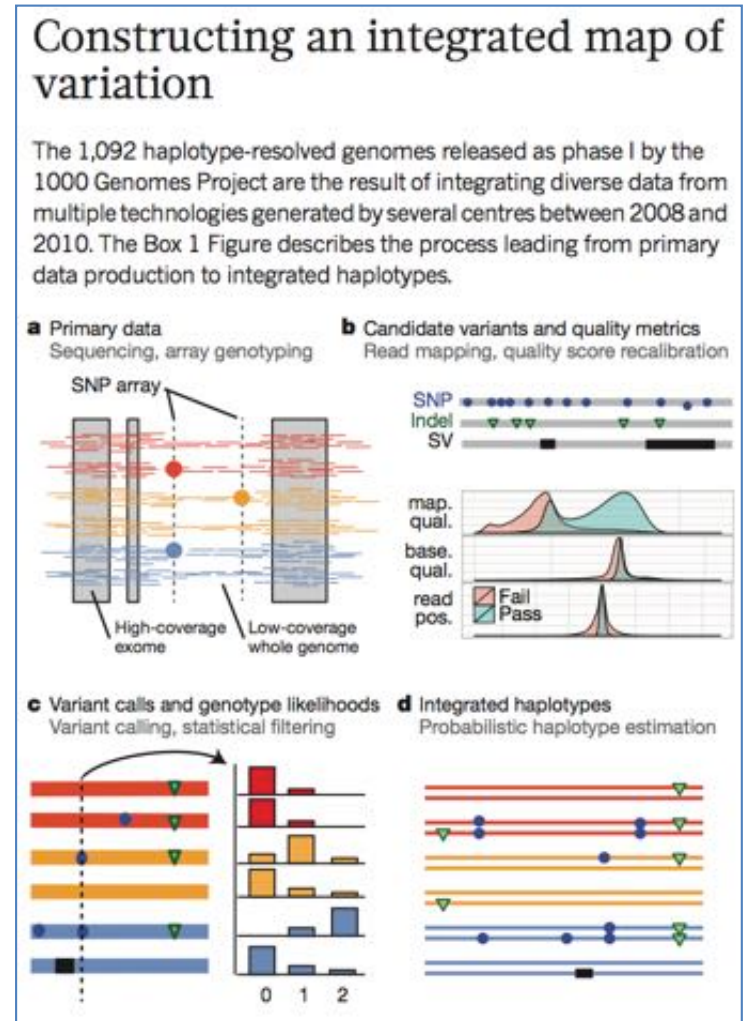
Part 2: (Healthy) Modern Humans

1000 Genomes Populations



1000 Genomes: Human Mutation Rate

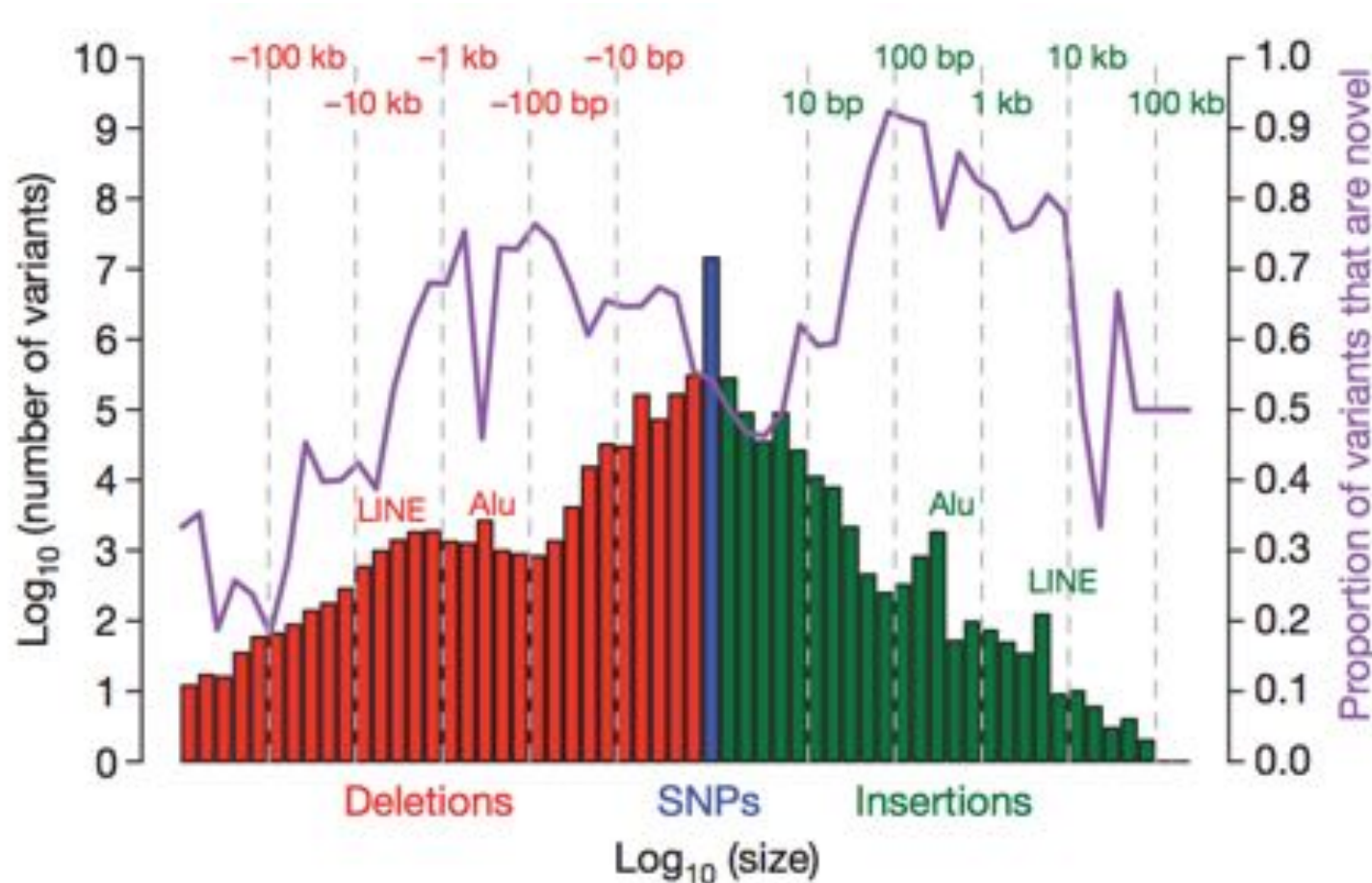
- Phase I Release
 - 1092 individuals from 14 populations
 - Combination of low coverage WGS, deep coverage WES, and SNP genotype data
- Overall SNP rate between any two people is $\sim 1/1200\text{bp}$ to $\sim 1/1300$
 - $\sim 3\text{M}$ SNPs between me and you (.1%)
 - $\sim 30\text{M}$ SNPs between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

