

## Supplementary information to: Genomics is failing on diversity (Comment in *Nature* 538, 161–164; 2016)

Alice B. Popejoy & Stephanie M. Fullerton

### Introduction

Biomedical research is intertwined with genomics, as investigations of genetic causes for disease unravel the biological mechanisms responsible for wellness and illness. The genome-wide association study (GWAS) has been the primary tool for this discovery phase, and despite its limitation to detecting primarily common variation linked to common diseases, it has been very powerful for elucidating previously opaque biological mechanisms. Thousands of significant associations between genetic variants and phenotypes, or biological traits, have been found through GWAS that have contributed to foundational knowledge in the genetic architecture of disease pathways. Unfortunately, GWAS have been performed using samples from individuals of primarily European descent, and as such have missed a substantial portion of the genetic variation that is present in the human population. In order to develop a comprehensive understanding of the ways that genes and environments interact to determine risk for disease and impact drug responses or adverse drug reactions, it will be necessary to study individuals from a wide variety of social and ancestral backgrounds, living in diverse settings and communities.

### A. Querying the GWAS Catalog

To query the GWAS Catalog sample descriptions, A.B. Popejoy wrote a Python script to parse the list and description of studies that can be downloaded from the GWAS Catalog website: <https://www.ebi.ac.uk/gwas/docs/downloads>. The parser identifies the number of independent publications in the Catalog, and reports the distribution of ancestral populations sampled across all studies. It also produces lists of the number of individuals from each unique study population, and can be tailored to produce lists of PubMed ID numbers of studies with different ancestral populations. The program and documentation are available upon request, but may require manual updates as additional populations are added to the Catalog, as it was designed to parse particular phrases and patterns of phrases included in the current sample descriptions.

### B. Sample Descriptions and Ancestry Categories

As sample descriptions in the Catalog are not standardized, decisions needed to be made about how to categorize each sample with respect to broader ancestry bins for the aggregate analysis. These decisions were based on ancestry of origin and/or geographic proximity. For example, Afrikaners from South Africa are included in the “European ancestry” bin because of their European (Dutch) ancestral origins, which is discordant from their current residence on the continent of Africa. Singapore and Malaysia are included in the “Asian ancestry” bin while the Philippines and other island nations are included in the “Pacific Islander” bin based on the geographic connection to or separation from the continental landmass of Asia. A complete list of sample descriptions from the GWAS Catalog and their ancestral bins is provided below.

<b>“African ancestry”</b>	<b>“Arab/Middle Eastern”</b>	Dai Chinese ancestry	Singaporean ancestry
African American	Afghanistan ancestry	Dravidian ancestry	Singaporean Chinese ancestry
African American ancestry	Arab ancestry	East Asian ancestry	South Asian ancestry
African ancestry	Arab-Israeli ancestry	East Asian cases	South East Asian ancestry
African cases	Arab-Israeli founder	Han Chinese ancestry	South Indian
African Caribbean ancestry	Iranian ancestry	Han Chinese cases	Southern Indian ancestry
Afro-Caribbean	Jewish Israeli ancestry	Han Chinese controls	Sri Lankan Sinhalese ancestry
Afro-Caribbean cases	Jewish Israeli controls	Han Chinese individuals	Taiwanese ancestry
Afro-Caribbean controls	Jewish-Israeli ancestry	Hong Kong Chinese ancestry	Taiwanese
Afro-Caribbean individuals	Lebanese ancestry	Hui Chinese ancestry	Thai ancestry
Black cases	Israeli/Arab ancestry	Indian ancestry	Thai-Chinese ancestry
Black child cases	Israeli/Arab controls	Indian Asian ancestry	Tibetan ancestry
Black controls	Israeli/Arab cases	Japanese ancestry cases	Uighur cases
Black individuals	Middle Eastern ancestry	Japanese ancestry	Uighur controls
Gambian ancestry	Middle Eastern Arab ancestry	Japanese controls	Uygur Chinese ancestry
Malawian ancestry	Pakistani ancestry	Japanese ancestry	Uygur-Kazakh Chinese
Moroccan ancestry	Saudi Arabian ancestry	Jingpo Chinese ancestry	Vietnamese ancestry
Nigerian ancestry	Turkish ancestry cases	Korean ancestry	Vietnamese-Korean ancestry
North African ancestry	Turkish ancestry	Malay ancestry	
Seychelles female individuals	Turkish cases	Malaysian ancestry	<b>“European ancestry”</b>
Seychelles male individuals	Turkish controls	Malaysian Chinese ancestry	Amish cases
Seychellois ancestry	Turkish uveitis cases	Mongolian ancestry	Amish controls
Sub-Saharan African		Nepalese ancestry	Amish individuals
Tanzanian ancestry	<b>“Asian ancestry”</b>	North Indian ancestry	Ashkenazi Jewish cases
Tunisian ancestry	Asian ancestry	Oriental ancestry	Ashkenazi Jewish controls
West African ancestry	Bangladeshi ancestry	Punjabi Sikh ancestry	Ashkenazi
Yoruban ancestry	Chinese ancestry	She Chinese ancestry	Bulgarian ancestry
		Silk Road individuals	

