

Identifying genetic variants that affect viability in large cohorts

Hakhamanesh Mostafavi¹, Tomaz Berisa², Molly Przeworski^{1,3,*}, Joseph K. Pickrell^{1,2,*}

¹ Department of Biological Sciences, Columbia University, New York, NY, USA

² New York Genome Center, New York, NY, USA

³ Department of Systems Biology, Columbia University, New York, NY, USA

*: These authors co-supervised this project.

Abstract

A number of open questions in human evolutionary genetics would become tractable if we were able to directly measure evolutionary fitness. As a step towards this goal, we developed a method to test whether individual genetic variants, or sets of genetic variants, currently influence viability. The approach consists in testing whether the frequency of an allele varies across ages, accounting for variation in ancestry. We applied it to the Genetic Epidemiology Research on Aging (GERA) cohort and to the parents of participants in the UK Biobank. In the GERA cohort, the top signal is the *APOE* ε4 allele ($P < 10^{-15}$), whereas in the UK Biobank, the strongest signals are detected in males only, and are for variants near *CHRNA3* ($P \sim 4 \times 10^{-8}$) as well as set of genetic variants that influence heart disease and lipid levels. We suggest that gene-by-environment interactions have altered the genetic architecture of viability in these two cohorts.

Introduction

A number of central questions in evolutionary genetics remain open, in particular for humans. Which types of variants affect fitness? Which components of fitness do they affect? What is the relative importance of viability and balancing selection in shaping genetic variation? Part of the difficulty is that our understanding of selection pressures acting on the human genome is based either on experiments in fairly distantly related species or cell lines, necessarily biased towards cases with simpler genetic architectures, or on indirect statistical inferences from patterns of genetic variation (Sabeti et al. 2006; Nielsen et al. 2007; Fu and Akey 2013).

The statistical inferences rely on patterns of genetic variation in present day samples (or very recently, in ancient samples (Mathieson et al. 2015)) to identify regions of the genome that appear to carry the footprint of positive selection (Nielsen et al. 2007). For example, a commonly used class of methods asks whether rates of non-synonymous substitutions between humans and other species are higher than expected from putatively neutral sites, in order to detect recurrent changes to the same protein (Yang and Bielawski 2000). Another class instead relies on polymorphism data and looks for various footprints of adaptation involving single changes of large effect (Maynard Smith and Haigh 1974). These approaches detect adaptation over different timescales and, likely as a result, suggest quite distinct pictures of human adaptation (Sabeti et al. 2006). For example, approaches that are sensitive to selective pressures acting over millions of years have identified individual chemosensory and immune-related genes (e.g., (Nielsen et al. 2005)). In contrast, approaches that contrast patterns in different human populations and are most sensitive to selective pressures active over thousands or tens of thousands of years revealed strong selective pressures on individual genes that influence human pigmentation (e.g., (Field et al. 2016; Williamson et al. 2007; Voight et al. 2006)), diet (Bersaglieri et al. 2004; Tishkoff et al. 2007; Perry et al. 2007), as well as sets of variants that shape height (Turchin et al. 2012; Robinson et al. 2015; Berg and Coop 2014). Even more recent still, studies of contemporary populations have suggested that natural selection has influenced life history traits like age at first childbirth as well as educational attainment over the course of the last century (Stearns et al. 2010; KAar et al. 1996; Byars et al. 2010; Milot et al. 2011; Troup et al. 2015; Beauchamp 2016).

Because these approaches are designed (either explicitly or implicitly) to be sensitive to a particular mode of adaptation, they provide a partial and potentially biased picture of what variants in the genome are under selection. In particular, most have much higher power to adaptations that involve strongly beneficial alleles that were rare in the population when first favored and will tend to miss selection on standing variation or adaptation involving many loci with small beneficial effects (e.g., (Przeworski, Coop, and Wall 2005; Pennings and Hermisson 2006; Teshima, Coop, and Przeworski 2006; Coop et al. 2009)). Moreover, even where these methods identify a beneficial allele, they are not informative about possible fitness trade-offs between sexes or across ages.

Here, we aim to develop a more direct and, in principle, comprehensive way to study adaptation in humans, focusing on *current* viability selection. Similar to the approach that

Alison took in comparing frequencies of the sickle cell allele in newborn and adults living in malarial environments (Allison 1964), our approach is to directly observe the effects of genotypes on survival by taking advantage of the recent availability of genotypes from a large cohort of individuals of different ages. Specifically, we test for differences in the frequency of an allele across individuals of different ages, correcting for changes in ancestry and possible batch effects. This approach is sensitive to variants that impact survival to a given age. Any genetic variant that affects survival by definition has a fitness cost, even if the cost is too small to be effectively selected against (which depends on the effective population size, the age structure of the population and the age at which the variant exerts its effects (Charlesworth 1994)). Of course a genetic variant can influence fitness without influencing survival, through effects on reproduction or inclusive fitness. Thus, our approach considers only one of the components of fitness that are likely important for human adaptation.

As a proof of principle, we apply our approach to two recent datasets: to 57,696 individuals of European ancestry from the Resource for Genetic Epidemiology Research on Aging (GERA) Cohort (Banda et al. 2015; Kvale et al. 2015) and, by proxy (Joshi et al. 2016; Pilling et al. 2016), to the parents of 95,513 individuals of European ancestry surveyed as part of the UK Biobank (UK Biobank). We do so for individual genetic variants, then jointly for sets of variants previously found to influence one of 40 polygenic traits (Pickrell et al. 2016).

Results

Our method for testing for differences in allele frequencies across age bins

If a genetic variant does not influence viability, its frequency should be the same in individuals of all ages. We therefore test for changes in allele frequency across individuals of different ages, while accounting for systematic differences in the ancestry of individuals of different ages (for example, as a result of migration patterns over decades) and genotyping batch effects. We use a logistic regression model in which we regress each individual's genotype on their age bin, their ancestry as determined by principal component analysis (PCA) (Figure S1), and the batch in which they were genotyped (see Materials and Methods for details). In this model, we treat age bin as a categorical variable; this allows us to test for a relationship between age and the frequency of an allele regardless of the functional form of this relationship. We also test a model with an interaction between age and sex, to assess whether a variant affects survival differently in the two sexes.

