### Lecture 19. Ancient and Modern Humans Michael Schatz

April 10, 2017 JHU 600.749: Applied Comparative Genomics





### Part I: Ancient Hominds

### **Our Origins**







### Sequencing ancient genomes Janet Kelso Max-Planck Institute

G



#### 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**





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

Green et al. 2010

~35 Illumina flow cells











### Neanderthal Interbreeding



As modern humans migrated out of Africa, they apparently interbred with Neanderthal's 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





![](_page_19_Picture_0.jpeg)

#### **Extraordinary preservation**

![](_page_20_Figure_1.jpeg)

**Sequence length** 

![](_page_20_Figure_3.jpeg)

![](_page_21_Figure_0.jpeg)

![](_page_22_Picture_1.jpeg)

Map after Pickrell et al., 2009

![](_page_23_Picture_1.jpeg)

![](_page_24_Picture_1.jpeg)

![](_page_25_Picture_1.jpeg)

![](_page_26_Picture_1.jpeg)

![](_page_27_Picture_1.jpeg)

### We have always mixed!

![](_page_29_Picture_0.jpeg)

Cite as: B. Vernot et al., Science 10.1126/science.aad9416 (2016).

## Excavating Neandertal and Denisovan DNA from the genomes of Melanesian individuals

#### Benjamin Vernot,<sup>1</sup> Serena Tucci,<sup>1,3</sup> Janet Kelso,<sup>3</sup> Joshua G. Schraiber,<sup>1</sup> Aaron B. Wolf,<sup>1</sup> Rachel M. Gittelman,<sup>1</sup> Michael Dannemann,<sup>3</sup> Steffi Grote,<sup>3</sup> Rajiv C. McCoy,<sup>1</sup> Heather Norton,<sup>4</sup> Laura B. Scheinfeldt,<sup>5</sup> David A. Merriwether,<sup>6</sup> George Koki,<sup>7</sup> Jonathan S. Friedlaender,<sup>8</sup> Jon Wakefield,<sup>9</sup> Svante Pääbo,<sup>2\*</sup> Joshua M. Akey<sup>1\*</sup>

<sup>1</sup>Department of Genome Sciences, University of Washington, Seattle, Washington, USA. <sup>2</sup>Department of Life Sciences and Biotechnology, University of Ferrara, Italy, <sup>1</sup>Department of Evolutionary Genetics, Max-Planck-Institute for Evolutionary Anthropology, Leipzig, Germany, <sup>4</sup>Department of Anthropology, University of Cincinnati, Cincinnati, OH, USA. <sup>1</sup>Coriell Institute for Medical Research, Camden, NJ, USA. <sup>1</sup>Department of Anthropology, Binghamton University, Binghamton, NY, USA. <sup>1</sup>Institute for Medical Research, Goroka, Eastern Highlands Province, Papua New Guinea. <sup>1</sup>Department of Anthropology, Temple University, Philadelphia PA, USA. <sup>1</sup>Department of Statistics, University of Washington, Seattle, Washington, USA.

\*Corresponding author. E-mail: paabol@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.

![](_page_30_Figure_0.jpeg)

### **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 substitutions87 affect splice sites

#### **Changes in Non-coding & regulatory sequences**

26 affect well-defined motifs inside regulatory regions

### **Enrichment analysis**

| Nonsymonymous | None                | Giant melanosomes in melanocytes (p-6.77e-6; FWER=0.091;  |  |  |  |  |  |
|---------------|---------------------|---|--|--|--|--|--|
|               | skin pigmentation   |   |  |  |  |  |  |
| Splice sites  | skin pignentation   |   |  |  |  |  |  |
| 3' UTR        | None                | <ul> <li>1-3 toe syndactyly (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> <li>1-5 toe syndactyly (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> <li>Aplasia/Hypoplasia of the distal phalanx of the thumb (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> <li>Bifid or hypoplastic epiglottis (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> <li>Central polydactyly (feet) (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> </ul>                   |  |  |  |  |  |
| skeletal      | morphologies (lim   | b length, digit development)  |  |  |  |  |  |
|               |                     | FDR=0.0887928)<br>- Dysplastic distal thumb phalanges with a central hole (p=1.34288e-05;   |  |  |  |  |  |
| morphol       | ogies of the laryna | x and the epiglottis FWER-0.538;  |  |  |  |  |  |
|               |                     | <ul> <li>FDR=0.0887928)</li> <li>Laryngeal cleft (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> <li>Midline facial capillary hemangioma (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> <li>Preductal coarctation of the aorta (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> <li>Radial head subluxation (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> <li>Short distal phalanx of the thumb (p=1.34288e-05; FWER=0.538; FDR=0.0887928)</li> </ul> |  |  |  |  |  |

