

Use of Race, Ethnicity, and Ancestry as Population Descriptors in Genomics Research

Prof. Melinda Mills, University of Oxford
Director, Leverhulme Centre for Demographic Science

The National Academies of Sciences, Engineering,
Medicine, 04 April 2022

REVIEW ARTICLE

<https://doi.org/10.1038/s42003-018-0261-x>

OPEN

A scientometric review of genome-wide association studies

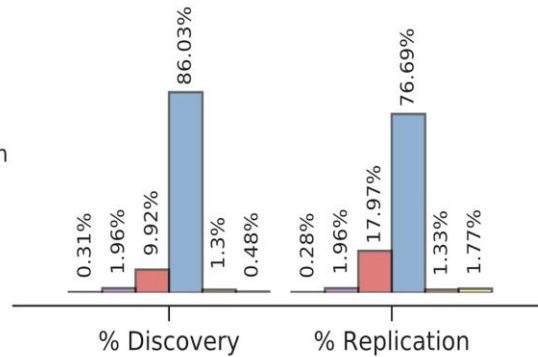
Melinda C. Mills ¹ & Charles Rahal ¹

This scientometric review of genome-wide association studies (GWAS) from 2005 to 2018 (3639 studies; 3508 traits) reveals extraordinary increases in sample sizes, rates of discovery and traits studied. A longitudinal examination shows fluctuating ancestral diversity, still predominantly European Ancestry (88% in 2017) with 72% of discoveries from participants recruited from three countries (US, UK, Iceland). US agencies, primarily NIH, fund 85% and women are less often senior authors. We generate a unique GWAS H-Index and reveal a tight social network of prominent authors and frequently used data sets. We conclude with 10 evidence-based policy recommendations for scientists, research bodies, funders, and editors.

