### GWAS across Time and Space What have we learned?

Arcadi Navarro

### BIOSTEC

March the 4<sup>th</sup> 2014



### Warning



### Warning



### Warning



#### Where?



#### Where?



#### Who?

#### The Navarro lab, 2013





- Longish Introduction. Why Genomics? Association studies?
  Genome Wide Association Studies (GWAS)
- 2. Candidate-Gene studies and their problems
- 3. Genome-Wide association studies
- 4. Whatever happened to personalized medicine

#### The Human Genome and Disease

The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!" ... It was the reason under the Human Genome Project



#### **The Human Genome and Disease**

## The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!" ... It was the reason under the Human Genome Project

-----

BMC 0	Genomics 2008, 9:472	http://www.biomedcentral.com/1471-2164/9/472						
Table I	· Conomics Funding by Operation and Yoak in US	s (millio	<b>n</b> r)					
Rank	Organization	2003	2004	2005	2006	Source		
n/a	Department of Health and Human Services (DHHS) National Institutes of Health (NIH) – GENETICS <sup>O</sup>	\$4236	\$4535	\$4840	\$4878	[28]		
I	NIH: National Cancer Institute (NCI) + National Human Genome Research Institute (NHGRI)	\$562	\$593	\$593	\$571	See Tables 2, 3 and text		
2	European Commission	\$459	\$462	\$466	\$468	Personal Communication, Indridi Benediktsson, September 2006		
3	European Commission Matching Funds	\$459	\$462	\$466	\$468	Personal Communication, Indridi Benediktsson, September 2006		
4	United Kingdom Wellcome Trust <sup>o</sup>	\$194	\$194	\$208	\$199	[52,53]		
5	Department of Energy (DOE) Office of Biological and Environmental Research <sup>O</sup>	\$129	\$152	\$154	\$158	[33-36]		
6	National Science Foundation (NSF) Biological Sciences Directorate <sup>O</sup>	\$124	\$129	\$134	\$141	Personal Communication, Vernon Ross, April 2007		
7	Japan Ministry of Education, Culture, Sports, Science and Technology (MEXT) <sup>A</sup>	\$84.5	\$99.2	\$119	\$125	Personal Communication, Kazuko Shinohara, January 2007		
8	United Kingdom Biotechnology and Biological Sciences Research Council (BBSRC) <sup>A</sup>	\$121	\$127	\$128	\$117	Personal Communication, Clare Nixon, January 2007		
9	Genome Canada <sup>A</sup>	\$67.3	\$65.7	\$71.6	\$106	Personal Communication, Genny Cardin, July 2006		
10	China (Ministry of Science and Technology, National Natural Science Foundation of China, and Chinese Academy of Sciences)	\$80	\$80	\$80	no report	Personal Communication, Anonymous, October 2007		
н	Germany Nationales Genomforschungsnetz (NGFN)	\$71.6	\$61.4	\$62.0	\$64.8	Personal Communication, Uta Strasser, September 2006		
12	Department of Defense (DOD) Congressional Directed Medical Research Programs (CDMRP) <sup>O</sup>	\$102	\$86.8	\$53.5	\$54.9	[45]		
13	Cancer Research UK <sup>A</sup>	\$34.1	\$45.2	\$48.1	\$51.0	Personal Communication, Lynne Davies, January 2007		
14	Netherlands Genomics Initiative <sup>J</sup>	\$48.8	\$18.1	\$51.7	\$45.8	[31,87,88]		
15	South Korea Ministry of Science and Technology (MOST) <sup>j</sup>	\$39.7	\$35.0	\$41.5	\$44.3	Personal Communication, Jeongheui Lim, January 2007		
16	United States Department of Agriculture (USDA) Agricultural Research Service <sup>D</sup>	\$32.5	\$38.8	\$41.9	\$43.1	Personal Communication, Peggy DelCollo and Joe Garbarino, November 2006		
17	Ireland Higher Education Authority	\$34.1	\$34.7	\$34.5	\$34.5	Personal Communication, Sorcha Carthy, August 2006		
18	Canada Natural Sciences and Engineering Research Council (NSERC) <sup>A</sup>	\$28.2	\$30.4	\$32.2	\$34.1	Personal Communication, Barney Laciak, October 2007		
19	Department of Homeland Security (DHS) <sup>O</sup>	\$13.4	\$25.9	\$32.8	\$27.2	Personal Communication, Elizabeth George, November 2006		

ble	I: Genomics Funding by Organization and Year, in US	\$ (millio	ns) (Cont	inued)		
20	Japan Ministry of Agriculture, Forestry and Fisheries (MAFF) <sup>A</sup>	\$17.9	\$22.9	\$21.7	\$21.7	Personal Communication, Kazuko Shinohara, January 2007
21	Canadian Biotechnology Strategy (CBS) <sup>A</sup>	\$16.1	\$15.9	\$15.9	\$16.2	[55]
22	Howard Hughes Medical Institute (HHMI) <sup>j</sup>	\$15.2	\$13.8	\$14.3	\$15.6	Personal Communication, Sherry White, August 2006
23	Canada National Research Council (NRC) Genomics and Health Initiative <sup>A</sup>	\$18.2	\$15.8	\$13.4	\$15.3	Personal Communication, Gary Fudge August 2007
24	Japan Ministry of Health, Labour and Welfare (MHLW)^A	\$12.8	\$13.8	\$18.6	\$14.9	Personal Communication, Kazuko Shinohara, January 2007
25	American Cancer Societyl	\$5.95	\$6.90	\$5.61	\$11.9	Personal Communication, Donella Wilson, July 2006
26	Spain Genoma Espanal	\$9.98	\$12.1	\$13.5	\$11.7	Personal Communications, Javier Montero Plata, June and September 2006
27	Australia National Health and Medical Research Council (NHMRC)*	\$7.54	\$6.34	\$6.24	\$7.65	Personal Communication, Marian Blake, July 2006
28	Department of Health and Human Services (DHHS) Centers for Disease Control and Prevention (CDC) <sup>O</sup>	\$3.85	\$4.53	\$6.99	\$6.95	[37,38]
29	Department of Defense (DOD) Defense Advanced Research Projects Agency (DARPA) Bio/Info/Micro Program <sup>O</sup>	\$34.0	\$13.4	\$13.2	\$6.5	[41-44]
30	South African Medical Research Council	\$1.73	\$1.84	\$2.04	\$2.24	Personal Communication, Clive Glass, October 2007
31	Ireland Science Foundation <sup>J</sup>	\$7.41	\$7.96	\$1.38	\$2.24	Personal Communication, Tracy Moloney, September 2006
32	Ireland Health Research Board	\$1.29	\$2.23	\$1.50	\$1.61	Personal Communication, Gillian Hastings, January 2007
33	South Africa National Research Foundation	\$0.658	\$0.560	\$0.999	no report	Personal Communication, Marna van Rooyen, October 2006
34	Japan Ministry of Economy, Trade and Industry (METI)A	\$0	\$4.51	\$0	\$0	Personal Communication, Kazuko Shinohara, January 2007
	TOTAL	\$2834	\$2878	\$2948	\$2881	

Genomics research funded, by organization, is shown in millions of US\$ per year. The total funding values each year were determined by summing the values from each organization, with the exception of the first row (NHH – Genetics), which is described in the text. Rankings were determined by ordering the 2006 values, where the average of the three previous years was used as a substitute for 2006 values when the actual 2006 data was unavailable. The start of each fiscal year is indicated by the superscript character after each organization, where J = January I, A = April I, and O = October I.