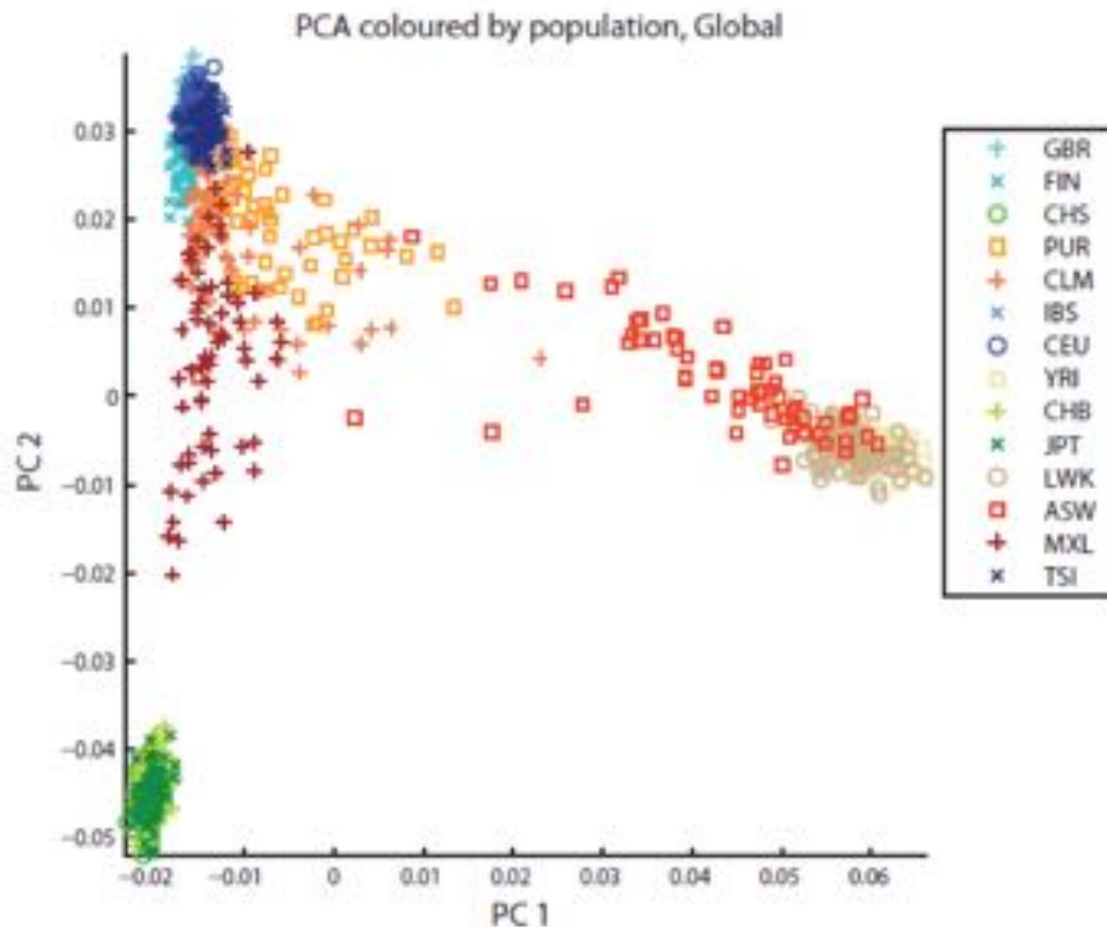


- Mutations follows a “log-normal” frequency distribution
 - Most mutations are SNPs followed by small indels followed by larger events

A map of human genome variation from population-scale sequencing

1000 genomes project (2010) *Nature*. doi:10.1038/nature09534

Variation across populations



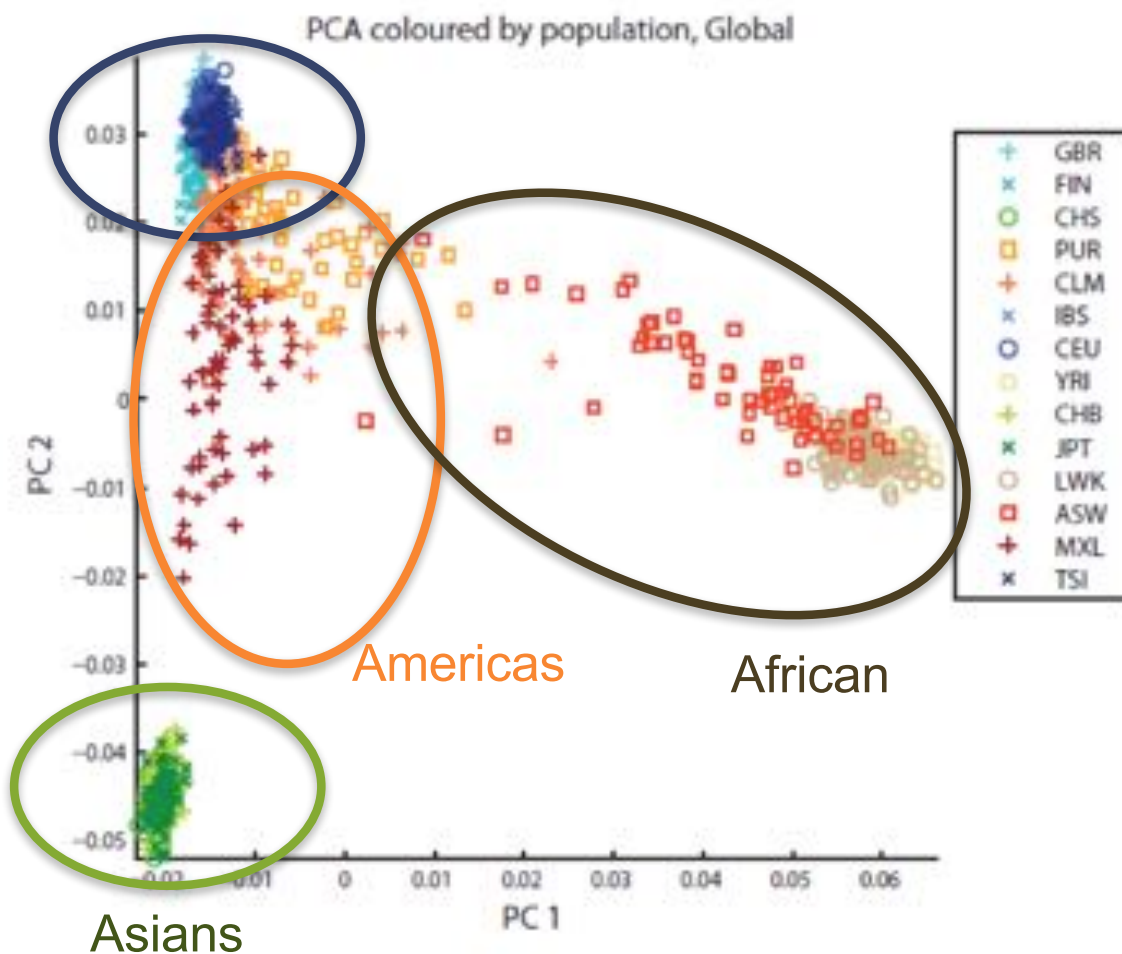
LEVEL	POP_PAIR	# of highly differentiated SNPs	% in transcribed regions*
AFR	ASW-LWK	258	46.8
AFR	LWK-YRI	251	50.2
AFR	ASW-YRI	213	45.8
ASN	CHS-JPT	275	48.1
ASN	CHB-JPT	176	43.7
ASN	CHB-CHS	79	38.7
EUR	FIN-TSI	343	42.6
EUR	CEU-FIN	201	40.7
EUR	FIN-GBR	197	43.2
EUR	GBR-TSI	100	38.9
EUR	CEU-TSI	57	53.8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