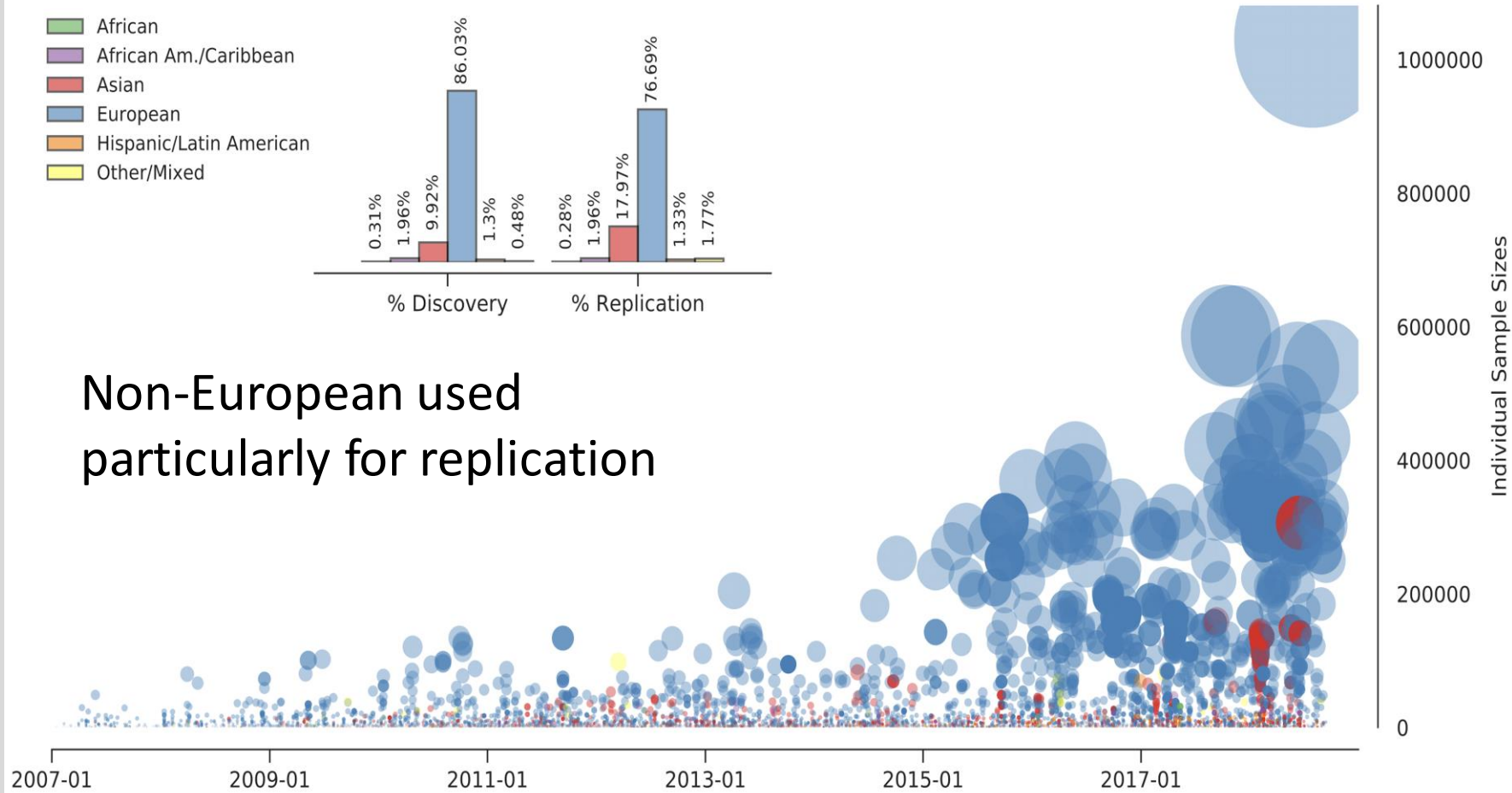


80-90% of genetic discovery European ancestry

- African
- African Am./Caribbean
- Asian
- European
- Hispanic/Latin American
- Other/Mixed



Non-European used particularly for replication



Defining race, ethnicity, ancestry

Race/ethnicity socially constructed not biological category – different from ancestry in genetics

[J Health Soc Behav](#). 2021 Jun 8;22:1465211018682. doi: 10.1177/00221465211018682. Online ahead of print.

Reconstructing Sociogenomics Research: Dismantling Biological Race and Genetic Essentialism Narratives

Pamela Herd ¹, Melinda C Mills ², Jennifer Beam Dowd ²

Affiliations [+](#) expand

PMID: 34100668 DOI: [10.1177/00221465211018682](#)

Abstract

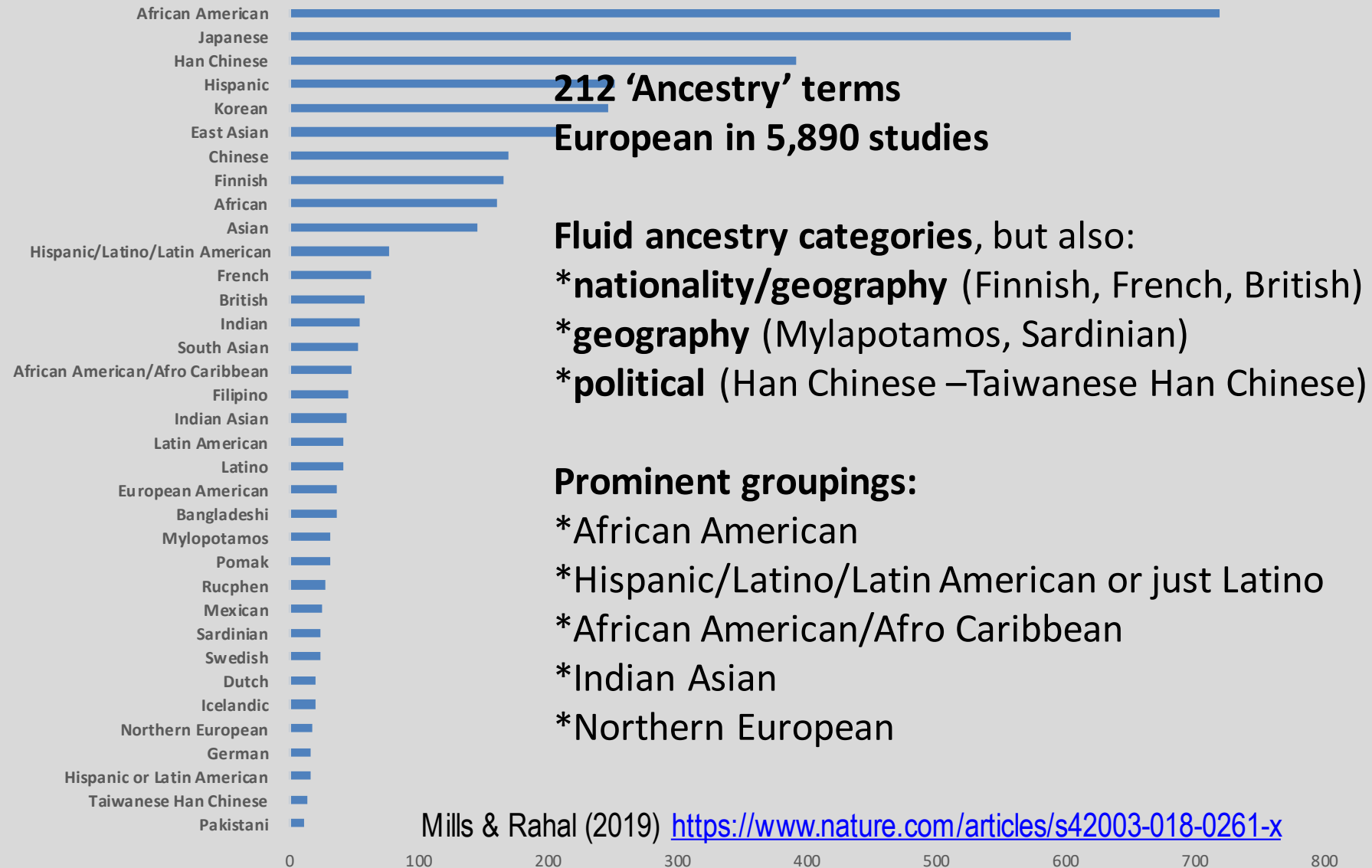
We detail the implications of sociogenomics for social determinants research. We focus on education and race because of how early twentieth-century scientific eugenic thinking facilitated a range of racist and eugenic policies, most of which helped justify and pattern racial and educational morbidity and mortality disparities that remain today, and are central to sociological research. Consequently, we detail the implications of sociogenomics research by unpacking key controversies and opportunities in sociogenomics as they pertain to the understanding of racial and educational inequalities. We clarify why race is not a valid biological or genetic construct, the ways that environments powerfully shape genetic influence, and risks linked to this field of research. We argue that sociologists can usefully engage in genetics research, a domain dominated by psychologists and behaviorists who, given their focus on individuals, have mostly not examined the role of history and social structure in shaping genetic influence.