\*Although the fiscal year for the Australian government begins July I, the original data was reported by calendar year.



#### A deluge of developments since then

The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!" ... has been the motivation under many advances.



ED Green et al. Nature 470, 204-213 (2011) doi:10.1038/nature09764

The promise in the 90s: "In 10 years we will unravel the genetic basis of complex diseases!!..."

- (1) HapMap & 1000 Genomes Project
- (2) Genotyping arrays
- (3) Large samples

Patterns of LD / variation in many human populations

Dense coverage of human genome

Case-Control or Family Linkage

上房医 下医医已病

"Superior Doctors Prevent the Disease. Mediocre Doctors Treat the Disease Before Evident. Inferior Doctors Treat the Full Blown Disease."

-Huang Dee: Nai - Ching (2600 B.C. 1st Chinese Medical Text

The promise in the 90s: "In 10 years we will unravel the genetic basis of complex diseases!!..."

(1) HapMap & 1000 Genomes Project	Patterns of LD / variation in many human populations
(2) Genotyping arrays	Dense coverage of human genome
(3) Large samples	Case-Control or Family Linkage

- Prevent disease from occurring
- Identify the cause of the disease
- Treat the cause of the disease rather than the symptoms
- Genomics may identify the cause of disease ("All medicine may become pediatrics" Paul Wise, Professor of Pediatrics, Stanford Medical School, 2008)
- Effects of environment, accidents, aging, penetrance ...
- Health care costs can be greatly reduced if
  - invests in preventive medicine
  - o one targets the cause of disease rather than symptoms

## The promise in the 90s: "In 10 years we will unravel the genetic basis of complex diseases!!..."



## The promise in the 90s: "In 10 years we will unravel the genetic basis of complex diseases!!..."



The remarkable advances in molecular biology during the past two decades have given man an understanding of the basic processes that shape his life and have placed within the realm of possibility medical achievements undreamed of a scant few years ago.

The way is being opened not only for permanent cures of genetic disease but also for drastic changes in man's genetic makeup.

http://www.time.com/time/covers/0,16641,19710419,00.html Buchanan AV et al, Int J Epidemiol, June 2006 (35: 593-596)

**Genetics**:

#### **COMPLEX GENETIC ARCHITECTURE OF COMPLEX DISEASES**

- ✓ Polygenic + Environmental
- ✓ High heritability ( $h^2 > 50\%$ ) and Familial aggregation ( $\lambda_s >> 1$ )
- ✓ But lack of Mendelian inheritance (lots of sporadic cases)



Manolio et al. (2008) J Clin Invest. 118(5):1590

#### **COMPLEX GENETIC ARCHITECTURE OF COMPLEX DISEASES**

The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!". Everything seemed to go well...

... and people were doing this using diverse approaches.

#### ✓ Linkage mapping:

Select an informative family.

More than 100,000 papers in 30 years (mostly in OMIM)

#### ✓ Candidate genes:

Select a few variants in one or a few genes.

More than 84000 papers in 20 years (see the GAD database)

✓ Genome-Wide Association Studies:

Hypothesis-free approach. One looks as as many variables as possible More than 1200 papers (see www.genome.gov)

#### THE GENETIC ARCHITECTURE OF MENDELIAN DISEASES

The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!". Everything seemed to go well...

### **OMIM Home Page**

#### http://www.ncbi.nlm.nih.gov/omim/



sorders, by genetics researchers, and by advanced students in science and medicine. While the OMIM datab

#### **COMPLEX GENETIC ARCHITECTURE OF COMPLEX DISEASES**

The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!" ... And it seemed easy to fulfil...



Odds Ratio: 3.6 95% CI = 1.3 to 10.4



Consortium

#### **Many Candidate Gene Studies**

The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!" ... And it seemed easy to fulfil...

Disease

Gene View

CH-SNP-HapMap Reference

invironmental Factor Gene Interaction



Odds Ratio: 3.6 95% CI = 1.3 to 10.4

#### The Genetic Association Database (GAD)

The Genetic Association Database is an archive of human genetic association studies of complex diseases and disorders. The goal of this database is to allow the user to rapidly identify medically relevant polymorphism from the large volume of polymorphism and mutational data, in the context of standardized nomenclature.

The data is from published scientific papers. Study data is recorded in the context of official human gene nomenclature with additional molecular reference numbers and links. It is gene centered. That is, each record is a record of a gene or marker. If a study investigated 6 genes for a particular disorder, there will be 6 records.

Next 25										All View	S	earch for DIABETES	Record found:	500	
	Last Update	Checked	CDC Index	1 - GAD 2 - CDC	Year	Assoc? YorN	Gene Symbol	омім	Gene Expert	Broad Phenotype (Disease)	Disease Expert	MeSH Disease Terms	Disease Class	Chr	Ch-Band
viev	13-APR-07		17626	2	2005		CTLA4	123890	Сx	diabetes, type 1			IMMUNE	2	2q33
viev	13-APR-07		16890	2	2005		GCK	138079	C <sub>x</sub>	diabetes, type 2			METABOLIC	7	7p15.3-p15.1
view	13-APR-07		8711	2	2006		AP0A5	606368		myocardial infarct; diabetes,			CANCER	11	11q23
viev	13-APR-07		16322	2	2005	N	D102	601413		body mass; diabetes, type 2; g		Diabetes Mellitus, Type 2 Insu	METABOLIC	14	14q24.2-q24.3
viev	13-APR-07		11253	2	2005		CTLA4	123890	<b>C</b> x	diabetes, type 1; thyroid dise			IMMUNE	2	2q33
viev	13-APR-07		15560	2	2005		CAPN10	605286	C <sub>x</sub>	diabetes, type 2; nephropathy,			METABOLIC	2	2 q37.3

#### **Candidate Gene Association Studies**

The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!". Everything seemed to go well...

(1) HapMap & 1000 Genomes Project	Patterns of LD / variation in many human populations
(2) Genotyping arrays	Dense coverage of human genome
(3) Large samples	Case-Control or Family Linkage
However Statistical Problems	Sample size = Power Multiple testing
Genetic problems	Ascertainment bias Locus & Allelic heterogeneity Population substructure

As a consequence, there were (are?) LACK OF REPLICATION problems!!

### **Replication validity of genetic association studies**

John P.A. Ioannidis<sup>1–3</sup>, Evangelia E. Ntzani<sup>1</sup>, Thomas A. Trikalinos<sup>1</sup> & Despina G. Contopoulos-Ioannidis<sup>1,4</sup>

nature genetics • volume 29 • november 2001



So the BULK of old candidate gene associations are not good

More than 100,000 papers registered in the GAD...but... the field is dominated by false positives (getting better recently)

The precision of many old papers is similar to this:

"I have nothing to offer but blood, toil, tears and sweat."

Mahatma Gandhi



#### But not everything was bad...

What about focusing on patters of lack of replicability rather than on good replications?



#### But not everything was bad...

Increased inconsistency of replicability ( $\phi$ ) with increased F<sub>ST</sub> (i.e. higher consistency between Europe and Asia with lower genetic distance)



Correlation between discordance in replicability and  $F_{ST}$  for the 37 associations from the Continental Set.