### **FOXP2** Analysis

![](_page_33_Figure_1.jpeg)

Molecular evolution of FOXP2, a gene involved in speech and language Enard et al (2002) Nature. doi:10.1038/nature01025

![](_page_34_Picture_0.jpeg)

### Part II: Modern Humans

ARTICLE

doi:10.1038/nature11632

# 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

![](_page_36_Figure_1.jpeg)

### **1000** Genomes Populations

| Population  | DNA<br>sequenced<br>from blood | Offspring<br>Samples from<br>Trios Available | Pilot<br>Samples | Phase 1<br>Samples | Final Phase<br>Discovery<br>Sample | Final<br>Release<br>Sample | Total |      |
|---|--------------------------------|--|------------------|--------------------|------------------------------------|----------------------------|-------|------|
| Chinese Dai in Xishuanghanna, China(CDX)                            | mi                             | 989  | 0                | 0                  | 99                                 | 93                         |       | - 99 |
| Han Chinese in Bejing, China (CHB)                                  | . 800                          | 80   | 91               | 97                 | 105                                | 000                        |       | 106  |
| Japanese in Tokyo, Japan (JPT)                                      |                                | 80   |                  | 80                 | -104                               | 504                        |       | 305  |
| Kinh in Ho Chi Minh City, Vietnam (KHV)                             | 344                            | 348  | 0                | 0                  | 101                                | 98                         |       | 301  |
| Southern Han Chinese, China (CHS)                                   | 80                             | 988  | .0               | 100                | 308                                | 105                        |       | 112  |
| Total East Asian Ancestry (EAS)                                     |                                |  | 185              | 284                | 515                                | .594                       | 523   |      |
| Bengali in Bangladesh (BEB)   | -                              | 948  | 0                | 0                  | 85                                 | 86                         |       | 85   |
| Gujarati Indian in Houston, TX (GIH)                                | -                              | 908  | 0                | 0                  | 106                                | 103                        |       | 106- |
| Indian Teluga in the UK (ITU)                                       | yes                            | pers   | 0                | 0                  | 105                                | 102                        |       | 103  |
| Punjabi in Lahore Pakistan (PJL)                                    | 3444                           | 344  | .0               | 0                  | 96                                 | 96.                        |       | 96   |
| Sri Lankan Tamil in the UK (STU)                                    | yes                            | 248  | 0                | 0                  | 103                                | 003                        |       | 303  |
| Total South Asian Ancestry (SAS)                                    |                                |  |                  | 1.1                | 494                                | 400                        | - 694 |      |
| African Ancestry in Southwest US (ASW)                              |                                | . 200  | 0                | 51                 | 66                                 | 61                         |       | - 66 |
| African Caribbean in Barbados (ACB)                                 | 705                            | 3428   | 0                | 0                  | - 96                               | 95                         |       | 96   |
| Esan in Nigeria (ESN)   | 888                            | yes  | 0                | 0                  | 99                                 | 95                         |       | 99   |
| Gambian in Western Division, The Gambia (GWD)                       | 100                            | yes  | .0               | 0                  | 123                                | 113                        |       | 113  |
| Luhya in Webuye, Kenya (LWK)  | ine .                          | pes  | 182              | 87                 | 301                                | 90                         |       | 110  |
| Mende in Siema Leone (MSL)  |                                | 348  | 0                |                    | 85                                 | 85                         |       | 85   |
| Yoruba in Ibadan, Nigeria (YRI)                                     | 89                             | 508  | 106              | 85                 | 109                                | 108                        |       | 116  |
| Total African Ancestry (AFR)  |                                |  | 208              | 346                | 649                                | 662                        | 671   |      |
| British in England and Scotland (GBR)                               |                                | 988  | 0                | 89                 | 82                                 | 90                         |       | 94   |
| Finnish in Finland (FIN)  | -                              | 80   | 0                | 93                 | 99                                 | 99                         |       | 300  |
| Iberian populations in Spain (IBS)                                  |                                | 348  | 0                | LH.                | 897                                | 103                        |       | 307  |
| Toscani in Italy (TSI)  | 100                            | 80   | 66               | 56                 | 108                                | 1007                       |       | 110  |
| Utah residents with Northern and Western European<br>ancestry (CEU) | -                              | 344  | 94               | 85                 | 99                                 | 98                         |       | 303  |
| Total European Ancestry (EUR)                                       |                                |  | 190              | 379                | 508                                | .540                       | 514   |      |
| Colombian in Medellin, Colombia (CLM)                               | inter (                        | 2000   | 0.0              | 60                 | 94                                 |                            |       | 95   |
| Mexican Ancestry in Los Angeles, California (MXL)                   | 80                             | Side   | 0                | 65                 | 67                                 | 64                         |       | 69   |
| Peruvian in Lima, Peru (PEL)  | 765                            | 7428   | 0                | 0                  | 56                                 | 85                         |       | 86   |
| Puerto Rican in Puerto Rico (PUR)                                   | yes                            | jan  | 0                | 55                 | 105                                | 104                        |       | 305  |
| Total Americas Ancestry (AMR)                                       |                                |  |                  | 181                | 312                                | M                          | 355   |      |
| Total   |                                |  | A10              | 3892               | 2536                               | 2994                       | 2577  |      |