Keywords: education; gene-environment interactions; race; sociogenomics.

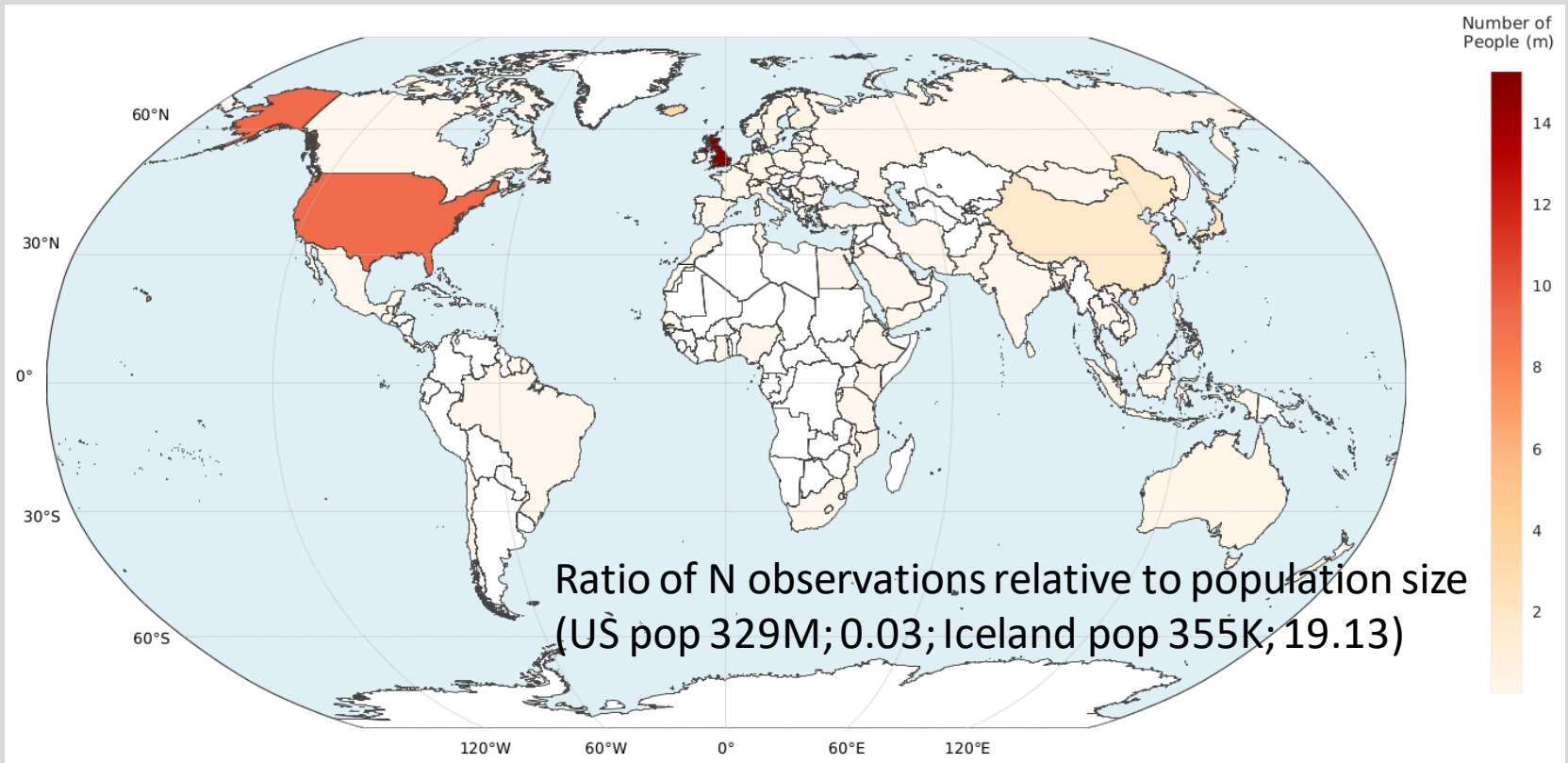


Top terms used in GWAS

Below terms after European (N=5,890 studies) removed, Mills & Rahal 2019



Country of Recruitment – 72% come from 3 countries (UK, US, Iceland)