Marigorta et al. BMC Genomics 2011, 12:55 http://www.biomedcentral.com/1471-2164/12/55



**Open Access** 

#### RESEARCH ARTICLE

## Recent human evolution has shaped geographical differences in susceptibility to disease

Urko M Marigorta<sup>1</sup>, Oscar Lao<sup>2</sup>, Ferran Casals<sup>1</sup>, Francesc Calafell<sup>13</sup>, Carlos Morcillo-Suárez<sup>1,4</sup>, Rui Faria<sup>1,5</sup>, Elena Bosch<sup>1,3</sup>, François Serra<sup>6</sup>, Jaume Bertranpetit<sup>1,3</sup>, Hernán Dopazo<sup>6</sup>, Arcadi Navarro<sup>1,4,7\*</sup>

#### Abstract

Background: Searching for associations between genetic variants and complex diseases has been a very active area of research for over two decades. More than 51,000 potential associations have been studied and published, a figure that keeps increasing, especially with the recent explosion of aray-based Genome-Wide Association Studies. Even if the number of true associations described so far is high, many of the putative risk variants detected so far have failed to be consistently replicated and are widely considered false positives. Here, we focus on the worldwide patterns of replicability of published association studies.

Results: We report three main findings. First, contrary to previous results, genes associated to complex diseases present lower degrees of genetic differentiation among human populations than average genome-wide levels. Second, also contrary to previous results, the differences in replicability of disease associated-loci between Europeans and East Asians are highly correlated with genetic differentiation between these populations. Finally, highly replicated genes present increased levels of high-frequency derived alleles in European and Asian populations when compared to African populations.

Conclusions: Our findings highlight the heterogeneous nature of the genetic etiology of complex disease, confirm the importance of the recent evolutionary history of our species in current patterns of disease susceptibility and could cast doubts on the status as false positives of some associations that have failed to replicate across populations.

#### WORLD COLONIZATION BY HUMANS (last ~100 kyears)



#### **COMPLEX GENETIC ARCHITECTURE OF COMPLEX DISEASES**

The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!" ... And it seemed easy to fulfil...



Odds Ratio: 3.6 95% CI = 1.3 to 10.4



Consortium

#### **Plus many GWAS**

## The promise in the 90s: "In 10 years we will unravel the genetic bases of complex diseases!!" ... And it seemed easy to fulfil...



0000

1185

The CARDIA-DENEVA Stud



Consortium









 $\bigcirc$ Coffee consumption  $\bigcirc$ Cognitive function Abdominal aortic aneurysm Acute lymphoblastic leukemia  $\bigcirc$ Conduct disorder  $\bigcirc$ Colorectal cancer Corneal thickness  $\bigcirc$ Coronary disease Age-related macular degeneration Cortical thickness  $\bigcirc$  $\bigcirc$ Creutzfeldt-Jakob disease  $\bigcirc$ Crohn's disease  $\bigcirc$ Crohn's disease and celiac disease Cutaneous nevi Cystic fibrosis severity Amyotrophic lateral sclerosis Dermatitis Angiotensin-converting enzyme activity DHEA-s levels Diabetic retinopathy Dilated cardiomyopathy  $\bigcirc$ Drug-induced liver injury Drug-induced liver injury (amoxicillin-clavulanate)  $\bigcirc$ Endometrial cancer Endometriosis Attention deficit hyperactivity disorder  $\bigcirc$  $\bigcirc$ Eosinophil count Eosinophilic esophagitis  $\bigcirc$ Epirubicin-induced leukopenia Erectile dysfunction and prostate cancer treatment Erythrocyte parameters Esophageal cancer Essential tremor Exfoliation glaucoma Eve color traits F cell distribution Fibrinogen levels Folate pathway vitamins Follicular lymphoma Fuch's corneal dystrophy Freckles and burning  $\bigcirc$ Gallstones Butyrylcholinesterase levels  $\bigcirc$ Gastric cancer Glioma  $\bigcirc$ Glycemic traits  $\bigcirc$ Graves disease O Hair color Cardiovascular risk factors Hair morphology Handedness in dyslexia Carotenoid/tocopherol levels O HDL cholesterol O Heart failure  $\bigcirc$ Celiac disease and rheumatoid arthritis Heart rate Cerebral atrophy measures O Height O Hemostasis parameters Chronic lymphocytic leukemia Hepatic steatosis

O Hepatitis

 $\bigcirc$ 

 $\bigcirc$ 

 $\bigcirc$ 

 $\bigcirc$ 

 $\bigcirc$ 

 $\bigcirc$ 

 $\bigcirc$ 

 $\bigcirc$ Asthma

 $\bigcirc$ 

 $\bigcirc$ Autism

 $\bigcirc$ 

 $\bigcirc$ 

Adhesion molecules

Adiponectin levels

AIDS progression

Alcohol dependence

Alopecia areata

Amyloid A levels

Arterial stiffness

Atrial fibrillation

Basal cell cancer

Behcet's disease

Bitter taste response

Bleomycin sensitivity

Blond or brown hair

Blue or green eyes

BMI, waist circumference

Cardiac structure/function

Carotid atherosclerosis

Chronic myeloid leukemia

Bipolar disorder

Biliary atresia

Bladder cancer

Blood pressure

Bone density

Breast cancer

C-reactive protein

Calcium levels

Carnitine levels

Celiac disease

Cleft lip/palate

Bilirubin

Birth weight

Alzheimer disease

Ankylosing spondylitis

Asparagus anosmia

Atherosclerosis in HIV

Hepatitis B vaccine response Neuroblastoma Hepatocellular carcinoma  $\bigcirc$ O Hirschsprung's disease  $\bigcirc$ Obesity O HIV-1 control  $\bigcirc$ Hodgkin's lymphoma  $\bigcirc$ Homocysteine levels  $\bigcirc$ HPV seropositivity Ο O Hypospadias Idiopathic pulmonary fibrosis  $\bigcirc$ IFN-related cytopeni  $\bigcirc$ IgA levels IaE levels  $\bigcirc$ Pain Inflammatory bowel disease Insulin-like growth factors  $\bigcirc$  $\bigcirc$ Intracranial aneurysm  $\bigcirc$  $\bigcirc$ Iris color  $\bigcirc$ Iron status markers Ischemic stroke  $\bigcirc$  $\bigcirc$ Ο Juvenile idiopathic arthritis  $\bigcirc$ Keloid  $\bigcirc$ Kidney stones  $\bigcirc$ LDL cholesterol  $\bigcirc$ Leprosy Leptin receptor levels  $\bigcirc$ Liver enzymes  $\bigcirc$  $\bigcirc$ Longevity  $\bigcirc$ LP (a) levels  $\bigcirc$  LpPLA(2) activity and mass  $\bigcirc$ Lung cancer  $\bigcirc$ Magnesium levels  $\bigcirc$ Major mood disorders  $\bigcirc$ Malaria  $\bigcirc$ O Male pattern baldness O Mammographic density  $\bigcirc$  Matrix metalloproteinase levels  $\bigcirc$ O MCP-1  $\bigcirc$ Melanoma  $\bigcirc$ Menarche & menopause Meningioma  $\bigcirc$ Meningococcal disease  $\bigcirc$ Metabolic syndrome  $\bigcirc$ Migraine  $\bigcirc$ Movamova disease  $\bigcirc$ Multiple sclerosis Myeloproliferative neoplasms Myopia (pathological) Ο N-glycan levels  $\bigcirc$ O Narcolepsy  $\bigcirc$ O Nasopharyngeal cancer