We first evaluated the power of this method using simulations. We considered three possible trends in allele frequency with age: (i) a constant frequency up to a given age followed by a steady decrease, i.e., a variant that affects survival after a given age (e.g., variants contributing to late-onset disorders), (ii) a steady decrease across all ages for a variant with detrimental effect throughout life, and (iii) a U-shape pattern in which the allele frequency decreases to a given age but then increases, reflecting trade-offs in the effects at young and old ages, as hypothesized by the antagonistic pleiotropy theory of aging (Williams 1957) or as may be seen if there are protective alleles whose effects are visible late in life (Bergman et al. 2007) (Figure

1). In all simulations, we used sample sizes and age distributions that matched the GERA cohort (Figure S2). For simplicity, we also assumed no population structure or batch effects across age bins (Materials and Methods). For all trends, we set a maximum of 20% change in the allele frequency from the value in the first age bin (Figure 1).

Because of the age distribution of individuals in the GERA cohort (Figure S2), our power to detect the trend is greater when most of the change in allele frequency occurs at middle ages (Figure 1). For example, for an allele with an initial allele frequency of 15% that begins to decrease in frequency among individuals at age 20, age 50, or age 70 years, there is around 20%, 90% and 60% power, respectively, to detect the trend at $P < 5 \times 10^{-8}$, the commonly-used criterion for genome-wide significance (Pe'er et al. 2008). We also experimented with a version of the model where the age bin is treated as an ordinal variable; as expected, this model is more powerful if there is a linear relationship between age and allele frequency (Materials and Methods). Since in most cases, we do not know the functional form of the relationship between age and allele frequency *a priori*, we used the categorical model for all analyses, unless otherwise noted.

In the UK Biobank, all individuals were 45-69 years old at enrollment, so the age range of the participants is restricted and our method has low power. However, the UK Biobank participants reported the age at death of their parents; following two recent studies (Joshi et al. 2016; Pilling et al. 2016), we therefore used these values (when reported) instead in our model. In this situation, we are testing for correlations between an allele frequency and the age at which the father or mother died. This approach obviously comes with the caveat that children inherit only 50% of their genome from each parent and so power is reduced (e.g., (Liu, Erlich, and Pickrell 2016)). Further, the patterns expected when considering individuals who have died differ subtly from those generated among surviving individuals. Notably, when an allele begins to decline in frequency starting at a given age (Figure 1A), there should be an *increase* in the allele frequency among individuals who died at that age, followed by a decline in frequency, rather than the steady decrease expected among surviving individuals (Figure S3, see Materials and Methods for details).

We also adapted this model to allow us to test for changes in frequency at sets of genetic variants. Many phenotypes of interest, from complex disease risk to anthropomorphic traits such as age of menarche, are polygenic (Visscher et al. 2012; He and Murabito 2014). If a polygenic trait has an effect on fitness, either directly or indirectly (i.e., through pleiotropic effects), the individual loci that influence the trait may be too subtle in their survival effects to be detectable with current sample sizes. We therefore investigated whether there is a shift across ages in sets of genetic variants that were identified as influencing a trait in genome-wide association studies (GWAS) (Table S1). Specifically, for a given trait, we calculated a polygenic score for each individual based on trait effect sizes of single variants previously estimated in GWAS and then test whether the scores vary significantly across age bins (see Materials and Methods for details). These scores are calculated under an additive model, which appears to provide a good fit to GWAS data (Cantor, Lange, and Sinsheimer 2010). If a polygenic trait is under stabilizing selection (e.g., human birth weight (Karn and Penrose 1951)), i.e., an

intermediate polygenic score is optimal, no change in the mean value of polygenic score across different ages is expected. However, if extreme values of a trait are associated with lower chance of survival, the spread of the polygenic scores should decrease with age. To consider this possibility, we tested whether the squared difference of the polygenic scores from the population mean varies significantly across age bins (see Materials and Methods for details).

Testing for changes in allele frequency at individual genetic variants

We first applied the method to the GERA cohort, using 9,010,280 filtered genotyped and imputed autosomal single-nucleotide polymorphisms (SNPs). We focused on a subset of filtered 57,696 individuals confirmed to be of European ancestry by PCA (see Materials and Methods, Figures S4 and S5). The ages of these individuals were reported in bins of 5 year intervals (distribution shown in Figure S2). We tested for significant changes in allele frequencies across these bins. For each SNP, we obtained a P value comparing a model in which the allele frequency changes with age to a null model (quantile-quantile plot shown in Figure S6A). All variants that reached genome-wide significance ($P < 5 \times 10^{-8}$) reside on chromosome 19 near the *APOE* gene (Figure 2A and Figure S7). This locus has previously been associated with longevity in multiple studies (Christensen, Johnson, and Vaupel 2006; Murabito, Yuan, and Lunetta 2012). The $\epsilon 4$ allele of the *APOE* gene is known to increase the risk of late-onset Alzheimer's disease (AD) as well as of cardiovascular diseases (Corder et al. 1993; Bennet et al. 2007). We observed a monotonic decrease in the frequency of the T allele of the $\epsilon 4$ tag SNP rs6857 (C, protective allele; T, risk allele) beyond the age of 70 years old (Figure 2B). This trend is observed for both the heterozygous and homozygous risk variants (Figure S8), and for both males and females (Figure S9). No variant reached genome-wide significance testing for age by sex interactions (quantile-quantile plot shown in Figure S6B).

We further investigated the trends in frequency with age for the other two major *APOE* alleles defined by rs7412 and rs429358 SNPs: $\epsilon 2$ (rs7412-T, rs429358-T) and $\epsilon 3$ (rs7412-C, rs429358-T), while $\epsilon 4$ is (rs7412-C, rs429358-C) (Liu et al. 2013). Unlike the $\epsilon 4$ allele, $\epsilon 2$ carriers are suggested to be at lower risk of Alzheimer's disease, cardiovascular disease, and mortality relative to the $\epsilon 3$ carriers (Liu et al. 2013; Christensen, Johnson, and Vaupel 2006). We focused on a subset of 38,703 individuals with unambiguous counts of each *APOE* allele. There is a significant change in the frequency of the $\epsilon 4$ allele with age in this subset ($P \sim 6 \times 10^{-12}$), similar to the trend observed for the tag SNP rs6857 (Figure S10). The $\epsilon 3$ allele showed the reverse trend, with a significant, monotonic increase in frequency beyond age of 70 years old ($P \sim 2 \times 10^{-8}$) (Figure S10). The enrichment of the $\epsilon 3$ allele in elderly individuals can be explained by the corresponding depletion of the $\epsilon 4$ allele, however, so does not necessarily imply an independent, protective effect of $\epsilon 3$. The frequency of the $\epsilon 2$ allele did not change significantly with age ($P \sim 0.2$), which may reflect low power for the allele frequency of ~ 0.06 (Figure S10).