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

#### Constructing an integrated map of variation

The 1,092 haplotype-resolved genomes released as phase I by the 1000 Genomes Project are the result of integrating diverse data from multiple technologies generated by several centres between 2008 and 2010. The Box 1 Figure describes the process leading from primary data production to integrated haplotypes.

![](_page_38_Figure_12.jpeg)

#### An integrated map of genetic variation from 1,092 human genomes 1000 genomes project (2012) *Nature*. doi:10.1038/nature11632

### Human Mutation Types

![](_page_39_Figure_1.jpeg)

- 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, <sup>1,2\*</sup> Suganthi Balasubramanian, <sup>3,4</sup> Adam Frankish, <sup>1</sup> Ni Huang, <sup>1</sup> James Morris, <sup>1</sup> Klaudia Walter, <sup>1</sup> Luke Jostins, <sup>1</sup> Lukas Habegger, <sup>3,4</sup> Joseph K. Pickrell, <sup>5</sup> Stephen B. Montgomery, <sup>6,7</sup> Cornelis A. Albers, <sup>1,8</sup> Zhengdong D. Zhang, <sup>9</sup> Donald F. Conrad, <sup>10</sup> Gerton Lunter, <sup>11</sup> Hancheng Zheng, <sup>12</sup> Qasim Ayub, <sup>1</sup> Mark A. DePristo, <sup>13</sup> Eric Banks, <sup>13</sup> Min Hu, <sup>1</sup> Robert E. Handsaker, <sup>13,14</sup> Jeffrey A. Rosenfeld, <sup>15</sup> Menachem Fromer, <sup>13</sup> Mike Jin, <sup>3</sup> Xinmeng Jasmine Mu, <sup>3,4</sup> Ekta Khurana, <sup>3,4</sup> Kai Ye, <sup>16</sup> Mike Kay, <sup>1</sup> Gary Ian Saunders, <sup>1</sup> Marie-Marthe Suner, <sup>1</sup> Toby Hunt, <sup>1</sup> If H. A. Barnes, <sup>1</sup> Clara Amid, <sup>1,17</sup> Denise R. Carvalho-Silva, <sup>1</sup> Alexandra H. Bignell, <sup>1</sup> Catherine Snow, <sup>1</sup> Bryndis Yngvadottir, <sup>1</sup> Suzannah Bumpstead, <sup>1</sup> David N. Cooper, <sup>18</sup> Yali Xue, <sup>1</sup> Irene Gallego Romero, <sup>1,5</sup> 1000 Genomes Project Consortium, Jun Wang, <sup>12</sup> Yingrui Li, <sup>12</sup> Richard A. Gibbs, <sup>19</sup> Steven A. McCarroll, <sup>13,14</sup> Emmanouil T. Dermitzakis, <sup>7</sup> Jonathan K. Pritchard, <sup>5,20</sup> Jeffrey C. Barrett, <sup>1</sup> Jennifer Harrow, <sup>1</sup> Matthew E. Hurles, <sup>1</sup> Mark B. Gerstein, <sup>3,4,21</sup>† Chris Tyler-Smith<sup>1</sup>†