NHLBI	30127	9950	7043	5218	2144	781	865	9849	486	2301	1997	225	3249	3220	3157	1102	1422	1288	142	1781	125	692	345	27
NCI	9638	1857	2733	729	838	814	383	1786	256	456	166	4343	199	140	45	567	285	280	1179	41	81	107	11	17
NIA NIH HHS	8577	1922	1369	548	1849	711	122	3189	2711	361	130	27	265	562	396	84	400	531	140	225	33	68	69	20
MRC	9975	1186	2633	648	1018	321	378	3327	794	1102	368	184	278	230	199	98	403	249	377	206	319	189	37	70
NIDDK NIH HHS	7388	1627	2069	1163	417	285	211	2279	48	980	335	5	927	208	181	186	243	223	852	179	220	520	42	20
NIMH	4444	641	461	340	625	453	59	943	2158	198	56	80	39	20	20	133	273	211	13	51	95	58	6	2
Wellcome Trust	4687	346	1289	182	355	168	177	1171	314	368	321	109	136	206	46	58	123	156	163	127	218	378	9	10
NHGRI	3791	989	671	359	338	199	229	914	308	304	225	100	279	252	211	97	186	158	86	186	50	129	58	2
NCRR NIH HHS	2935	773	438	468	353	229	100	948	160	206	409	40	237	269	196	242	71	156	138	126	113	104	36	12
NIAID NIH HHS	2492	1487	47	1674	859	312	300	1082	16	8	1693	3	79	5	9	1386	4	13	80	82	199	1	87	0
	European	African Am./Caribbean	Asian	Hispanic/Latin American	In Part Not Recorded	Other/Mixed	African	Other Meas.	Neurological	Body Meas.	Other Disease	Cancer	Lipid/Lipoprotein Meas.	Cardiovascular Disease	Cardiovascular Meas.	Response To Drug	Biological Process	Other Trait	Digestive System	Hematological Meas.	Immune System	Metabolic	Inflammatory Meas.	Liver Enzyme Meas.
	Ancestry							Broad EFO Category																

The GWAS Diversity Monitor tracks diversity by disease in real time

To the Editor — The Genome-wide association study (GWAS) is a primary tool for the discovery of associations between genetic variants and complex phenotypes, cataloged by the National Human Genome Research Institute–European Bioinformatics Institute (NHGRI-EBI) GWAS Catalog, which currently contains information on more than 4,346 published studies across more than 4,933 diseases and traits. Although there has been a considerable

2019. Our cumulative estimate at the time of writing currently stands at 88.45%, even despite the recent launch of initiatives such as H3Africa, the African Genome Variation Project and GenomeAsia 100k. Geographic and demographic diversity is also limited, and other estimates suggest that 72% of participants are recruited from just three countries (the United States, the United Kingdom and Iceland)¹.

The transferability of GWAS results

age or sex, and socioeconomic status of individuals⁶. With the move toward the use of PRSs derived from GWAS for clinical applications⁷, most PRSs derived from GWAS would exacerbate existing global health inequalities⁸. Substantially more genetic variation exists in non-European populations, and this variation can provide a rich resource for finding new genetic associations (Supplementary Note 1).

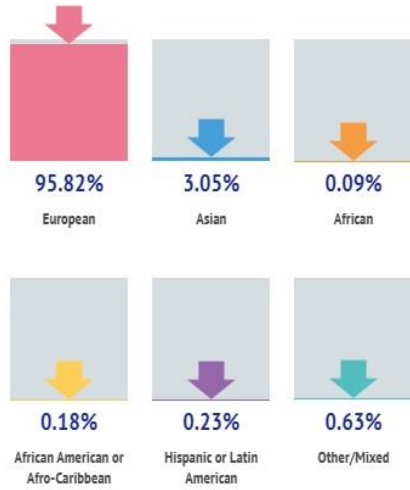
GWAS often fail to identify variants that