Carlantino individuals	Korculan individuals	Caribbean Hispanic controls	Latino male controls
Carlantino female individuals	Korkula individuals	Costa Rican ancestry	Mexican American
Caucasian Eastern	Korkulan individuals	Dominican Republic ancestry	Mexican ancestry
Mediterranean ancestry	Northern Finnish founder individuals	Hispanic ancestry	Surinamese ancestry
Cilento individuals	Old Order Amish individuals	Hispanic and unknown ancestry	<b>“Native peoples”</b>
Erasmus Rucphen Family individuals	Orcadian female individuals	Hispanic asthmatic individuals	American Indian ancestry
Erasmus Ruchpen individuals	Orcadian individuals	Hispanic cases	Bashkir ancestry
European	Romanian founder cases	Hispanic child cases	Cape Verdian cases
European American cases	Russian ancestry	Hispanic child controls	Martu Australian Aboriginal ancestry
European ancestry	Sardinian cases	Hispanic controls	Native Hawaiian ancestry
European ancestry individuals	Sardinian controls	Hispanic female individuals	Native American ancestry
European ancestry cases	Sardinian individuals	Hispanic incident cases	Pima Indian ancestry
European ancestry controls	Sorbian individuals	Hispanic individual cases	Plains American Indian ancestry
European cases	South African Afrikaner ancestry	Hispanic individuals	
European child controls	Southern European ancestry	Hispanic male individuals	<b>“South Pacific Islander”</b>
European controls	Talana adult individuals	Hispanic newborn cases	Filipino ancestry
European individuals	Talana individuals	Hispanic newborn controls	Filipino female individuals
Finland founder cases	Tatar ancestry	Hispanic prevalent cases	Filipino male individuals
Finland founder controls	Val Borbera individuals	Latin American cases	Kosraen individuals
Finnish Saami individuals	Vis individuals	Latin American controls	Micronesia ancestry
French Canadian individuals	Western European ancestry	Latin American individuals	Oceania ancestry
Friuli Venezia Giulia individuals	<b>“Hispanic and Latin American”</b>	Latino cases	Papua New Guinean ancestry
Hutterite adult individuals	Brazilian ancestry	Latino child cases	Solomon Islander ancestry
Hutterite individuals	Brazilian individuals	Latino controls	
Italian isolated population individuals	Caribbean Hispanic cases	Latino current smoker	
Jewish cases		Latino female cases	
Jewish controls		Latino female controls	
		Latino individuals	
		Latino male cases	

### C. Replicating Need & Goldstein (2009)

As in the 2009 analysis conducted by David Goldstein and Anna Need, the numbers and proportions of samples presented here need to be understood not as true individual samples, but as potential replicates from the same cohorts and databases (sampling with replacement, and a growing population to choose from over time). It is possible that a large portion of the GWAS conducted simply re-sampled the same individuals from a few large databases for different studies. As such, the numbers of individuals reported in both 2009 and 2016 represent the number of times any individual from a particular ancestral background was sampled, not the number of times independent individuals were sampled from those populations. As a result of re-sampling the same individuals, there may be statistical problems with results across studies investigating the same phenotypes, but this multiple comparisons issue is outside the scope of this paper.

**Table S1. Number and proportion of GWAS participants by ancestral group in 2016 and as reported by Need & Goldstein (2009).**

	African ancestry only	Arab & Mid.E. <sup>a</sup> only	Asian ancestry only	European & Jewish Only	Hispanic & L.A. <sup>b</sup> only	Native Peoples only	South Pacific Isl. <sup>c</sup> only	Mixed & multiple groups <sup>d</sup>	Total*
<b>2016</b>									
<b>Group-based No. Studies</b>	58 (2.31%)	3 (0.12%)	349 (13.90%)	1461 (58.18%)	19 (7.57%)	5 (0.20%)	9 (0.36%)	484 (19.28%)	2388 (95.10%)*
<b>Total Samples</b>	1045224 (3.08%)	27040 (0.08%)	4795132 (14.13%)	27435555 (80.82%)	184265 (0.54%)	17929 (0.05%)	94043 (0.28%)	347928 (1.02%)	33947116 (100%)
<b>2009</b>									
<b>Group-based No. Studies</b>	1 (0.27%)	0 (0.00%)	26 (6.97%)	322 (86.33%)	3 (0.80%)	2 (0.54%)	1 (0.27%)	Mixed <sup>d</sup> Ancestry 11 (2.95%)	366 (98.12%)*
<b>Total No. Samples</b>	9840 (0.57%)	0 (0.00%)	52877 (3.18%)	1581776 (96.37%)	1019 (0.06%)	1102 (0.06%)	2622 (0.15%)	92437 (5.32%)	1741673 (105.7%)*

**a** Arab and Middle Eastern; **b** Hispanic and Latin American; **c** South Pacific Islander; **d** Mixed group-based studies refer to studies conducted on multiple single-ancestry groups, including original and replication samples; mixed ancestry samples refer to individuals with mixed (more than one) ancestry.

\* Totals are >100% due to samples overlapping ancestral categories, and <100% because studies are compared to the total pool, including those without ancestry information.

**Table S1. Legend Details.**

In Need & Goldstein’s (2009) report, the “Mixed” category describes samples with multiple ethnicities (multi-ethnic individuals), so the total percentage of samples across ancestral categories is greater than 100%. In this [2016] analysis, the “Mixed & multiple groups” category has two meanings: for GWAS (studies), it refers to studies that were done using multiple different ancestral populations and studies that included multi-ethnic individuals. Mixed ancestry samples refer to the proportion of samples that were of mixed ancestry (mirroring the 2009 analysis). Studies in the “Mixed & multiple groups” category sampled multiple single-ethnicity populations, which often included Europeans; and among the multi-ethnic populations sampled, most individuals had European mixed with some other ancestry.

In both analyses (2009, 2016) the studies without ancestry information are excluded from the table and thus the total percentages (<100%) reflect the true total number of studies (2,511 for 2016 and 373 for 2009), not the subset of studies with ancestry information. Inclusion of samples without ancestry information in the total comparison pool could bias the resulting percentages if one particular ancestral group happens to be over-represented in the group of samples without this information. Thus, studies and samples without ancestry information were excluded in the pie chart graphic and table.