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.

(2012) Science. doi: 10.1126/science.1215040

### Homozygous LoF Mutations

#### LETTER

doi:10.1038/nature22034

#### Human knockouts and phenotypic analysis in a cohort with a high rate of consanguinity

Danish Saleheen<sup>1,3</sup>\*, Pradeep Natarajan<sup>3,4</sup>\*, Irina M. Armean<sup>4,3</sup>, Wei Zhao<sup>3</sup>, Asif Rasheed<sup>2</sup>, Sumeet A. Khetarpal<sup>6</sup>, Hong-Hee Won<sup>7</sup>, Konrad J. Karczewski<sup>4,3</sup>, Anne H. O'Donnell-Luria<sup>4,3,8</sup>, Kaitlin E. Samocha<sup>4,3</sup>, Benjamin Weisburd<sup>4,5</sup>, Namrata Gupta<sup>4</sup>, Mozzam Zaidi<sup>2</sup>, Maria Samuel<sup>2</sup>, Atif Imran<sup>2</sup>, Shahid Abbas<sup>9</sup>, Faisal Majeed<sup>2</sup>, Madiha Ishaq<sup>2</sup>, Saba Akhtar<sup>2</sup>, Kevin Trindade<sup>6</sup>, Megan Mucksavage<sup>6</sup>, Nadeem Qamar<sup>10</sup>, Khan Shah Zaman<sup>10</sup>, Zia Yaqoob<sup>10</sup>, Tahir Saghir<sup>10</sup>, Syed Nadeem Hasan Rizvi<sup>10</sup>, Anis Memon<sup>10</sup>, Nadeem Hayyat Mallick<sup>11</sup>, Mohammad Ishaq<sup>12</sup>, Syed Zahed Rasheed<sup>12</sup>, Fazal-ur-Rehman Memon<sup>13</sup>, Khalid Mahmood<sup>14</sup>, Naveeduddin Ahmed<sup>13</sup>, Ron Do<sup>16,17</sup>, Ronald M. Krauss<sup>18</sup>, Daniel G. MacArthur<sup>4,3</sup>, Stacey Gabriel<sup>4</sup>, Eric S. Lander<sup>4</sup>, Mark J. Daly<sup>4,5</sup>, Philippe Frossard<sup>2</sup>8, John Danesh<sup>19,20</sup>8, Daniel J. Rader<sup>6,21</sup>8 & Sekar Kathiresan<sup>3,4</sup>8

A major goal of biomedicine is to understand the function of every across 14,345 autosomal genes were annotated as pLoF mutations (that gene in the human genome<sup>1</sup>. 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 proteincoding 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 SLC9A3RI, with mediators of calcium and phosphate signalling. Heterozygous deficiency of APOC3 has been shown to protect against coronary heart disease<sup>4,5</sup>; 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

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 × 10<sup>-16</sup>) (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

#### (2017) Nature. doi:10.1038/nature22034

### Variation across populations

![](_page_42_Figure_1.jpeg)

| LEVEL | POP_PAIR | # of Highly<br>differentiated<br>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

![](_page_43_Figure_2.jpeg)