Natriuretic peptide levels

Nicotine dependence Open angle glaucoma Open personality Optic disc parameters Osteoarthritis Osteoporosis Otosclerosis Other metabolic traits Ovarian cancer Pancreatic cancer Paget's disease Panic disorder Parkinson's disease Periodontitis Peripheral arterial disease Personality dimensions Phosphatidylcholine levels Phosphorus levels Photic sneeze Phytosterol levels Platelet count Polycystic ovary syndrome Primary biliary cirrhosis Primary sclerosing cholangitis PR interval Progranulin levels Progressive supranuclear palsy Prostate cancer Protein levels PSA levels Psoriasis Psoriatic arthritis Pulmonary funct. COPD QRS interval QT interval Quantitative traits Recombination rate Red vs non-red hair Refractive error Renal cell carcinoma Renal function Response to antidepressants Response to antipsychotic therapy Response to carbamazepine Response to clopidogrel therapy Response to hepatitis C treat  $\bigcirc$ Response to interferon beta therapy

Response to metaformin Response to statin therapy Restless legs syndrome Retinal vascular caliber  $\bigcirc$ Retinol levels Rheumatoid arthritis  $\bigcirc$ Ribavirin-induced anemia Schizophrenia  $\bigcirc$  $\bigcirc$ Serum metabolites  $\bigcirc$ Skin pigmentation Smoking behavior  $\bigcirc$  $\bigcirc$ Speech perception Sphingolipid levels  $\bigcirc$ Statin-induced myopathy  $\bigcirc$ Stevens-Johnson syndrome  $\bigcirc$ Stroke  $\bigcirc$  $\bigcirc$ Sudden cardiac arrest  $\bigcirc$ Suicide attempts  $\bigcirc$ Systemic lupus erythematosus Systemic sclerosis Ο  $\bigcirc$ T-tau levels  $\bigcirc$ Tau AB1-42 levels Telomere length Testicular germ cell tumor Thyroid cancer  $\bigcirc$  $\bigcirc$ Thyroid volume Tooth development Total cholesterol Triglycerides  $\bigcirc$ Tuberculosis Type 1 diabetes Type 2 diabetes  $\bigcirc$  $\bigcirc$ Ulcerative colitis  $\bigcirc$ Urate  $\bigcirc$ Urinary albumin excretion Urinary metabolites  $\bigcirc$ Uterine fibroids Venous thromboembolism Ventricular conduction **VEGF** levels 0 Vertical cup-disc ratio Vitamin B12 levels Vitamin D insuffiency Vitamin E levels  $\bigcirc$ Vitiligo Warfarin dose Weight  $\bigcirc$ White cell count White matter hyperintensity  $\bigcirc$ YKL-40 levels

#### Problem. Where did the heritability go?

### The missing heritability problem

#### Nature Feature Nov 2008



#### The case of the missing heritability

When scientists opened up the human genome, they expected to find the genetic components of common traits and diseases. But they were nowhere to be seen. **Brendan Maher** shines a light on six places where the missing loot could be stashed away.

When scientists opened up the human genome, they expected to find the genetic components of common traits and diseases. But they were nowhere to be seen. Brendan Maher shines a light on six places where the missing loot could be stashed away.

The case of the missing heritability

Manolio et al Nature October 2009

#### Vol 461 8 October 2009 doi:10.1038/nature08494

#### nature

#### REVIEWS

### Finding the missing heritability of complex diseases

Teri A. Manolio<sup>1</sup>, Francis S. Collins<sup>2</sup>, Nancy J. Cox<sup>3</sup>, David B. Goldstein<sup>4</sup>, Lucia A. Hindorff<sup>5</sup>, David J. Hunter<sup>6</sup>, Mark I. McCarthy<sup>7</sup>, Erin M. Ramos<sup>5</sup>, Lon R. Cardon<sup>8</sup>, Aravinda Chakravarti<sup>9</sup>, Judy H. Cho<sup>10</sup>, Alan E. Guttmacher<sup>1</sup>, Augustine Kong<sup>11</sup>, Leonid Kruglyak<sup>12</sup>, Elaine Mardis<sup>13</sup>, Charles N. Rotimi<sup>14</sup>, Montgomery Slatkin<sup>15</sup>, David Valle<sup>9</sup>, Alice S. Whittemore<sup>16</sup>, Michael Boehnke<sup>17</sup>, Andrew G. Clark<sup>18</sup>, Evan E. Eichler<sup>19</sup>, Greg Gibson<sup>20</sup>, Jonathan L. Haines<sup>21</sup>, Trudy F. C. Mackay<sup>22</sup>, Steven A. McCarroll<sup>23</sup> & Peter M. Visscher<sup>24</sup>

Genome-wide association studies have identified hundreds of genetic variants associated with complex human diseases and traits, and have provided valuable insights into their genetic architecture. Most variants identified so far confer relatively small increments in risk, and explain only a small proportion of familial clustering, leading many to question how the remaining, 'missing' heritability can be explained. Here we examine potential sources of missing heritability and propose research strategies, including and extending beyond current genome-wide association approaches, to illuminate the genetics of complex diseases and enhance its potential to enable effective disease prevention or treatment.

traits, and have provided valuable insights into their genetic architecture. Most variants identified so far confer relatively small increments in risk, and explain only a small proportion of familial clustering, leading many to question how the remaining, "missing" heritability can be explained. Here we examine potential sources of missing heritability and propose research strategies, including and extending beyond current genome-wide association approaches, to illuminate the genetics of complex diseases and enhance its potential to enable effective disease prevention or treatment.

#### Problem. Where did the heritability go?

### The missing heritability problem

Disease	Number of loci	Proportion of heritability explained
Age-related macular degeneration <sup>72</sup>	5	50%
Crohn's disease <sup>21</sup>	32	20%
Systemic lupus erythematosus73	6	15%
Type 2 diabetes <sup>74</sup>	18	6%
HDL cholesterol <sup>75</sup>	7	5.2%
Height <sup>15</sup>	40	5%
Early onset myocardial infarction <sup>76</sup>	9	2.8%
Fasting glucose <sup>77</sup>	4	1.5%

#### Problem. Where did the heritability go?

#### Consequences of the "missing heritability" problem



Part of the heritability has been discovered...but most remains "missing"!

### Does this challenge the CV/CD paradigm? Are they rare variants? Are GWAS results artifacts?

### How to explain all this?

Think Complex Diseases:



**Figure 1** | **Feasibility of identifying genetic variants by risk allele frequency and strength of genetic effect (odds ratio).** Most emphasis and interest lies in identifying associations with characteristics shown within diagonal dotted lines. Adapted from ref. 42.



#### How to explain all this?

Rare Variants, rare CNVs, epigenetics....?

#### ... but the signal may be there!! Heritability may be Hidden, not Missing

# Common SNPs explain a large proportion of the heritability for human height

Jian Yang<sup>1</sup>, Beben Benyamin<sup>1</sup>, Brian P McEvoy<sup>1</sup>, Scott Gordon<sup>1</sup>, Anjali K Henders<sup>1</sup>, Dale R Nyholt<sup>1</sup>, Pamela A Madden<sup>2</sup>, Andrew C Heath<sup>2</sup>, Nicholas G Martin<sup>1</sup>, Grant W Montgomery<sup>1</sup>, Michael E Goddard<sup>3</sup> & Peter M Visscher<sup>1</sup>

SNPs discovered by genome-wide association studies (0 account for only a small fraction of the genetic variation complex traits in human populations. Where is the remaindent heritability? We estimated the proportion of variance for human height explained by 294,831 SNPs genotyped or 3,925 unrelated individuals using a linear model analysi validated the estimation method with simulations based the observed genotype data. We show that 45% of varia can be explained by considering all SNPs simultaneous most of the heritability is not missing but has not previo been detected because the individual effects are too sm to pass stringent significance tests. We provide evidence that the remaining heritability is due to incomplete link disequilibrium between causal variants and genotyped exacerbated by causal variants having lower minor allel frequency than the SNPs explored to date.