Mills, M.C. & C. Rahal (2020), *Nature Genetics*, <https://rdcu.be/b2BX4>

gwasdiversitymonitor.com

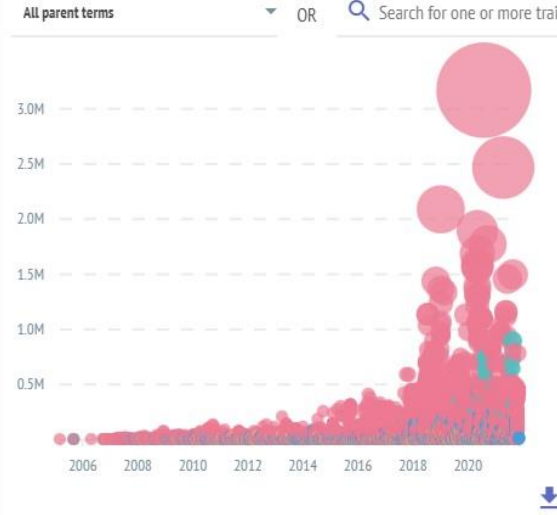
Total GWAS participants diversity

Version 1.0.0. Last check for data: 2022-01-18 09:54:08.



Ancestry over time by parent term

Discovery Stage

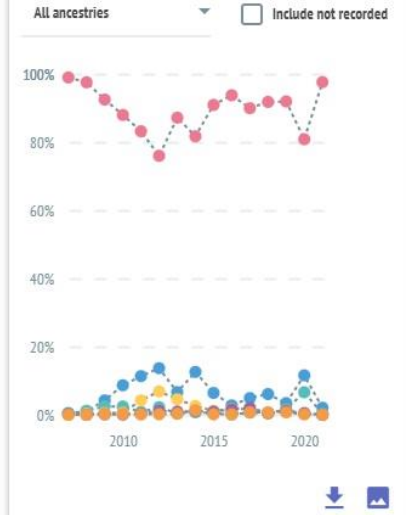


VIEW ALL

- European
- Asian
- African
- African American or Afro-Caribbean
- Hispanic or Latin American
- Other/Mixed

Participants across all parent terms

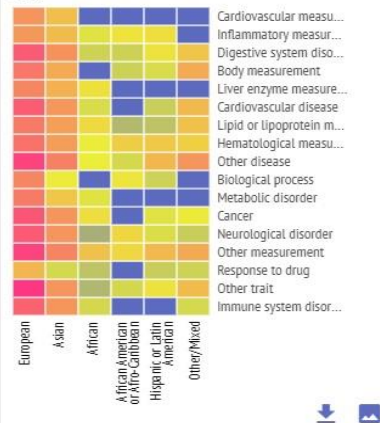
Discovery Stage



Parent term by 'broader' ancestry

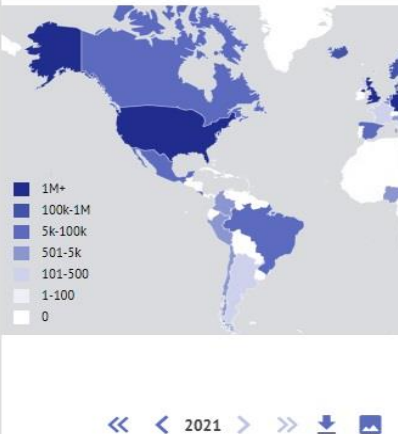
By Participants - Discovery Stage

<< < 2021 > >>



Participants by country (all parent terms)