- 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

![](_page_44_Figure_1.jpeg)

- The aggregation and analysis of highquality 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 nearcomplete depletion of predicted protein-truncating

#### Analysis of protein-coding genetic variation in 60,706 humans

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

### dbSNP

| S NCBI   | dbSNP<br>Short Genetic Variations                                   |                  |                             |  |   | 野  |                            |                              |   |   |
|--|---|------------------|-----------------------------|--|---|--|----------------------------|------------------------------|---|---|
| PubMed Nucleotid   | e Protein Genome  | Structu          | re PopSe                    | t Taxono                                       | тy  | OMIM Books SN                                      | P                          |                              |   |   |
| Parent Entres ( 140  | Search for SNP  | on NCBI          | Referenc                    | e Assembi                                      | Y   |  | -                          |                              |   |   |
| Search Entrez Sor  | = ) 101   | S                |                             | Go   |   |  |                            |                              |   |   |
|  | dbSNP Summary   |                  |                             |  |   |  |                            |                              |   |   |
| Have a question<br>about db5NP7 Try  | RELEASE: NC   | BI dbS           | NP Buil                     | ld 141   |   |  |                            |                              |   |   |
| searching the SMP  |   |                  |                             |  |   |  |                            |                              |   |   |
| the section of the se | doonin Component  | Availabi         | iny Dates:                  |  |   |  |                            |                              |   |   |
| Ge   | Compon  | ent              | Date                        | e available                                    |   |  |                            |                              |   |   |
| Children of the second   | dbSNP web query   | for build        | 141: May                    | 21, 2014                                       |   |  |                            |                              |   |   |
| ENERAL   | ftp data for build 1  | 41:              | May                         | 21, 2014                                       |   |  |                            |                              |   |   |
| antact Us  | Entrez Indexing for   | r build 14       | H: May                      | 21, 2014                                       |   |  |                            |                              |   |   |
| te Map   | BLAST database for  | or build 1       | 41: Mm                      | 21, 2014                                       |   |  |                            |                              |   |   |
| bSNP Homepage  |   |                  |                             |  |   |  |                            |                              |   |   |
| bSNP Summary   | <ul> <li>The complete data</li> <li>All formats and corr</li> </ul> | for build        | 141 are av<br>are describ   | allable at the                                 | to no   | nobi nim nih govisno<br>bi nim nih govisno/00      | o/ in m                    | ultiple form                 | nats.   |   |
| TP Download  |   |                  | and George                  |  |   |  |                            | THE OCH                      |   |   |
| UMAN VARIATION   | <ul> <li>Please address an</li> </ul>                               | y question       | ns or comm                  | nents regard                                   | sing t  | he data to snp-admin                               | C PRICE                    | si nim nih g                 | 204   |   |
| OCUMENTATION   | New Submission si   | nce prev         | ious build                  |  |   |  |                            |                              |   |   |
| EARCH  | Organism  | Current<br>Build | t New Submissions<br>(ss#s) |  | New RefSNP Clusters<br>(rs#'s) ( # validated) |  | New ss# with N<br>Genotype |                              | New ss# with<br>Frequency                                 |   |
|  | Homo sepiens  | 141              | 20,708,470                  |  | 137 (0)                                       |  |                            |                              | 4   |   |
|  | Total: 1 Organisms  |                  | 20,708,470                  |  | 137 ()  |  |                            |                              | 4   |   |
|  | *Submissions receive<br>BUILD STATISTICS<br>Organism                | dbSNP<br>Build   | Genome<br>Build             | of current b<br>Number<br>Submissie<br>(ss#'s) | of<br>ons                                     | Number of<br>RefSNP Cluster<br>(rs#'s) ( # validab | ≢ clus<br>rs<br>ed)        | Number<br>(rs#'s)<br>in gene | next build<br>of Number of<br>(ss3"s)<br>with<br>genotype | Number of<br>(ss#'s)<br>with<br>frequency |
|  | Homo sapiens  | 141              | 38.1                        | 260,570  | 204   | 62,387,983 (43,737                                 | ,321)                      | 29.901.11                    | 73,909,256  | 35,997,943                                |
|  | Total 4 Conscious   |                  | 0                           | 000 000  |   |  | -                          |                              |   |   |

genome

- 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.