We considered the possibility that some unobserved confounding variable was driving the strength of this signal at *APOE*. Since we observe two genotyped SNPs with signals similar to rs6857 within the locus, genotyping error seems unlikely to be driving the pattern (Figure S7). Another concern might be a form of ascertainment bias, in which individuals with Alzheimer's disease are underrepresented in the Kaiser Permanente Medical Care Plan. However, there is no correlation in these data between the amount of time that an individual has been enrolled in this insurance plan and the individual's *APOE* genotype (Figure S11). These observations, along with previously reported associations at this locus, argue that the allele frequency trends in Figure 2B are driven by effects of *APOE* genotype on mortality (or severe disability). Moreover, the effects that we identify are concordant with epidemiological data on the peak age of onset of Alzheimer's disease given 0 to 2 copies of *APOE* ϵ 4 (Corder et al. 1993). Thus, this case illustrates how our approach provides resolution about age effects of deleterious variants.

We estimated that we have ~93% power to detect the trend in allele frequency with age as observed for rs6857 (at a genome-wide significance level; see Materials and Methods). Using both versions of the model treating age bin as either a categorical or an ordinal variable, we have similar power to detect other potential trends considered in Figure 1, for variants as common as rs6857 and with similar magnitude of effect on survival. Yet across the genome, only *APOE* variants show a significant change in allele frequency with age for both versions of the model (Figure 2 and Figure S12), suggesting that there are few common variants in the genome with an effect on survival as strong as seen in *APOE* region.

We then turned to the UK Biobank data. We applied our method to individuals of European ancestry whose data passed our filters; of these, 88,595 had death information available for their father and 71,783 for their mother. We analyzed 590,437 genotyped autosomal variants, applying similar quality control measures as with the GERA dataset (see Materials and Methods). We tested for significant changes in allele frequencies with father's age at death and mother's age at death stratified in eight 5-year interval bins (quantile-quantile plots shown in Figure S13).

Consistent with recent studies (Joshi et al. 2016; Pilling et al. 2016), the variants showing a genome-wide significant change in allele frequency with father's age at death ($P < 5 \times 10^{-8}$) reside within a locus containing the nicotine receptor gene *CHRNA3* (Figure 3A). The A allele of the *CHRNA3* SNP rs1051730 (G, major allele; A, minor allele) has been shown to be associated with increased smoking quantity among individuals who smoke (Tobacco and Genetics Consortium 2010). We observed a linear decrease in the frequency of the A allele of rs1051730 throughout almost all age ranges (Figure 3B). This allele did not show a significant trend with age in GERA ($P \sim 0.45$, Figure S14).

For mother's age at death, a SNP in a locus containing the *MEOX2* gene reached genome-wide significance (Figure 3C). The C allele of rs4721453 (T, major allele; C, minor allele) increases in frequency in the age bin centered at 76 years old (Figure S15), i.e., there is an enrichment among individuals that died at 74 to 78 years of age, which corresponds to a deleterious effect

of the C allele in this period. The trend is similar and nominally significant for other genotyped common SNPs in moderate linkage disequilibrium with rs4721453 (Figure S15). Also, the signal for rs4721453 remains nominally significant when using subsets of individuals genotyped on the same genotyping array: 44,552 individuals on the UK Biobank Axiom array ($P \sim 7 \times 10^{-5}$) and 25,231 individuals on the UK BiLEVE Axiom array ($P \sim 10^{-4}$). These observations suggest that the result is not due to genotyping errors, but it is not reproduced in GERA ($P \sim 0.17$, Figure S16) and so it remains to be replicated. *APOE* variants were among the top nominally significant variants ($P < 10^{-7}$) (Figure 3C). At the *APOE* SNP rs769449 (G, major allele; A, minor allele), there is an increase in the frequency of A allele at around 70 years old before subsequent decrease (Figure 3D). This pattern is consistent with our finding in GERA (of a monotonic decrease beyond 70 years of age), considering the difference in patterns expected between allele frequency trends with age among survivors versus individuals who died (Figure S3).

We note that by considering parental age at death of the UK Biobank participants (as done also in (Joshi et al. 2016; Pilling et al. 2016)), we introduce a bias towards older participants (who are more likely to have deceased parents, Figure S17A). We confirmed that our top signals are not significantly affected by such potential bias, observing similar trends in allele frequency with parental age at death when conditioning on age of the participants (Figure 17B).

We further tested for trends in allele frequency with parental age at death that differ between fathers and mothers focusing on 62,719 individuals with age at death information for both parents. No SNP reached genome-wide significance level (Figure 18A). The rs4721453 near the *MEOX2* gene and *APOE* variant rs769449 showed nominally significant sex effects ($P \sim 7 \times 10^{-8}$ and $P \sim 2 \times 10^{-3}$, respectively), with stronger effects in females (Figure 18B). Variants near the *CHRNA3* locus were nominally significant when using the model with parental age at deaths treated as ordinal variables (rs11858836, $P \sim 6 \times 10^{-4}$), with stronger effects in males (Figure 18B).

Testing for changes in allele frequency at trait-associated variants

We next applied our model to polygenic traits, rather than individual genetic variants. We focused on 40 polygenic traits, including disease risk and anthropomorphic traits of evolutionary importance such as age at menarche (AAM), for which a large number of common variants have been mapped in GWAS (see Table S1 for the list of traits and number of loci) (Pickrell et al. 2016). For each individual and each trait, we calculated a polygenic score, and then tested whether this polygenic score, or its squared difference from the mean in the case of stabilizing selection, differs among individuals in different age bins.

In the GERA cohort, no trait reached statistical significance, after accounting for multiple tests (Figure 4A). The strongest signal is a suggestive increase in the polygenic score for age at menarche in older ages ($P \sim 3.4 \times 10^{-3}$, without correction for multiple testing; Figure 4B), consistent with epidemiological studies suggesting early puberty timing to be associated with various adverse health outcomes (Day, Elks, et al. 2015). We did not exclude males for analysis

of age at menarche, because of the strong genetic correlation between the timing of puberty in males and females (Day, Bulik-Sullivan, et al. 2015). For disease traits potentially decreasing the chance of survival with increasing age, a monotonic decrease in polygenic score with age is plausible. Therefore, we also applied our version of the model with age treated as an ordinal variable, for which Alzheimer's disease (AD) (excluding the *APOE* locus) showed the strongest signal ($P \sim 5.1 \times 10^{-3}$, see Figure 4C), indicative of a decrease in chance of survival with increased genetic risk of AD (Figure 4D). No trait showed significant sex by age effect (Figure S19), or significant change in the squared difference of polygenic score from the mean with age (Figure S20).