Both Stages

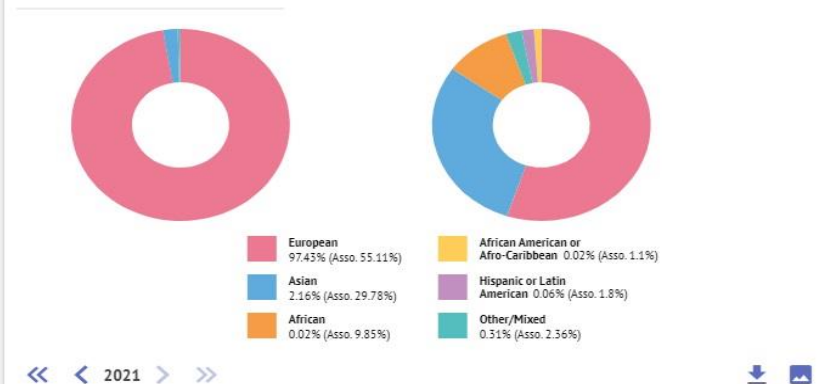


Participants by ancestry

Discovery Stage

Click to hide associations discovered

All parent terms



Count of all associations discovered

Discovery Stage

>>

Hide

Real-time search 5500+ phenotypes by ancestry

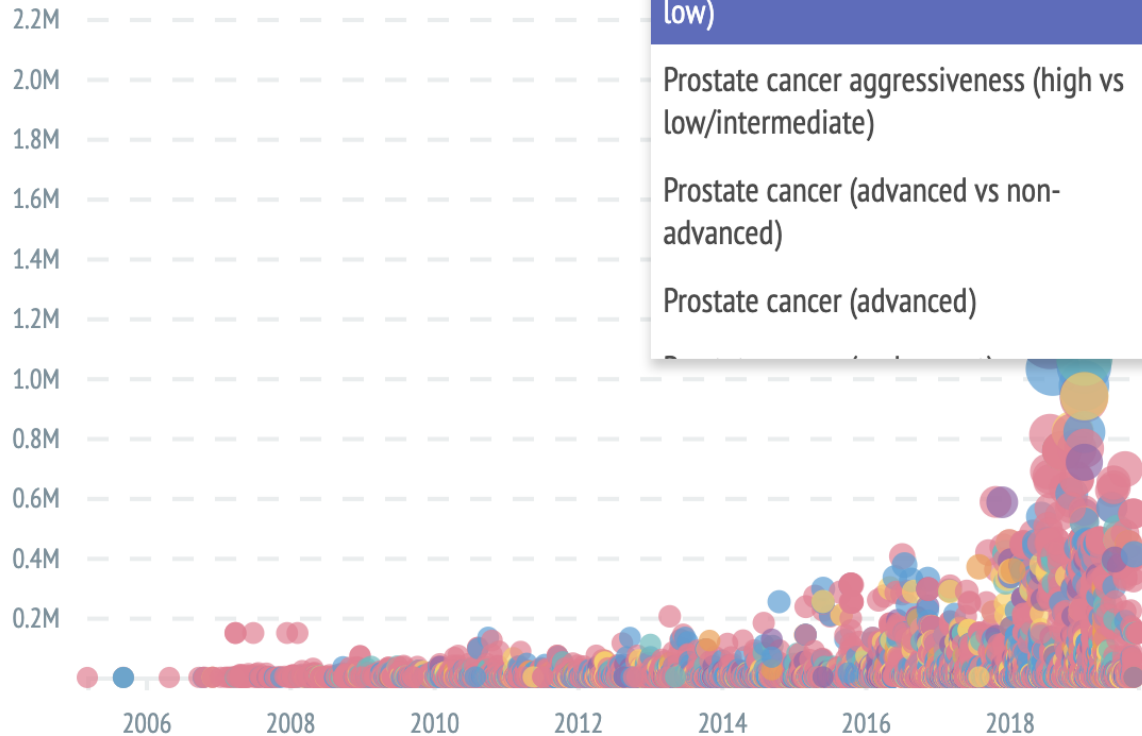
Ancestry over time by parent term

Discovery Stage

All parent terms

OR

Prostate



Prostate cancer aggressiveness (high vs low)

Prostate cancer aggressiveness (high vs low/intermediate)

Prostate cancer (advanced vs non-advanced)

Prostate cancer (advanced)

VIEW ALL

European

African

African American or Afro-Caribbean

Other/Mixed

Asian

Hispanic or Latin American

Mills, M.C. & C. Rahal (2020) The GWAS Diversity Monitor tracks diversity by disease in real time, *Nature Genetics*

Top 10 datasets

Cohorts	
	1. European ancestry
Rotterdam	2. Mostly industrialized countries (NL, US, UK, Germany) – share similar disease prevalence and population profiles
Cooperativa de Aseguramiento	3. Older populations
Framingham	4. Sex ratio imbalance, more women
Atherosclerosis Risk in Communities	5. Non-representative samples (UKBB)
Cardiovascular Health Study	<ul style="list-style-type: none"> genetic associations are modifiable, bias GWAS estimates towards over-represented group; association observed in one study dependent on exposure-outcome relationship in discovery & target population, Keyes & Westreich 2019)
British 1947	
UK Adult Cohort	
European Cancer Cohort	
Nurses Health Study	
Study of Women's Health	

Coded largest 1,250 largest GWAS as of August 2018 to generate list most used datasets

Full list of 2,000+ data used:

https://github.com/crahal/GWASReview/blob/master/tables/Manually_Curated_Cohorts.csv

Mills & Rahal (2019) <https://www.nature.com/articles/s42003-018-0261-x>

Hidden heritability due to heterogeneity across seven populations

Felix C. Tropf^{1*}, S. Hong Lee², Renske M. Verweij³, Gert Stulp³, Peter J. van der Most⁴, Ronald de Vlaming^{5,6}, Andrew Bakshi⁷, Daniel A. Briley⁸, Charles Rahal¹, Robert Hellpap¹, Anastasia N. Iliadou⁹, Tõnu Esko¹⁰, Andres Metspalu¹⁰, Sarah E. Medland¹¹, Nicholas G. Martin¹¹, Nicola Barban¹, Harold Snieder⁴, Matthew R. Robinson^{7,12} and Melinda C. Mills¹