![](_page_46_Picture_0.jpeg)

### Part III:

### Genetic Privacy

![](_page_47_Picture_0.jpeg)

Identifying Personal Genomes by Surname Inference Melissa Gymrek *et al. Science* **339**, 321 (2013); DOI: 10.1126/science.1229566

![](_page_47_Picture_3.jpeg)

![](_page_47_Picture_4.jpeg)

### What are microsatellites

#### • Tandemly repeated sequence motifs

- Motifs are I 6 nt long
- So far, min. 8 nt length, min. 3 tandem repeats for our analyses

#### Ubiquitous in human genome

>5.7 million uninterrupted microsatellites in hg19

#### • Extremely unstable

- Mutation rate thought to be  $\sim 10^{-3}$  per generation in humans

#### Unique mutation mechanism

- Replication slippage during mitosis and meiosis

#### • May be under neutral selection

 $\mathsf{cCTCTCTCTCTCTCTCTCTCTCTCA} \rightarrow (\mathsf{CT})_{13} \qquad \mathsf{tCAACAACAACAACAACAAAA} \rightarrow (\mathsf{CAA})_7$ 

 $tTTGTCTTGTCTTGTCTTGTCTTGTCC \rightarrow (TTGTC)_{6} \quad cCATTCATTCATTCATTa \rightarrow (CATT)_{4}$ 

**Microsatellites: Simple Sequences with Complex Evolution** Ellegren (2004) *Nature Reviews Genetics*. doi:10.1038/nrg1348

### Replication slippage

#### Out-of-phase re-annealing

- Nascent and template strands dissociate and re-anneal out-of-phase
- Loops repaired by mismatch repair machinery (MMR)
  - Very efficient for small loops
  - Possible strand-specific repair

#### Stepwise process

- Nascent strand gains or loses full repeat units
- Typically single unit mutations
- Varies by motif length, motif composition, etc.

![](_page_49_Figure_10.jpeg)

Ellegren (2004) Nature Reviews Genetics. doi:10.1038/nrg1348

![](_page_49_Figure_12.jpeg)

#### **lobSTR Algorithm Overview**

![](_page_50_Figure_1.jpeg)

**IobSTR:A short tandem repeat profiler for personal genomes** Gymrek et al. (2012) *Genome Research*. doi:10.1101/gr.135780.111

# Why should we care about microsatellites?

- Polymorphism and mutation rate variation
- Disease
  - Huntington's Disease
  - Fragile X syndrome
  - Friedrich's ataxia
- Mutations as lineage
  - Organogenesis/embryonic development
  - Tumor development

![](_page_51_Figure_9.jpeg)

#### **Phylogenetic fate mapping**

Salipante (2006) PNAS. doi: 10.1073/pnas.0601265103

![](_page_52_Picture_0.jpeg)

The Contineed DNA trabes System, or CODS, teents forenell adamos and computer technology risk a loof for Intellig extent ormes. It analities federal, state, and local forenell adamos and computer technology risk a loof for Intellig extent ormes. It analities federal, state, and local forenell adamos and computer technology risk a loof for Intellig extent ormes. It analities federal, state, and local forenell indications to exchange and computer technology risk a loof for Intellig extent ormes. It analities federal, state, and local forenell indications to exchange and computer technology risk and one officients. Using the features DNA findes System of CODIS, the features DNA fores also helps identify missing and undertified individuals.

#### Overview

CODID generates investigative leads in cases where lookgival enteriors is recovered from the crime scene. Makhes made among profiles in the Porenaic Index can be come scenes logither, prevailty identifying senal affenders. Saled upon a match, police from multiple pradiciture can occurring their respective investigations and share the loads they developed independently. Motifies made between the Porenaic and Offender Indexes provide investigations with the identity of auspected perpendents. Since names and other personally identifiable internation are not stored at XDIS, qualified DNA analysis in the laborationes sharing matching profiles contact salts other to confirm the candidate match.