nature

disequilibrium between causal variants and genotyped Si exacerbated by causal variants having lower minor allele frequency than the SNPs explored to date. genetics

### Genome partitioning of genetic variation for complex traits using common SNPs

ANALYSIS

ANALYSIS

Jian Yang<sup>1\*</sup>, Teri A Manolio<sup>2</sup>, Louis R Pasquale<sup>3</sup>, Eric Boerwinkle<sup>4</sup>, Neil Caporaso<sup>5</sup>, Julie M Cunningham<sup>6</sup>, Mariza de Andrade<sup>7</sup>, Bjarke Feenstra<sup>8</sup>, Eleanor Feingold<sup>9</sup>, M Geoffrey Hayes<sup>10</sup>, William G Hill<sup>11</sup>, Maria Teresa Landi<sup>12</sup>, Alvaro Alonso<sup>13</sup>, Guillaume Lettre<sup>14</sup>, Peng Lin<sup>15</sup>, Hua Ling<sup>16</sup>, William Lowe<sup>17</sup>, Rasika A Mathias<sup>18</sup>, Mads Melbye<sup>8</sup>, Elizabeth Pugh<sup>16</sup>, Marilyn C Cornelis<sup>19</sup>, Bruce S Weir<sup>20</sup>, Michael E Goddard<sup>21,22</sup> & Peter M Visscher<sup>1</sup>

disequilibrium between causal variants and genotyped SNPs, michael E Goddard<sup>21,22</sup> & Peter M VisScher<sup>1</sup> - 50 variants that are

#### From almost no heritability to a lot in a 3 years!

#### From 2009 to 2012 (Manolio et al to Visscher et al)

Table 1   Estimates of heritability and number of loci for several complex traits							
Disease	Number of loci	Proportion of heritability explained					
Age-related macular degeneration <sup>72</sup>	5	50%					
Crohn's disease <sup>21</sup>	32	20%					
Systemic lupus erythematosus73	6	15%					
Type 2 diabetes <sup>74</sup>	18	6%					
HDL cholesterol <sup>75</sup>	7	5.2%					
Height <sup>15</sup>	40	5%					
Early onset myocardial infarction <sup>76</sup>	9	2.8%					
Fasting glucose <sup>77</sup>	4	1.5%					
* Residual is after adjustment for age, gender, diabetes	L						
* Residual is after adjustment for age, gender, diabeter							
Fasting glucose"							

Table 1. Population Variation Explained by GWAS for a Selected        Number of Complex Traits							
Trait or Disease	h <sup>2</sup> Pedigree Studies	h <sup>2</sup> GWAS Hits <sup>a</sup>	h <sup>2</sup> All GWAS SNPs <sup>b</sup>				
Type 1 diabetes	0.9 <sup>98</sup>	0.6 <sup>99 ,c</sup>	0.3 <sup>12</sup>				
Type 2 diabetes	0.3-0.6 <sup>100</sup>	$0.05 - 0.10^{34}$					
Obesity (BMI)	0.4–0.6 <sup>101,102</sup>	$0.01 - 0.02^{36}$	0.2 <sup>14</sup>				
Crohn's disease	0.6-0.8 <sup>103</sup>	$0.1^{11}$	0.4 <sup>12</sup>				
Ulcerative colitis	0.5 <sup>103</sup>	$0.05^{12}$					
Multiple sclerosis	$0.3 - 0.8^{104}$	$0.1^{45}$					
Ankylosing spondylitis	>0.90 <sup>105</sup>	$0.2^{106}$					
Rheumatoid arthritis	0.6 <sup>107</sup>						
Schizophrenia	$0.7 - 0.8^{108}$	0.01 <sup>79</sup>	0.3 <sup>109</sup>				
Bipolar disorder	0.6-0.7 <sup>108</sup>	0.0279	0.4 <sup>12</sup>				
Breast cancer	0.3 <sup>110</sup>	$0.08^{111}$					
Von Willebrand factor	0.66-0.75 <sup>112,113</sup>	0.13 <sup>114</sup>	0.2514				
Height	0.8 <sup>115,116</sup>	$0.1^{13}$	0.5 <sup>13,14</sup>				
Bone mineral density	0.6-0.8 <sup>117</sup>	$0.05^{118}$					
QT interval	0.37-0.60119,120	$0.07^{121}$	0.214				
HDL cholesterol	0.5 <sup>122</sup>	0.157					
Platelet count	$0.8^{123}$	0.05-0.158					

<sup>a</sup> Proportion of phenotypic variance or variance in liability explained by genome-wide-significant and validated SNPs. For a number of diseases, other parameters were reported, and these were converted and approximated to the scale of total variation explained. Blank cells indicate that these parameters have not been reported in the literature.

<sup>b</sup> Proportion of phenotypic variance or variance in liability explained when all GWAS SNPs are considered simultaneously. Blank cell indicate that these parameters have not been reported in the literature.

<sup>c</sup> Includes pre-GWAS loci with large effects.

includes pre-owns for with large effects.









#### **GWAS** across Space: must be Shared Variants

#### We just proved high correlation between European / East Asian variants

Marigorta and Navarro, Plos Genetics 2013



#### **GWAS** across Space: must be Shared Variants

#### And that effective replicability depends on statistical power

Marigorta and Navarro, Plos Genetics 2013



#### Which means that GWAS result so far must be due to high-frequency disease variants that are shared by all humankind

50.21% Private European-Americans Shared Private African-Americans

Europeans\* Europeans & Chinese<sup>†</sup> European & African<sup>‡</sup> 100 80 between populations (%) 60 40 20 Rare variants (1%) Common variants (15%)

CD Bustamante et al. Nature 475 (2011)

- **Common variants** are (usually) shared among continental populations (MAF > 0.1)
- **Rare variants** are mostly population-specific (MAF < 0.01)



Casals and Bertranpetit Science 6090 (2012)

#### But the raging debate is still open in many aspects!!

#### The mystery of missing heritability: Genetic interactions create phantom heritability

Or Zuk<sup>a</sup>, Eliana Hechter<sup>a</sup>, Shamil R. Sunyaev<sup>a,b</sup>, and Eric S. Lander<sup>a,1</sup>

Broad Institute of MIT and Harvard, Cambridge, MA 02142; and <sup>b</sup>Genetics Division, Brigham and Women's Hospital, Harvard Medical School, Boston, MA 02115

Contributed by Eric S. Lander, December 5, 2011 (sent for review October 9, 2011)

Human genetics has been haunted by the mystery of "missing heritability" of common traits. Although studies have discovered > 1,200 variants associated with common diseases and traits, these variants typically appear to explain only a minority of the heritability. The proportion of heritability explained by a set of variants is the ratio of (i) the heritability due to these variants (numerator), estimated directly from their observed effects, to (ii) the total heritability (denominator), inferred indirectly from population data. The prevailing view has been that the explanation for missing heritability lies in the numerator-that is, in as-yet undiscovered variants. While many variants surely remain to be found, we show here that a substantial portion of missing heritability could arise from overestimation of the denominator, creating "phantom heritability." Specifically, (i) estimates of total heritability implicitly assume the trait involves no genetic interactions (epistasis) among loci; (ii) this assumption is not justified, because models with interactions are also consistent with observable data; and (iii) under such models, the total heritability may be much smaller and thus the proportion of heritability explained much larger. For example, 80% of the currently missing heritability for Crohn's disease could be due to genetic interactions, if the disease involves interaction among three pathways. In short, missing heritability need not directly correspond to missing variants, because current estimates of total heritability may be significantly inflated by genetic interactions. Finally, we describe a method for estimating heritability from isolated populations that is not inflated by genetic interactions.