Meta-analyses of genome-wide heritability estimates obtained from seven sampling populations and time periods. We show that the heritability estimates obtained for education, 40% for age at first birth, and 40% for height are more likely to reflect heterogeneity across populations. These findings have substantial implications for the interpretation of behavioural phenotypes and the



RESEARCH ARTICLE



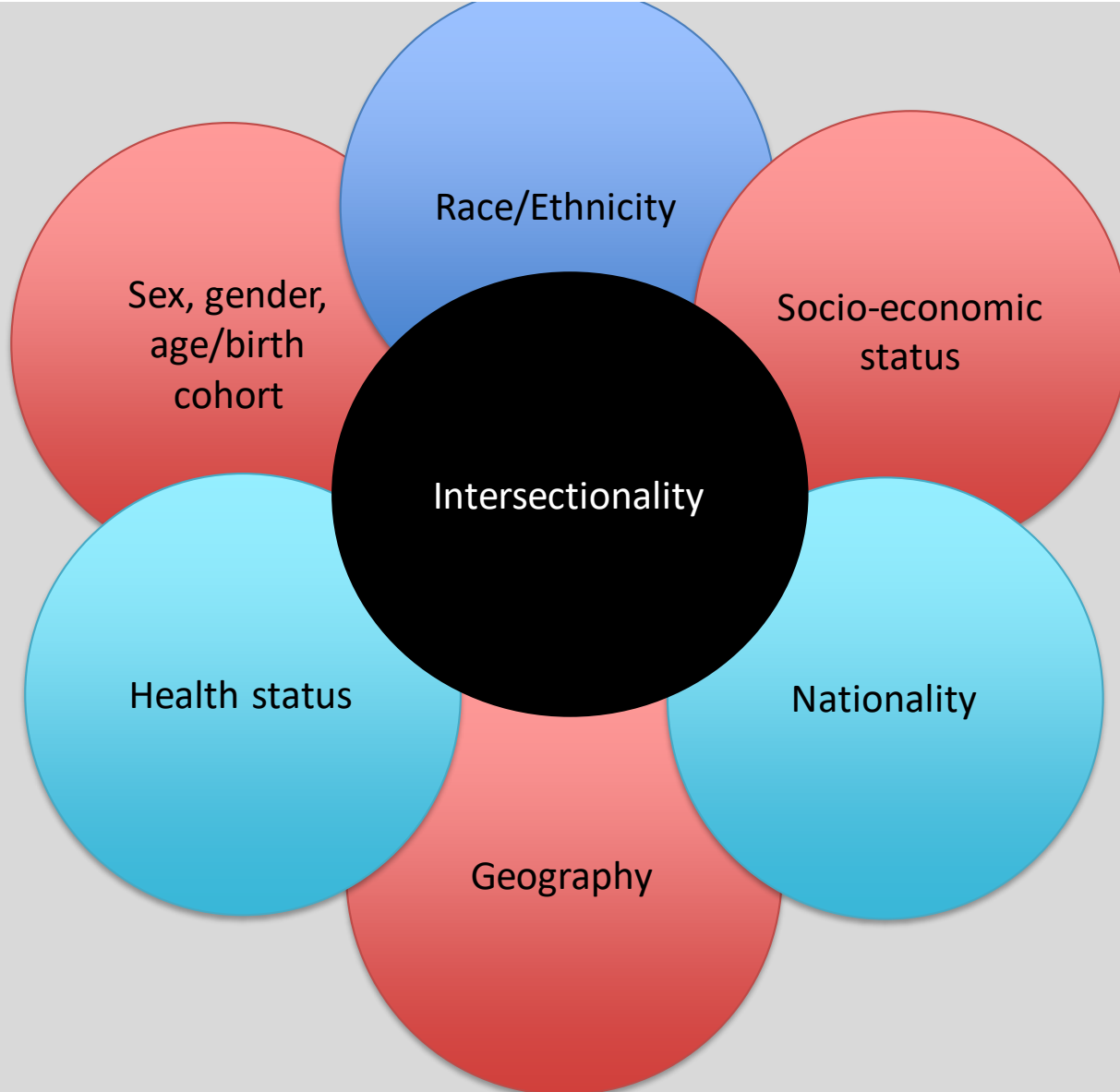
Variable prediction accuracy of polygenic scores within an ancestry group

Hakhamanesh Mostafavi^{1†*}, Arbel Harpak^{1†*}, Ipsita Agarwal¹, Dalton Conley^{2,3}, Jonathan K Pritchard^{4,5,6}, Molly Przeworski^{1,7*}

¹Department of Biological Sciences, Columbia University, New York, United States; ²Department of Sociology, Princeton University, Princeton, United States; ³Office of Population Research, Princeton University, Princeton, United States; ⁴Department of Genetics, Stanford University, Stanford, United States; ⁵Department of Biology, Stanford University, Stanford, United States; ⁶Howard Hughes Medical Institute, Stanford University, Stanford, United States; ⁷Department of Systems Biology, Columbia University, New York, United States