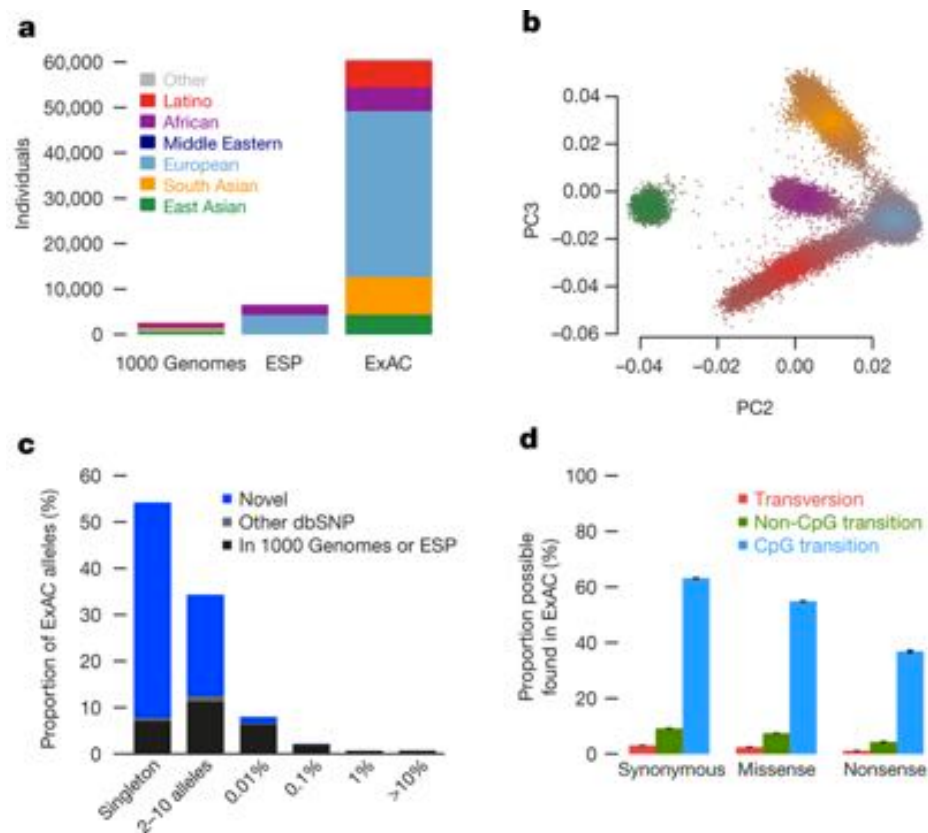


LEVEL	POP_PAIR	# of highly differentiated SNPs	% in transcribed regions*
AFR	ASW-LWK	258	46.8
AFR	LWK-YRI	251	50.2
AFR	ASW-YRI	213	45.8
ASN	CHS-JPT	275	48.1
ASN	CHB-JPT	176	43.7
ASN	CHB-CHS	79	38.7
EUR	FIN-TSI	343	42.6
EUR	CEU-FIN	201	40.7
EUR	FIN-GBR	197	43.2
EUR	GBR-TSI	100	38.9
EUR	CEU-TSI	57	53.8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