In the UK Biobank, consistent with two recent studies (Pilling et al. 2016; Marioni et al. 2016), several cardiovascular disease risk traits showed significant change in polygenic score with father's age at death (Figure 5A): total cholesterol ($P < 10^{-6}$), coronary artery disease ($P < 10^{-5}$), low-density lipoproteins ($P < 10^{-4}$), and body mass index ($P < 10^{-4}$). The decline in score with age for these traits is approximately linear with age (Figure 5B-E). Testing for sex by age interaction, high-density lipoproteins showed the strongest signal, with seemingly distinct trends in males and females ($P \sim 3 \times 10^{-4}$, Figures S21). Using the model with parental age at deaths treated as ordinal variables, total cholesterol showed stronger effects in males ($P \sim 7 \times 10^{-4}$, Figures S21).

No trait showed significant change in polygenic score with mother's age at death (Figure 5F). Also, no trait showed significant change in the squared difference of polygenic score from the mean with father's or mother's age at death (Figure S22). We confirmed that our different findings between the GERA cohort and the UK Biobank are not driven by using slightly different sets of trait-associated SNPs (considering that trait associated SNPs passing our filters were not identical), finding similar results when using SNPs that passed quality control steps in both datasets (Figure S23).

Discussion

We introduced a new approach to identify genetic variants that affect survival to a given age and thus to directly observe viability selection ongoing in humans. Attractive features of the approach include that we do not need to make a decision a priori about which traits matter to viability and focus not on an endpoint (e.g., lifespan) but on any shift in allele frequencies with age, thereby learning about the precise ages at which effects are manifest and possible differences between sexes.

To illustrate the potential of our approach, we performed a scan for genetic variants that impact age-specific mortality in the GERA and the UK Biobank cohorts. We only found a few variants, the majority of which were identified in previous studies. This result is in some ways expected: available data only provide high power to detect effects of common variants (>15-20%) on survival (Figure 1), yet if these variants were under viability selection, we would not expect them to be common, short of strong balancing selection due to trade-offs between sexes, ages or environments. As sample sizes increase, however, the approach introduced here

should provide a much more comprehensive picture of viability selection in humans. To illustrate this point, we repeated our power simulation with 500,000 samples, and found that we should have high power to detect the trends for alleles at a couple percent frequency (Figure S24).

Already, however, this application raises a number of interesting questions about the nature of viability selection in humans. Notably, we discovered only a few variants influencing viability in the two cohorts, all of which exert their effect late in life. On first thought, it may be concluded that such variants are neutrally-evolving. We would argue that if anything, our findings of only a few common variants with effects on survival only late in life suggest the opposite: that even variants with late onset effects have been weeded out by purifying selection. Indeed, unless the number of loci in the genome that could give rise to such variants (i.e., the mutational target size) is tiny, other variants such as *APOE* must often arise. That they are not observed when we have very high power to detect them suggests they are kept at lower frequency by purifying selection. Why might they be selected despite affecting survival only at old ages? Possible explanation include that they decrease the direct fitness of males sufficiently to be effectively selected (notably given the large, recent effective population size of humans (Gazave et al. 2014)) or that they impact the inclusive fitness of males and females. If this explanation is correct, it raises the question of why *APOE* $\epsilon 4$ has not been weeded out. We speculate that the environment today has changed in such a way that has made this allele more deleterious recently. For example, it has been proposed that the evolution of this allele has been influenced by changes in physical activity (Raichlen and Alexander 2014).

Also interesting is the fact that our results differ markedly between the GERA cohort of individuals from California and the parents of the individuals in the UK Biobank. These differences are most notable at the *CHRNA3* locus and the sets of polygenic traits; both show strong signals among fathers of the individuals in the UK Biobank but not in GERA. We cannot rule out that these differences are simply due to power, as comparisons between the two cohorts are complicated by a number of factors. However, the analysis of mothers and fathers of individuals in the UK Biobank should have similar power and there too, we see marked differences. Together, these observations point to strong gene-environment interactions of lifespan, such that the environments of males in the UK in the mid to late 1900s, females in this same period in the UK, and California in the late 1900s were different enough to change the genetic architecture of lifespan.

The *CHRNA3* locus, in which variants are associated with the amount of smoking among smokers, even suggests a specific relevant environmental factor. Smoking prevalence has decreased significantly over the past few decades in both the UK and California, however, men in the UK consistently smoked more than women in the UK and people from California: from 1970 to 2000, smoking prevalence decreased from around 70% to 36% in middle-aged men from UK, compared to from around 50% to 28% in middle-aged women from UK, and from around 40% to 20% in Californians (Peto et al. 2000; Pierce et al. 2010). These epidemiological patterns are potentially consistent with our observation of more pronounced effect on male than female age at death among parents of UK Biobank participants, and the lack of significant

pattern in either sex in GERA, particularly that GERA individuals are around a generation time younger than parents of UK Biobank participants. In any case, these results highlight the utility of cohorts with both genetic and environmental information, and ideally a range of ages.

Moving forward, application of our approach to the millions of samples in the pipeline (such as the UK Biobank (Sudlow et al. 2015), the Precision Medicine Initiative Cohort Program (Collins and Varmus 2015), and the Vanderbilt University biobank (BioVU)(Roden et al. 2008)), in which the viability effects of rare as well as common alleles can be examined, should provide a comprehensive answer to the question of which loci affect survival, helping to address long-standing open questions such as the relative importance of viability selection in shaping genetic variation and the extent to which genetic variation is maintained by fitness trade-offs between sexes or across ages.

Materials and Methods

1. Datasets

1.1. GERA cohort

We performed our analyses on the data for 62,318 participants of the Kaiser Permanente Northern California multi-ethnic Genetic Epidemiology Research on Adult Health and Aging (GERA) cohort, self-reported to be “White-European American”, “South Asian”, “Middle-Eastern” or “Ashkenazi” but no other ethnicities, among a list of 23 choices on the GERA survey, and genotyped on a custom array at 670,176 SNPs designed for Non-Hispanic White individuals (Banda et al. 2015; Kvale et al. 2015). We determined the age of the participants and the number of years they were enrolled in the Kaiser Permanente Medical Insurance Plan at the time of the survey (year 2007).