#### History

The FBI Laboratory's CODIB began as a pilot tofheare project in 1998, anning 14 alute and tocal laboratories. The DHA lobetthcator. Act of 1984 formatized the FBI's authority to saturate a National DHA Index Bystein (NDIB) for law antiprocess. Today, over 190 public law enforcement laboratories perforgate in NDIB across the United Blates. Internationally, more than 90 law antiproximent laboratories in over 50 countries use the CODIB achieve for their sam detailese.

#### Mission

The CODID Unit manages CODID and NDID. It is responsible for identegong, providing, and supporting the CODID program to federal, stude, and tool orms factorians in the United Dates and selected international law enforcement online teboretories to federa the exchange and companium of Romatic DNA endance from vicent, other Investigations. The CODID Unit.

#### **Genealogy Databases**

![](_page_53_Picture_1.jpeg)

#### GENETICS

**CELL REPOSITORIES** 

### Genealogy Databases Enable Naming Of Anonymous DNA Donors

#### Surname Inference

![](_page_54_Picture_1.jpeg)

#### Whose sequence reads are these?

![](_page_54_Picture_3.jpeg)

#### Identifying Personal Genomes by Surname Inference

Gymrek et al (2013) Science. doi: 10.1126/science.1229566

# Step 1. Profile Y-STRs from the individual's genome.

#### DYS458: 17 repeats

The human reference genome contains 16 copies of "TTTC". Venter has an extra copy of "TTTC", giving him a genotype of "17" at this marker. In a similar way, we can profile all other genealogical STR markers on the Y-chromosome where we know Venter's genome sequence to get the value of a whole panel of these markers.

# Step 2. Search for a surname hit in online genetic genealogy databases.

| DYS 303     | 0115 390         | DYS 19384             | DVS 196*           | 10 10               |                  | 015 3850   | 0Y5-428    | 0VE 348  | 0Y5 430        |
|-------------|------------------|-----------------------|--------------------|---------------------|------------------|------------|------------|----------|----------------|
| DYS 365-1** | 0V0 392          | 016 366 2             | 018 458            | 0115 450a<br>9 450a | 078.450          | DVS 455*** | DYS 454*** | DVD 447  | 075 437<br>- 0 |
| DYS 448     | DV5 443          | DV9 464a              | 070 4040<br>- 4040 | DYS 484             | DY3 4043         | DYS 4044*  | 01V2 4040* | DYS Hong | CY13 460       |
|             | VCA.14***        | VCA IIb <sup>ma</sup> | DVS-458            | DYS 607             | 045 576          | 0145 570   | COY .      | 00Y 8    | 0VS 442        |
| DVS 438     | 275 531<br>32 \$ | 015 575               | DYS 396014         | CVS 385519          | 01/13 500        | 01/5 537   | DY3.641    | DVS 472  | DYS-40651      |
| OV5 511     | CVS 425          | DYS-4134              | 0Y5-4120           | 0V5 117             | DVS 504          | DVS 436    | 0115-000   | DVS 534  | DYS 450        |
| DY19 444    | 1745 481<br>22 8 | DY'8 520              | DYB 440            | DYS 817             | 0115 588<br>41 4 | 0 487      | 010 572    | CY15 642 | DYE 452        |
| 0Y8 M5      | DYS 401-         | 015 452               | 0.478.410<br>0 4   | 015 ES              | -                | DVS 441    | DYS-445    | 010 A12  | 015 463        |
| 010 45H     | 0111 435         | 095 495               | DYS 454            | DYS 485             | 0+1 605          | 0+5 527    | 0vil \$10  | 0VS 540  | GY6 586        |
| oveisre     | CYB MH           | DVS 836               | DV8-638            | DV5.643             | 25 4             | DVS 718    | tri 210    | DV8 728  |                |