(frequency <1%) with large effects (3-9). We will discuss the frequency spectrum of disease-related variants in our second paper in this series.

Here we explore the possibility that a significant portion of the missing heritability might not reflect missing variants at all. The basic idea is easy to state: Current studies use estimators of  $h_{ad}^2$ that are not consistent (that is, converge to the wrong answer); they may seriously overestimate the denominator  $h_{ad}^2$  and thus underestimate  $\pi_{explained}$ . As a result, even when all variants affecting the trait are discovered,  $\pi_{explained}$  may fall far short of 100%. We refer to this gap as "phantom heritability."

Quantitative geneticists have long known that genetic interactions can affect heritability calculations (10). However, human genetic studies of missing heritability have paid little attention to the potential impact of genetic interactions. A few authors have constructed mathematical examples (11, 12), but these abstract models have not been related to biologically plausible mechanisms, and the studies have not considered whether the presence of genetic interactions would be readily detected, thereby preventing geneticists from being fooled by phantom heritability. The prevailing view among human geneticists appears to be that interactions play at most a minor part in explaining missing heritability.

Here we show that simple and plausible models can give rise to substantial phantom heritability. Biological processes often depend on the rate-limiting value among multiple inputs, such as the levels of components of a molecular complex required in stoichiometric ratios, reactants required in a biochemical pathway, or proteins required for transcription of a gene. We thus On top of that, the "Missing heritability" may be just non-existent. It may be "Phantom heritability" caused by ignoring epistasis, which inflates family based estimates of heritability

Estimated



A raging debate!!

### **IN SUMMARY:**

As to risk prediction:

(1) Missing heritability

(2) Hidden heritability

(3) Phantom heritability

(4)...

As to the causal variants:

(1) Common variants

(2) Rare variants

(3) Epigenetics

(4)...

#### So the information is there. What to do now?

(1) New information about variants

- (1) Larger sample sizes and families. Increasing pressure for a
  - **Global Alliance for Genomics and Health**

(2) New arrays and/or NGS sequencing(Human Omni5\_Quad and/or Hiseq2500)







### The current SOTA is either large samples or exomics

## Sequencing Evolution/Revolution





1990: thousand bases/day

























#### The current SOTA is either large samples or exomics

ARTICLES



# Exome sequencing identifies the cause of a mendelian disorder

Sarah B Ng<sup>1,10</sup>, Kati J Buckingham<sup>2,10</sup>, Choli Lee<sup>1</sup>, Abigail W Bigham<sup>2</sup>, Holly K Tabor<sup>2,3</sup>, Karin M Dent<sup>4</sup>, Chad D Huff<sup>5</sup>, Paul T Shannon<sup>6</sup>, Ethylin Wang Jabs<sup>7,8</sup>, Deborah A Nickerson<sup>1</sup>, Jay Shendure<sup>1</sup> & Michael J Bamshad<sup>1,2,9</sup>

We demonstrate the first successful application of exome sequencing to discover the gene for a rare mendelian disorder of unknown cause, Miller syndrome (MIM%263750). For four affected individuals in three independent kindreds, we captured and sequenced coding regions to a mean coverage of 40× and sufficient depth to call variants at ~97% of each targeted exome. Filtering against public SNP databases and eight HapMap exomes for genes with two previously unknown variants in each of the four individuals identified a single candidate gene, *DHODH*, which encodes a key enzyme in the pyrimidine *de novo* biosynthesis pathway. Sanger sequencing confirmed the presence of *DHODH* mutations in three additional families with Miller syndrome. Exome sequencing of a small number of unrelated affected individuals is a powerful, efficient strategy for identifying the genes underlying rare mendelian disorders and will likely transform the genetic analysis of monogenic traits.

Received 2 October; accepted 9 November; published online 13 November; corrected online 22 November 2009 (details online); doi:10.1038/ng.499



5 ۰

the genetic etiology of ASDs. ASDs are characterized by pervasive impairment in language, communication and social reciprocity and restricted interests or stereotyped behaviors<sup>1</sup>. Several new candidate loci for ASDs have recently been identified using genome-wide approaches that discover individually rare events of major effect2. A number of genetic syndromes with features of the ASD phenotype, collectively referred to as syndromic autism, have also been described<sup>4</sup>. Despite this progress, the genetic basis for the vast majority of cases remains unknown. Several observations support the hypothesis that the genetic basis for ASDs in sporadic cases may differ from that of families with multiple affected individuals, with the former being more likely to result from de novo mutation events rather than inherited variants<sup>1,5-7</sup>. In this events within coding sequence and three additional events mapping study, we sequenced the protein-coding regions of the genome (the to 3' untranslated regions (Table 2). A list of predicted variant sites exome)8 to test the hypothesis that de novo protein-altering muta- within these genes from the 1000 Genomes Pilot Project data15 is tions contribute substantially to the genetic basis of sporadic ASDs. provided for comparison (Supplementary Table 5).

20 probands, particularly among more severely affected

CNTNAP2 missense variant, and we provide functional

show that trio-based exome sequencing is a powerful

individuals, in FOXP1, GRIN2B, SCN1A and LAMC3. In the

FOXP1 mutation carrier, we also observed a rare inherited

support for a multi-hit model for disease risk3. Our results

approach for identifying new candidate genes for ASDs and

suggest that de novo mutations may contribute substantially to

identify a maternally inherited deletion (~350 kb) at 15g11.2 in one family (Supplementary Fig. 1). This deletion has been associated with increased risk for epilepsy10 and schizophrenia11,12 but has not been considered causal for autism.

Similar to researchers from a previous study13, who reported exome sequencing on ten parent-child trios with sporadic cases of moderate to severe intellectual disability, we performed exome sequencing on each of the 60 individuals separately by subjecting whole-blood derived genomic DNA to in-solution hybrid capture and Illumina sequencing (Online Methods). We obtained sufficient coverage to call variants for ~90% of the primary target (26.4 Mb) (Table 1). Genotype concordance with SNP microarray data was high (99.7%) (Supplementary Table 2), and on average, 96% of proband variant sites were also called in both parents (Supplementary Table 3). Given the expected rarity of true de novo events in the targeted exome (<1 per trio) (Supplementary Table 4)14, we reasoned that most apparently de novo variants would result from under calling in parents or systematic false positive calls in the proband. We therefore filtered variants previously observed in the dbSNP database, 1000 Genomes Pilot Project data15 and 1,490 other exomes sequenced at the University of Washington (Supplementary Fig. 2). We performed Sanger sequencing on the remaining de novo candidates (<5 per trio), validating 18

<sup>1</sup>Department of Genome Schnees, University of Washington School of MacLines Sauttle, Washington, USA, "Paulicines Trust Centre Neuran Genetics, Diseastive Oxford, DMI, UK, "Department of Psychiatry and Behavioral Sciences, University of Washington, Seattle, Washington, USA, "Language and Genetics Department, Max Flunck Institute for Psychiatry and Behavioral Sciences, University of Washington, USA, "Language and Genetics Department, Max Flunck Institute for Psychiatry and Retreministics," Howare Hughes Micclus Institute, Sauttle, Washington, USA, Correspondences should be an effective and the Sciences Sciences, University of Washington, Sauttle, Sauttle, Washington, USA, Correspondences should be and the Sciences Sciences Sciences Sciences Sciences, Sciences Sciences, Sci addressed to E.E.E. (eee@es.washington.edu) or J.S. (shendure@uw.edu).