1.2. UK Biobank

We performed our analyses on the data for 120,286 participants of the UK Biobank inferred to have European ancestry, and genotyped on the UK Biobank Axiom or the UK BiLEVE Axiom SNP arrays at total of 847,441 SNPs (UK Biobank 2015b; Sudlow et al. 2015).

2. Quality control (QC)

2.1. GERA cohort

We used PLINK v1.9 (Chang et al. 2015) to remove individuals with missing sex information or with a mismatch between genotype data and sex information, individuals with <96% call rate, and related individuals. We validated self-reported European ancestries using principal component analysis (PCA), see below, and removed individuals identified as non-European (Figures S4 and S5). In the end, 57,696 individuals remained.

Using PLINK, we removed SNPs with <1% minor allele frequency, SNPs with <95% call rate, and SNPs failing a Hardy-Weinberg equilibrium test with $P < 10^{-8}$ (filtering based on HWE test

could potentially exclude true signals of viability selection, if selection coefficients were very large (Meyer et al. 2012), but this possibility is much less likely than genotyping error). We additionally tested for a correlation between age (or sex) and missingness, which can induce artificial change in the allele frequencies as a function of age (or sex). We thus removed SNPs showing a significant age-missingness or sex-missingness correlation, defined as a chi-squared test with $P < 10^{-7}$. After these steps, 599,659 SNPs remained.

We imputed the genotypes of the filtered GERA individuals using post-QC SNPs, and using the 1000 Genomes phase 3 haplotypes as a reference panel (The 1000 Genomes Project Consortium 2015). We phased observed genotypes using EAGLE v1.0 software (Loh, Palamara, and Price 2016). The inferred haplotypes were then passed to IMPUTE2 v2.3.2 software for imputation in chunks of 1Mb, using the default parameters of the software (Howie, Donnelly, and Marchini 2009). To gain computational speed, variants with $<0.5\%$ minor allele frequency in the 1000 Genomes European populations were removed from the reference panel. This step should not affect our analysis because our statistical model is not well powered for rare variants, given the GERA data sample size. We called imputed genotypes with posterior probability >0.9 , and then filtered the imputed genotypes, removing variants with IMPUTE2 info score <0.5 and with minor allele frequency $<1\%$. We also used imputation with leave-one-out approach (Marchini et al. 2007) to impose a second stage of QC on genotyped SNPs, removing SNPs that were imputed back with high reported certainty (info score >0.5) and with $<90\%$ concordance between the imputed and the original genotypes. These yielded a total of 9,010,280 imputed and genotyped SNPs.

For our analysis of the *APOE* alleles ($\epsilon 2$, $\epsilon 3$ and $\epsilon 4$) which are defined by rs7412 and rs429358 SNPs (Liu et al. 2013), given the lack of tag SNPs for all three alleles, we kept a subset of 38,703 individuals with no poorly-imputed genotypes for these two SNPs, for whom the count of each *APOE* allele could be determined unambiguously.

2.2. UK Biobank

In the UK Biobank, we obtained sets of genotype calls and the output of imputation as performed by the UK Biobank researchers (UK Biobank 2015b, 2015a). We first applied QC metrics to the autosomal genotyped SNPs. We used PLINK to remove SNPs with $<1\%$ minor allele frequency, SNPs with $<95\%$ call rate, and SNPs failing a Hardy-Weinberg equilibrium test with $P < 10^{-8}$. These filters were applied separately to SNPs genotyped on the UK Biobank Axiom and the UK BiLEVE Axiom arrays. Then, we divided the genotyped SNPs into three sets (SNPs specific to either array and shared SNPs) and then performed additional QC on each set separately: we removed SNPs with significant allele frequency difference between genotyped and imputed calls (chi-squared test $P < 10^{-5}$) and SNPs showing a significant correlation between missingness and age or sex of the participants, as well as with participants' father's or mother's age at death (chi-squared test $P < 10^{-7}$). We then extracted this list of SNPs from the imputed genotype files available from the UK Biobank (we did not use the full set of imputed genotypes). From this set, we removed SNPs with $<1\%$ minor allele frequency, SNPs with $<95\%$ call rate, and SNPs failing a Hardy-Weinberg equilibrium test with $P < 10^{-8}$, yielding 590,437

SNPs. For variants influencing quantitative traits, we first extracted them from imputed genotypes, and then imposed the same QC measures as above.

Each participant was asked to provide the age at death of their father and their mother (if applicable) on each assessment visit. For each participant that reported an age at death of father and/or mother, we averaged over the ages reported at recruitment and any subsequent repeat assessment visits, and used PLINK to exclude individuals with >5 year variation in their answers across visits (around 800 individuals). We also removed adopted individuals, as well as individuals with a mismatch between genotype data and sex information, resulting in 88,595 individuals with age at death information for their father, 71,783 individuals for their mother, and 62,719 individuals for both parents.

3. Principal Component Analysis

We performed PCA, using the EIGENSOFT v6.0.1 package with the fastpca algorithm (Price et al. 2006; Galinsky, Bhatia, et al. 2016), for two purposes: (i) as a quality control on individuals to validate self-reported European ancestries, and (ii) to correct for population structure in our statistical model.

3.1. European ancestry validation

We used more stringent QC criteria specifically for the PCA, compared to the QC steps described above. We filtered a subset of 157,277 SNPs in GERA and 220,447 in the UK Biobank, retaining SNPs shared between the datasets and the 1000 Genomes phase 3 data, removing non-autosomal SNPs, SNPs with <1% minor allele frequency, SNPs with <99% call rate, and SNPs failing a Hardy-Weinberg equilibrium test with $P < 10^{-6}$. We then performed LD-pruning using PLINK with pairwise $r^2 < 0.2$ in windows of 50 SNPs shifting every 10 SNP. We used these SNPs to infer principal components for the 1000 Genomes phase 3 data (The 1000 Genomes Project Consortium 2015). We then projected individuals onto these PCs. In GERA, we observed that the majority of individuals have European ancestry, and marked individuals with PCs deviating from the population mean, for any of the first six PCs, as non-European (Figures S4 and S5). In the UK Biobank, the European ancestry of all individuals was confirmed (not shown).

3.2. Control for population structure

After the main QC stage, additional QC steps (as in *European ancestry confirmation*) were implemented for PCA. In the UK Biobank, we also removed inversion variants on chromosome 8 which otherwise dominate the PC2 (not shown). A subset of 156,721 SNPs in GERA and 207,657 SNPs in the UK Biobank was then used to infer PCs for individuals passing QC (Figure S1). The first 10 PCs were used as covariates in our statistical model.