http://www.ysearch.org

# Step 3. Search with additional metadata to narrow down the individual.

We enter the search information: Venter, CA, and 66: Tell Us Who You're Looking Forf Lives at Rasi Permat Inc. Rep. Chat North Resident print teaching our committee for another of 15 and post-industry to the ingents ALC: NO. OF THE OWNER. the surgery of Photo: Address Not loved in تكالبوه فسيلتها Borked at Property Regard 1. J Croig Vertier Las Anjaries, C.A. 100 and Woman La Monda, GA Line Torpey a second of the Catalan, CA. Your More Caritovila, MO Caritoriulu, 104 And Landson Passed Franker W inginité Verezi all Without Cusarninga CA Verhalt Resident Warried Cutterinotype Dartiera, CA Propage P Visiting Jaff Vietter Charlottain Long Beath, CA Ophia Variat Pellowship Tempical CA Lot Vensel 10ew More Lakewood, CJ. Rom Parate More Lacation

http://www.ussearch.com

#### Surname Inference

![](_page_58_Picture_1.jpeg)

#### It's Craig Venter!

![](_page_58_Picture_3.jpeg)

![](_page_58_Picture_4.jpeg)

#### Identifying Personal Genomes by Surname Inference

Gymrek et al (2013) Science. doi: 10.1126/science.1229566

#### Possible route for identity tracing

![](_page_59_Figure_1.jpeg)

- US population: ~313.9 million individuals
- log<sub>2</sub> 313,900,000 = 28.226 bits
- Sex ~ 1.0 information bits
- log<sub>2</sub> 156,950,000 = 27.226 bits

- Tracing attacks combine metadata and surname inference to triangulate the identity of an unknown individual.
- With no information, there are roughly 300 million matching individuals in the US, equating to 28.0 bits of entropy.
- Sex reduces entropy by 1 bit, state of residence and age reduces to 16, successful surname inference reduces to ~3 bits.

### The risks of big data?

#### Predicting Social Security numbers from public data

Alessandro Acquisti<sup>1</sup> and Ralph Gross

Carregie Mellon University, Pittsburgh, PA 15213

Communicated by Stephen E. Fienberg, Carnegie Mellon University, Pittsburgh, PA, May 5, 2009 (received for review January 18, 2009)

Information about an individual's place and date of birth can be exploited to predict his or her Social Security number (SSN). Using only publicly available information, we observed a correlation between individuals' SSNs and their birth data and found that for younger cohorts the correlation allows statistical inference of private SSNs. The inferences are made possible by the public availability of the Social Security Administration's Death Master

File and the widespread accessibility of person multiple sources, such as data brokers or pro working sites. Our results highlight the unexp sequences of the complex interactions are sources in modern information economies an risks associated with information revelation in

identity theft | online social networks | privacy | stati

n modern information economies, sensitive p plain sight amid transactions that rely on their their unhindered circulation. Such is the case v numbers in the United States: Created as iden tracking individual earnings (1), they have tu authentication devices (2), becoming one of the tion most often sought by identity thieves. T Administration (SSA), which issues them, has u keep SSNs confidential (3), coordinating with 1 their public exposure (4).\* After embarrassin sector entities also have attempted to strengthe their consumers' and employees' data (7).\* How have already left the barn: We demonstrate the number (SN). The SSA openly provides information about the process through which ANs, GNs, and SNs are issued (1). ANs are currently assigned based on the zipcode of the mailing address provided in the SSN application form [RM00201.030] (1). Low-population states and certain U.S. possessions are allocated 1 AN each, whereas other states are allocated sets of ANs (for instance, an individual arching from a zincode within

publish on social networking sites (10). Using this method, we identified with a single attempt the first 5 digits for 44% of DMF records of deceased individuals born in the U.S. from 1989 to 2003 and the complete SSNs with <1,000 attempts (making SSNs akin to 3-digit financial PINs) for 8.5% of those records. Extrapolating to the U.S. living population, this would imply the potential identification of millions of SSNs for individuals whose birth data were available. Such findings highlight the hidden privacy costs of widespread information dissemination and the complex interactions among multiple data sources in modern information economies (11), underscoring the role of public records as breeder documents (12) of more sensitive data.

then addressed

APPLICATION