Received 22 February; accepted 21 April; published online 15 May 2011; doi:10.1038/ng.835

NATURE GENETICS - VOLUME 43.1 NUMBER 6.1 JUNE 2011

585

identified four genes that are recurrently mutated; notch 1 (NOTCHI), exportin 1 (XPOI), myeloid differentiation primary response gene 88 (MYD88) and kelch-like 6 (KLHL6). Mutations in MYD88 and KLHL6 are predominant in cases of CLL with mutated immunoglobulin genes, whereas NOTCH1 and XPO1 mutations are mainly detected in patients with unmutated immunoglobulins. The patterns of somatic mutation, supported by functional and clinical analyses, strongly indicate that the recurrent NOTCH1, MYD88 and XPO1 mutations are oncogenic changes that contribute to the clinical evolution of the disease. To our knowledge, this is the first comprehensive analysis of CLL combining whole-genome sequencing with clinical characteristics and clinical outcomes. It highlights the usefulness of this approach for the identification of clinically relevant mutations in cancer.

To gain insights into the molecular alterations that cause CLL, we performed whole-genome sequencing of four cases representative of different forms of the disease: two cases, CLL1 and CLL2, with no mutations in the immunoglobulin genes (IGHV-unmutated) and two cases, CLL3 and CLL4, with mutations in these genes (IGHV-mutated) (Supplementary Table 1 and Supplementary Information). We used a according to their potential functional effect (Supplementary Informacombination of whole-genome sequencing and exome sequencing, as tion). We also searched for small insertions and deletions (indels) in well as long-insert paired-end libraries, to detect variants in chromo- coding regions: we found and validated five somatic indels, which somal structure (Supplementary Fig. 1 and Supplementary Tables 2-5). caused frameshifts in protein-coding regions (Supplementary Table 7).

in a CpG context (Fig. 1b and Supplementary Fig. 2). We also detected marked differences in the mutation pattern between CLL samples and these differences were associated with tumour subtype (Fig. 1b). Thus, IGHV-mutated cases showed a higher proportion of A>C/T>G muta tions than cases with unmutated IGHV (16  $\pm$  0.2% versus 6.2  $\pm$  0.1%). The base preceding the adenine in A to C transversions showed an over representation of thymine, when compared to the prevalence expected from its representation in non-repetitive sequences in the wild-type genome (P < 0.001, Fig. 1c), and there were fewer A to C substitution at GpA dinucleotides than would be expected by chance ( $P \le 0.001$ ). These differences between CLL subtypes might reflect the molecular mechanisms implicated in their respective development. The pattern and context of mutations are consistent with their being introduced by the error-prone polymerase η during somatic hypermutation in immunoglobulin genes". This indicates that polymerase  $\boldsymbol{\eta}$  could contribute to the high frequency of A > T to C > G transversions in cases with IGHV-mutated. It also extends the differences observed between these two CLL subtypes to the genomic level.

We classified the somatic mutations into three different classes

Departamento de Bioquímica y Biología Molecular, Instituto Universidario de Oncología, Universidad de Oviedo, 33006 Oviedo, Spein.<sup>9</sup>Unidad de Genómica, Institut d'Investigacions Biomieliques Augus Al Surger (1944): SOLIA Brenzes, Span, Tradical Perentgatings, Encode Antonio Henrice, Technologia Concellander, Solia Brenzes, Span, Tradical Perentgatings, Encode Antonio Henrice, State Concellander, S nemora presente nel la cue une cue conserve sub entre entre entre el semana conserve el se senago conserve sub entre el senare el sena El senare el se El senare el se

02011 Macmillan Publishers Limited. All rights reserved

Saan 11/Centro Nacional de Analies Gentinico, Parc Centric de Berceiron 08028 Barcelona, Spain 11/Centro Nacional de Britecinología, Consejo Superior de Freedigaciones Centralios, 25047 Madrid Spain <sup>19</sup>Wellcome Trust Sanger Institute, Hinkton CB10 1SA UK These authors contributed equally to this work

7 JULY 2011 | VOL 475 | NATURE | 10

#### The current SOTA is either large samples or exomics...or other designs

### Or **SOME** ingenious comparisons



#### Vol 464 29 April 2010 doi:10.1038/nature08990

#### nature

#### LETTERS

### Genome, epigenome and RNA sequences of monozygotic twins discordant for multiple sclerosis

Sergio E. Baranzini<sup>1</sup>, Joann Mudge<sup>2</sup>, Jennifer C. van Velkinburgh<sup>2</sup>, Pouya Khankhanian<sup>1</sup>, Irina Khrebtukova<sup>3</sup>, Neil A. Miller<sup>2</sup>, Lu Zhang<sup>3</sup>, Andrew D. Farmer<sup>2</sup>, Callum J. Bell<sup>2</sup>, Ryan W. Kim<sup>2</sup>, Gregory D. May<sup>2</sup>, Jimmy E. Woodward<sup>2</sup>, Stacy J. Caillier<sup>1</sup>, Joseph P. McElroy<sup>1</sup>, Refujia Gomez<sup>1</sup>, Marcelo J. Pando<sup>4</sup>, Leonda E. Clendenen<sup>2</sup>, Elena E. Ganusova<sup>2</sup>, Faye D. Schilkey<sup>2</sup>, Thiruvarangan Ramaraj<sup>2</sup>, Omar A. Khan<sup>5</sup>, Jim J. Huntley<sup>3</sup>, Shujun Luo<sup>3</sup>, Pui-yan Kwok<sup>6,7</sup>, Thomas D. Wu<sup>8</sup>, Gary P. Schroth<sup>3</sup>, Jorge R. Oksenberg<sup>1,7</sup>, Stephen L. Hauser<sup>1,7</sup> & Stephen F. Kingsmore<sup>2</sup>

#### Table 1 | SNP and indel genotypes and differences between siblings in three twin pairs

			Twin pair	041896			Twin pair 23017	8		Twin pair 041907			
Genotype change and individual	Platform	SNP genotypes	Replicated SNP genotype differences§	Indel genotypes	Replicated indel genotype difference	SNP genotypes	Replicated SNF genotype difference	P Indel genotypes	SNP genotypes	Replicated SNP genotype difference	Indel genotypes		
No change	Genome-Seq* SNP array (×2)	1,086,309 736,782	79,209 1,638 (98.3%)	26,908 NA	91 (91.9%) NA 91 (91.9%)	ND 783,189	NA 888 (95.3%)	ND NA 1.034	ND 796,870	NA 385 (98.0%)	ND NA 297		
Ref in -001 $\rightarrow$ bet in -101	Genome-Seq*†	202	0	1,514 3 NA	Ο ΝΔ	ND 36	NA	ND NA	ND 32	NA	ND NA		
101	mRNA-Seq <sup>†</sup> ‡	12	0	0	0	6	0	0	2	0	0		
Het in -001 → ref in -101	Genome-Seq*† SNP array (×2)	134 49	0 0	1 NA	0 NA	ND 31	NA 0	ND NA	ND 11	NA 0	ND NA		
Het in -001 $\rightarrow$	mRNA-Seq <sup>†</sup> ‡ Genome-Sea*†	5 1 513	0	0 128	0	9 ND	0 NA	0 ND	16 ND	0 NA	0 ND		
hom in -101	SNP array (×2)	29	0	NA	NA	24	0	NA	17	0	NA		
Hom in -001 $\rightarrow$	Genome-Seq*†	203 1,392	0	/ 81	0	ND	NA	ND	ND	NA.	5 ND		
het in -101	SNP array ( $\times$ 2)	16 102	0	NA 1	NA	62 429	0 Table 2	CpG site	s and clus	ters in mono	rygotic twins, normal	and cancer s	an