4. Quantitative Traits

We downloaded the list of variants contributing to 40 quantitative traits and their effect sizes recently described in Pickrell et al. (Pickrell et al. 2016), from: https://github.com/PickrellLab/gwas-pw-paper/tree/master/all_single. For age at menarche, we use a set of variants identified by running gwas-pw (Pickrell et al. 2016) on the Perry et al.

and 23andMe GWAS studies, effectively performing a meta-analysis. For all traits, we used variants that were genotyped/imputed with high quality in our data (see Table S1).

5. Statistical Model

5.1. An individual variant

Using a logistic regression we predict the genotype of individual j (the counts of an arbitrarily selected reference allele, $G_{ij} = 0, 1$ or 2) at variant i , using the individual's ancestry, the batch at which the individual was genotyped, and individual's age (as well as sex, see below) as explanatory variables. Specifically, distribution of G_{ij} is $\text{Bin}(2, p_{ij})$, where p_{ij} , the probability of observing the reference allele for individual j at variant i , is related to explanatory variables as:

$$\log\left(\frac{p_{ij}}{1 - p_{ij}}\right) = \alpha + \sum_{l=1}^{10} \beta_l PC_{lj} + \sum_m \gamma_m I_{j \in \text{BATCH}_m} + \sum_n \kappa_n J_{j \in \text{BIN}_n}$$

where β_l is the effect of principal component l (to account for population structure), γ_m is the effect of being in batch m (to account for potential systematic differences between genotyping packages), κ_n is the effect of being in age bin n , obtained by regression across individuals with non-missing genotypes at variant i , and I and J are indicator variables for the genotyping batch and age bin, respectively. In the version of the model in which we treat age as an ordinal variable, we replace J age bin variables with one age variable. In the GERA dataset, age binning is over the age of the participants in 14 categories, from age 19 onwards, in 5-year intervals, and for the UK Biobank, it is over 8 categories for the age at death of father or mother, from age 63 onwards, in 5-year intervals. In the UK Biobank, we included all ages at death below 63 in one age bin to minimize the potential noise caused by accidental deaths at young ages.

We tested for an effect of age categories by a likelihood ratio test with a null model using only the covariates (PCs and batch terms) ($H_0: \kappa_n = 0$) and an alternative also including age terms as predictors ($H_1: \kappa_n \neq 0$):

$$\begin{cases} H_0: \log\left(\frac{p_{ij}}{1 - p_{ij}}\right) = \alpha + \sum_{l=1}^{10} \beta_l PC_{lj} + \sum_m \gamma_m I_{j \in \text{BATCH}_m} \\ H_1: \log\left(\frac{p_{ij}}{1 - p_{ij}}\right) = \alpha + \sum_{l=1}^{10} \beta_l PC_{lj} + \sum_m \gamma_m I_{j \in \text{BATCH}_m} + \sum_n \kappa_n J_{j \in \text{BIN}_n} \end{cases}$$

To test for age by sex effects in GERA we included two sets of additional predictors. The first consists in two indicator variables for sex, K_{male} and K_{female} , which are included to capture possible sex effects induced by potential genotyping errors or mismapping of sex chromosome linked alleles (we note that because of Hardy-Weinberg equilibrium, mean allele frequency difference between males and females are not expected). The second set of predictors consists in age by sex terms, $J \times K$. We then compare a model with age and sex terms as predictors to a model also including age by sex terms. To test for sex effects in the UK Biobank, we compared a model with both father and mother age terms separately as predictors to a model with one set

of age categories for average age at death of both parents, only for individuals reporting the age at death for both parents. In all models PCs and batch terms were incorporated as covariates.

5.2. Set of variants

As for the model described above for an individual variant, we investigated age and age by sex effects on quantitative traits for which large number of large common genetic variants have been identified in genome-wide association studies (GWAS). For a given trait, we used a linear regression with the same covariates and predictors as for the model for an individual variant, to predict the polygenic score for individual j , S_j , by summing the previously estimated effect of single variants assuming additivity and that the effect sizes are similar in the GWAS panels and the cohorts considered here:

$$S_j = \alpha + \sum_{l=1}^{10} \beta_l PC_{lj} + \sum_m \gamma_m I_{j \in BATCH_m} + \sum_n \kappa_n J_{j \in BIN_n} + \varepsilon_j$$

S_j is calculated as $\sum a_i G_{ij} + \sum 2a_i q_i - \hat{S}$, where the first sum is across variants with non-missing genotypes, a_i is the effect size for the arbitrary selected reference allele at variant i , the second sum is across the variants with missing genotypes estimating their contribution assuming Hardy-Weinberg equilibrium where q_i is the frequency of the alternate allele, and \hat{S} is the expected polygenic score calculated for the 1000 Genomes European populations, subtracted to standardize the calculated scores. Likelihood ratio tests, as described above, were used to test for age and age by sex effects.

All Manhattan and quantile-quantile plots were generated using qqman (Turner 2014) and GWASTools (Gogarten et al. 2012) packages, respectively.

6. Power simulations

We ran simulations to determine the power of our statistical model to detect deviation of allele frequency trends with age across 14 age categories mimicking the GERA individuals (57,696 individuals with age distribution as in Figure S2) from a null model, which for simplicity was no change in frequency with age, i.e., no changes as a result of age-dependent variation in population structure and batch effects. For a given trend in frequency of an allele with age, we generated 1000 simulated trends where the distribution of the number of the alleles in age bin i is $Bin(2N_i, f_i)$, where N_i and f_i are the sample size and the sample allele frequency in bin i . We then estimated the power to detect the trend as the fraction of cases in which $P < 5 \times 10^{-8}$, by a chi-squared test.

7. Survival simulations

We ran simulations to investigate the relationship between allele frequency with age of the survived individuals and the age of the individuals who died in a cohort. We simulated 2×10^6 individuals going forward in time in 1 year increments. For each time step forward, we tuned the chance of survival of the individuals based on their count of a risk allele for a given SNP such

that the number of individuals dying in the increment complies with: (i) a normal distribution of ages at death with mean of 70 years and standard deviation of 13 years, roughly as is observed for parental age at deaths in the UK Biobank, and (ii) a given frequency of the risk allele among those who survive. Specifically, we modeled the survival rate of the population, S , as the weighted mean for 2 alleles carriers, S_2 , 1 allele carriers, S_1 , and non-carriers, S_0 :

$$S(x) = \sum_{i=0}^2 f_i S_i(x)$$

where f denotes the frequency of genotypes in the population and x denotes the age. S_i and S are related: $S_i(x) = S(x) f_i(x)/f_i$, where $f_i(x)$ is the genotype frequency among individuals survived up to age x . Given a trend in allele frequency with age, we calculated genotype frequencies with age assuming Hardy-Weinberg equilibrium, and then estimated genotype dependent chance of survival, $S_i(x)$, taking $S(x)$ as the survival function for $N(70, 13^2)$.