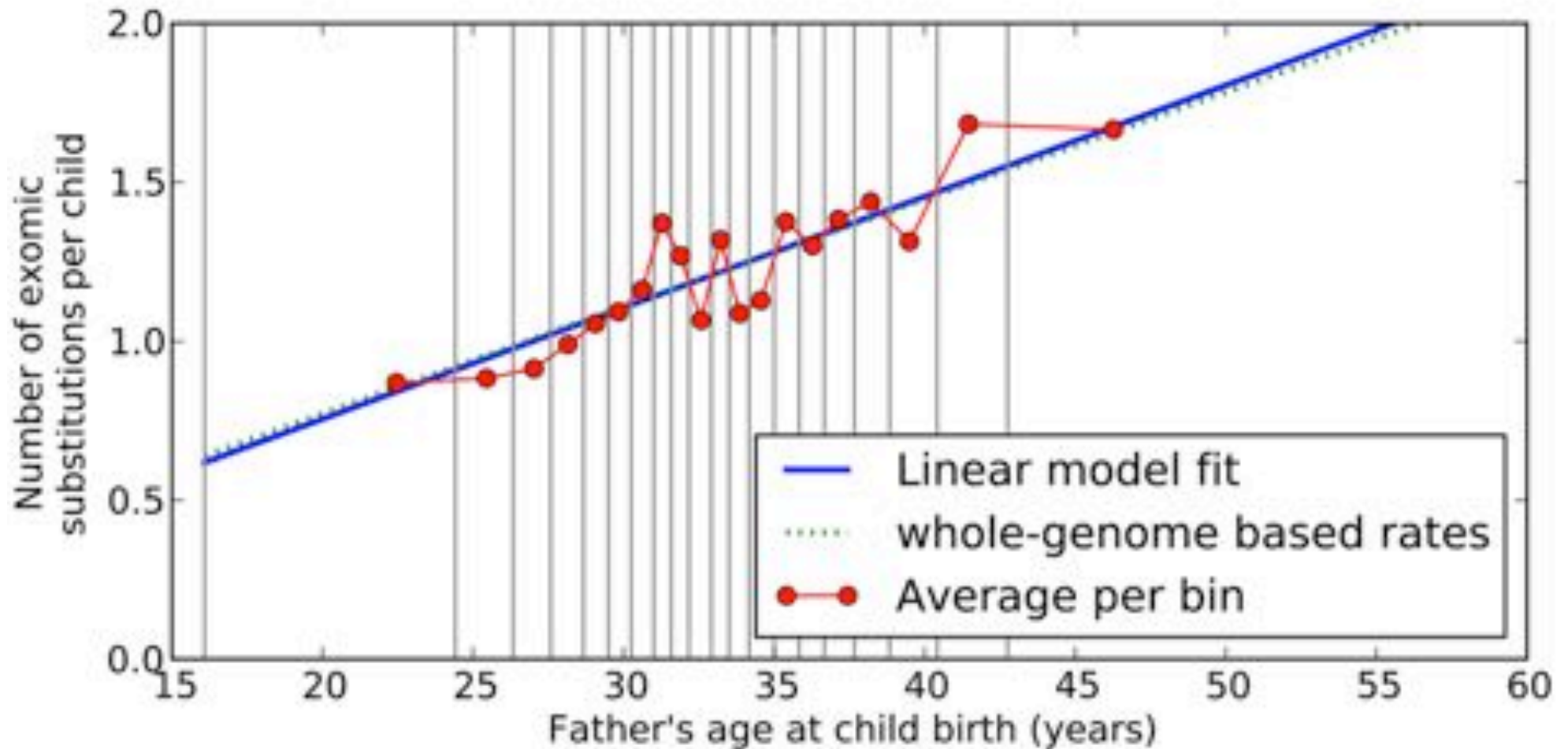
Lek et al (2016) Nature. doi:10.1038/nature19057

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.

De novo Mutations in Men



The contribution of de novo coding mutations to autism spectrum disorder
Iossifov et al (2014) *Nature*. doi:10.1038/nature13908