	mkinA-sed 1	102	0	T	0	429 0	Genomic DNA sample	CpG	CpG	Ratio of	CpGs	CpG	mCpG	Between sample	CpGs	CpG	mCpG
Genotype catego * Nucleotide geno	ries: homozygous re otyped if 11–44× co	ference (ref), ⊦ verage and Q ≧	neterozygous varian ≥ 20.	nt (het) and	homozygous varian	t (hom). NA, not app		sites*	clusters	CpGs to clusters	shared	clusters shared	unique to one sample†	comparison†	shared	clusters shared	unique to one sample†
† Genotypes dete ‡ Genotyped if pr & Detected by pla	rmined according to esent in >2 reads, > tform on correspond	frequency cut >1 uniquely alig	offs in Supplementa gning read and Q ≥ rated by platform lif	ary Table 8 20. sted on row	and differences call	ed if frequencies diffe	041896-001 T cell 041896-101 T cell	2,146,620	1,230,241 1,190,741	1.74 1.71	98.1%	98.2%	2‡ 0	041896- & 230178-001 T cell	97.4%	97.7%	522 305
s Detected by pla	control contespond	aing row, replic	ated by platform is	sted on row	v Delow.		230178-001 T cell 230178-101 T cell	1,636,285 1,917,131	1,038,787 1,155,024	1.58 1.66	97.8%	97.9%	3 7	041896-001 & 230178-101 T cell	96.5%	96.9%	445 362
							041907-001 T cell 041907-101 T cell	1,779,140 1.642,200	1,094,361 1.038.090	1.63 1.58	90.6%	92.7%	174 2	041896- & 041907-001 T cell	97.5%	98.1%	304 282
							Normal breast Breast cancer	1,829,855 2.010,173	1,086,405 1,192,180	1.68 1.69	96.7%	97.9%	696 861	041896-001 T cell & normal breast	97.3%	98.0%	5,620 1,560
							Normal lung Lung cancer	2,096,524 1,619,178	1,216,046 956,760	1.72 1.69	97.9%	98.8%	6,891 9,618	041896-001 T cell & normal lung	96.1%	97.0%	3,329 926

CpG sites and clusters were compared between CD4<sup>+</sup> lymphocytes from three pairs of monozygotic twins, breast and lung cancer and normal tissue samples.

\* >10 RRBS reads aligned by ELAND-extended and Q > 20.

+ CpG >80% methylated in one sample and <20% in other.

‡ Not replicated after RRBS read alignment with GSNAP.

### And this is only the scientific part of the problem

### Personalized Medicine: Vision vs. Re

- Disruptive developments in science and technology
- Convergence of molecular biology, genetics, advanced technology, bioinformatics, broadband
- "Team science"
- Transformational changes in medicine
  - Molecular-based products and services
  - Shift towards prevention
  - Reclassification of disease
  - Integration and coordination
  - IT solutions; Interoperability
  - Consumer-centered
  - Premise that knowledge will change behavior
- Huge public & private investments in R&D
- Health as a national asset
- Ethical, legal and policy issues addressed in parallel with the science



- Healthcare delivery focused on "sick care"
  - Standardization for quality improvement
- Fragmented, lack of coordination
- Costs growing and unsustainable
  - Pressures of expensive new technologies
  - Aging population in search of new services
  - Millions of Americans under- or uninsured
  - Employer-based system tenuous
  - No evidence of healthier citizenry
- Inefficient use of information
  - Lack of IT investment, connectivity
- Evidence base for medicine inadequate
  - Continuing debate about role of cost-effectiveness
- Huge provider knowledge gaps re genomics
- Complicated regulatory framework
- Reimbursement hurdles and uncertainties
- Powerful stakeholders in current system resist change

#### **Only starting**



#### **Only starting**







Biomarker	Application
Her-2/neu receptor	Select Herceptin (trastuzumab) for breast cancer
BRCA1/2	Breast and ovarian cancer inherited risk, prophylactic tamoxifen and surgery
Transcriptional profile – 21 genes	Avoid use of chemotherapy in breast CA patients with low risk of recurrence
CYP2D6/CYP2D19	Guide prescribing/ adjust dose of ~25% of commonly used drugs
VKOR/CYP2C9	Dosing of warfarin

Table 1 Examples of genetic and genomic testing in personalized medicine Pre-symptomatic risk assessment BRCA1/2 testing for breast cancer<sup>a</sup> Lynch syndrome testing for hereditary colon cancer<sup>b</sup> Long OT intervalc,d Spinal Muscular Atrophye Diagnosis Beta thalassemiaf Fusion genes and rearrangements including BCR-ABL, E2A-PBX1, TEL-AML1, and MLL in pediatric leukemiage Gene expression profiles define subtypes of breast cancer<sup>h</sup> Human Papilloma Virus detectioni Hepatitis C detection<sup>1</sup> PCR detection of micro-organisms (bacteria, fungi)k Prognosis Fragile X syndrome (number of trinucleotide repeats predicts severity)1 Gene expression signatures and prognosis in breast cancer<sup>m</sup> Gene expression analysis and lymphoma prognosis<sup>n</sup> Treatment and pharmacogenomics Therapies for targeted gene mutations in cancer<sup>o</sup> EGFR point mutations in lung cancer and glioblastoma and cetuximab, gefitinib, erlotinib, panitumumab, lapatinib treatment KIT, PDGFR mutations in sarcoma, glioma, liver and renal cancer, melanoma and imatinib, nilotinib, sunitinib, sorafenib treatment BRAF mutations in melanoma treated by RAF inhibitors BCR-ABL translocation in chronic myelogenous leukemia treated by imatinib KRAS wild-type status correlated with resistance to EGFR inhibition PARP inhibitors in BRCA mutant breast, ovarian, prostate and pancreatic cancer Herceptin (Trastuzumab) in HER2 + breast cancer Pharmacogenomic applications<sup>p</sup> CYP 2C19\*2 variant (rs4244285) associated with diminished clopidogrel response<sup>q</sup> Rs2395029 testing for HLA-B\*5701 allele, correlated with hypersensitivity to abacavir treatment for HIV+ patients<sup>r</sup>

<sup>a</sup> Robson and Offit (2007), <sup>b</sup> EGAPP (2009a), <sup>c</sup> Napolitano et al. (2005), <sup>d</sup> Lehnart et al. (2007), <sup>e</sup> Prior et al. (2008), <sup>f</sup> Galanello and Origa (2010), <sup>g</sup> Carroll et al. (2003), <sup>h</sup> Sorlie et al. (2001), <sup>i</sup> Nicol et al. (2010), <sup>j</sup> Pham et al. (2010), <sup>k</sup> Tsalik et al. (2010), <sup>1</sup> Sherman et al. (2005), <sup>m</sup> Kim and Paik (2010), <sup>n</sup> Rosenwald et al. (2002), <sup>o</sup> Macconaill and Garraway (2010), <sup>p</sup> U.S. Food and Drug Administration (2011), <sup>q</sup> Shuldiner et al. (2009), <sup>r</sup> Colombo et al. (2008)



© Francis Collins, 2008

Examples of going from individual genome-wide analyses to treatment are accumulating (e.g., Nic Volker, the Beery twins, Mike Snyder, John Lauerman).

Editorial, Molecular Systems Biology 9. January 2013

BRIEF REPORT



Nic Volker and the XIAP gene (Worthey et all. Genetics In Medicine 2011)



Science Transl Med 87 June 2011



### Take home message:



## Lots to do that we MUST be doing!!!