Figures

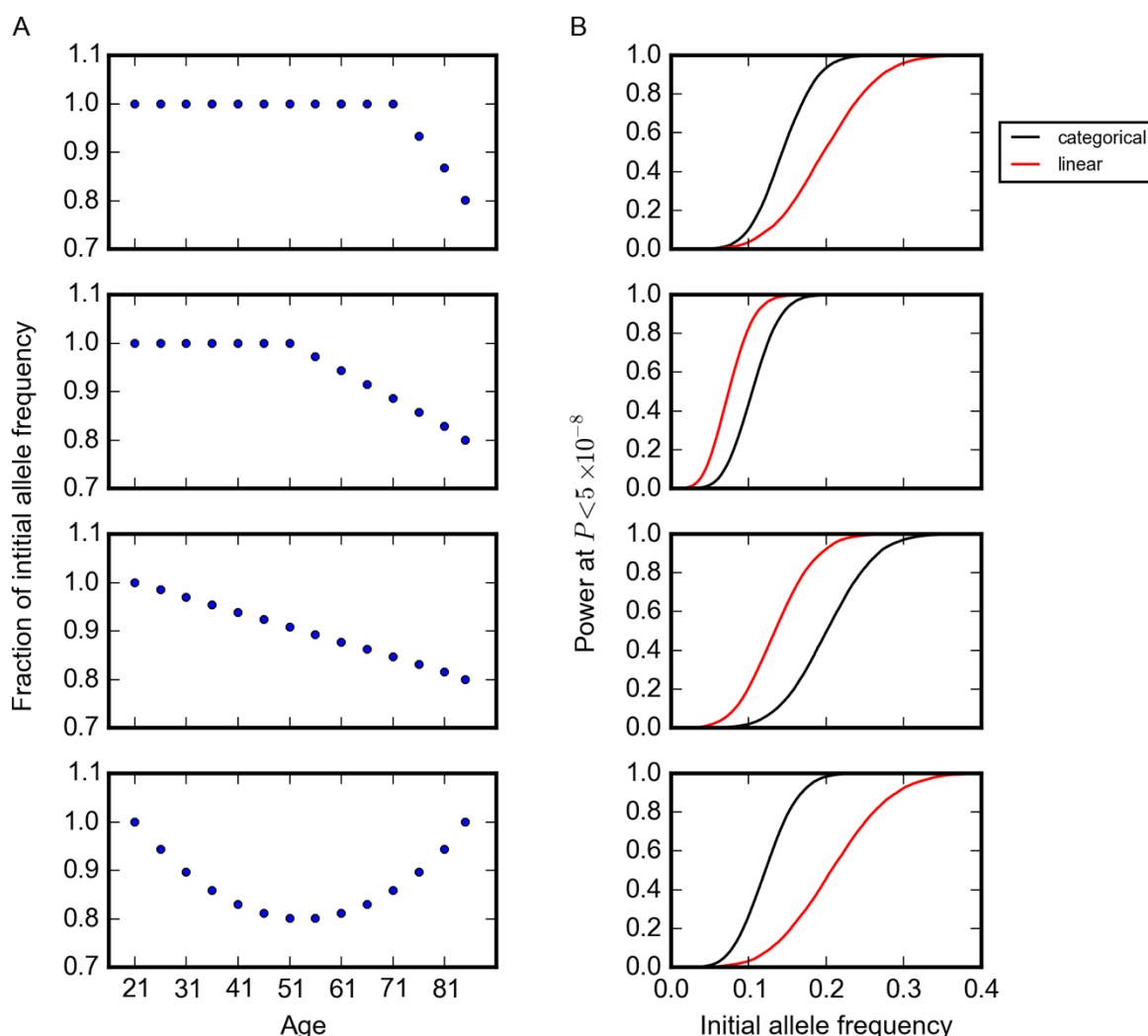


Figure 1. Power of the model to detect changes in allele frequency with age. (A) Trends in allele frequency with age considered in simulations. The y-axis indicates allele frequency normalized to the frequency in the first age bin. (B) Power to detect the trends in (A) at $P < 5 \times 10^{-8}$, given the sample size per age bin in the GERA cohort (Figure S2 and total sample size of 57,696). Shown are results using models with age treated as a categorical (black) or an ordinal (red) variable, assuming no change in population structure and batch effects across age bins. The curves show simulation results sweeping allele frequency values with an increment value of 0.001 (1000 simulations for each allele frequency) smoothed using a Savitzky-Golay filter using the SciPy package (Jones, Oliphant, and Peterson 2001).