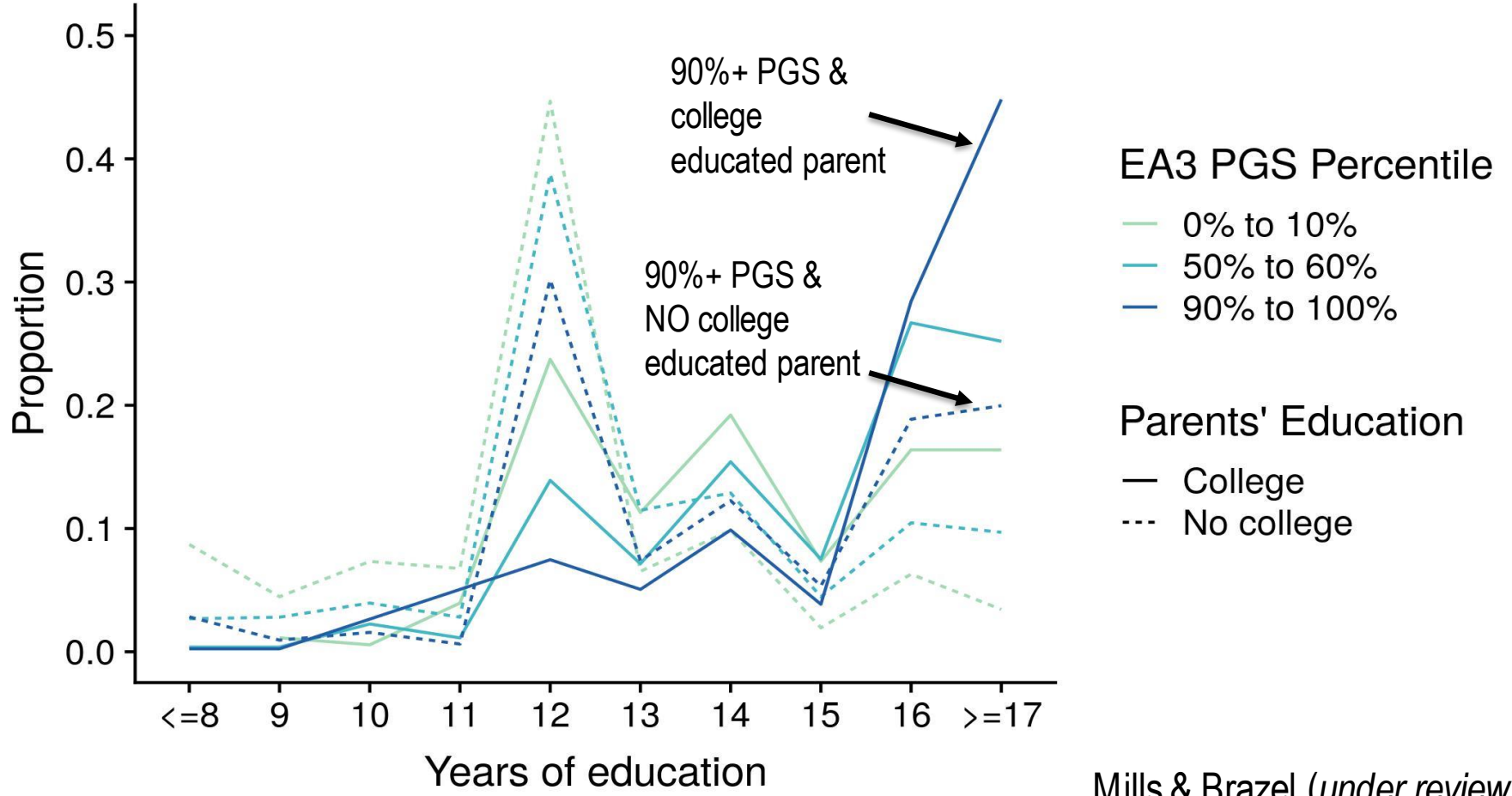
Abstract Fields as diverse as human genetics and sociology are increasingly using polygenic scores based on genome-wide association studies (GWAS) for phenotypic prediction. However,

Cumulative disadvantage & intersectionality



Highest PGS group (90-100%) highest education levels, but.....

GxE: children with at least **one college-educated parent** have **substantially higher educational attainment** in same PGS range than those without a college-educated parent



Future recommendations

Mills & Rahal (2019)

1. prioritize the inclusion of **multiple types of diversity** (socioeconomic status, sex)
2. **careful interpretation** of genetic differences
3. **participant and researcher** involvement
4. reduce **inequalities in authorship** and investigators
5. **reform incentive structures** role of authorship, data ownership, and data sharing
6. **coordinated governance, guidance** from multiple stakeholders
7. **monitoring** with funding consequences
8. **utilize influence** for the good of more people