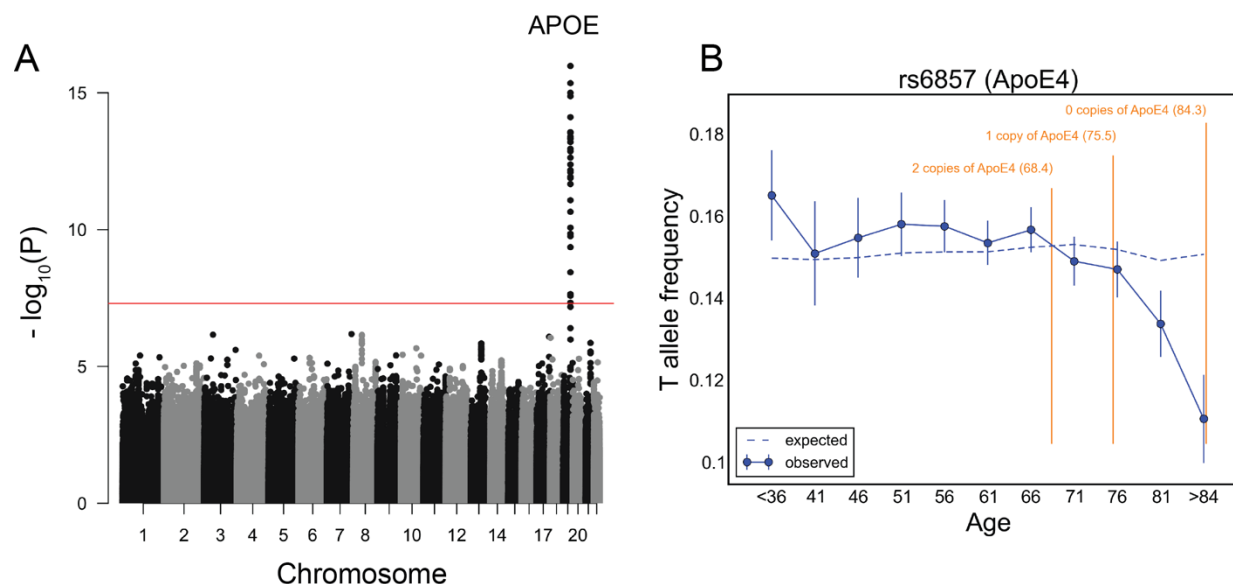


Figure 2. Testing for the influence of single genetic variants on age-specific mortality in the GERA cohort. (A) Manhattan plot for change in allele frequency with age P values. Red line marks the $P = 5 \times 10^{-8}$ threshold. (B) Allele frequency trajectory of rs6857, a tag SNP for *APOE* $\epsilon 4$ allele, with age. Data points are mean frequency of the risk allele within 5-year interval age bins (and 95% confidence interval), with the center of the bin indicated on the x-axis. Bins with ages below 36 years are merged into one bin because of the relatively small sample sizes per bin. The dashed line shows the expected frequency based on the null model accounting for confounding batch effects and changes in ancestry (see Materials and Methods). In orange are the mean age of onsets of Alzheimer's disease for carriers of 0, 1 or 2 copies of the *APOE* $\epsilon 4$ allele (Corder et al. 1993).

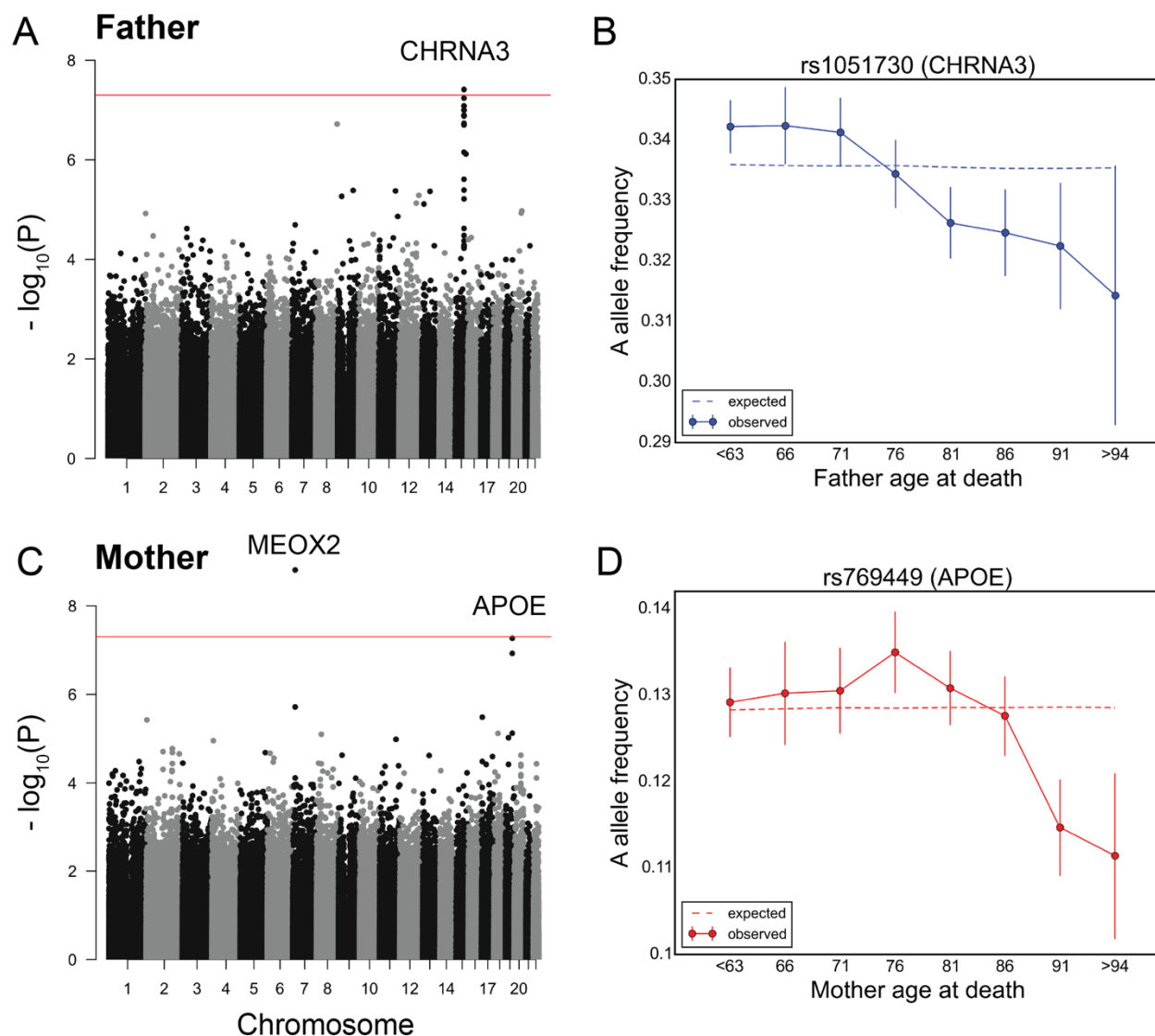


Figure 3. Testing for the influence of single genetic variants on age-specific mortality in the UK Biobank. (A) Manhattan plot for change in allele frequency with age at death of fathers P values. (B) Allele frequency trajectory of rs1051730, within *CHRNA3* locus, with father's age at death. (C) Manhattan plot for change in allele frequency with age at death of mothers P values. (D) Allele frequency trajectory of rs769449, within the *APOE* locus, with mother's age at death. Red lines in (A) and (C) mark $P = 5 \times 10^{-8}$ threshold. Data points in (B) and (D) are mean frequency of the risk allele within 5-year interval age bins (and 95% confidence interval), with the center of the bin indicated on the x-axis. The dashed line shows the expected frequency based on the null model, accounting for confounding batch effects and changes in ancestry (see Materials and Methods).

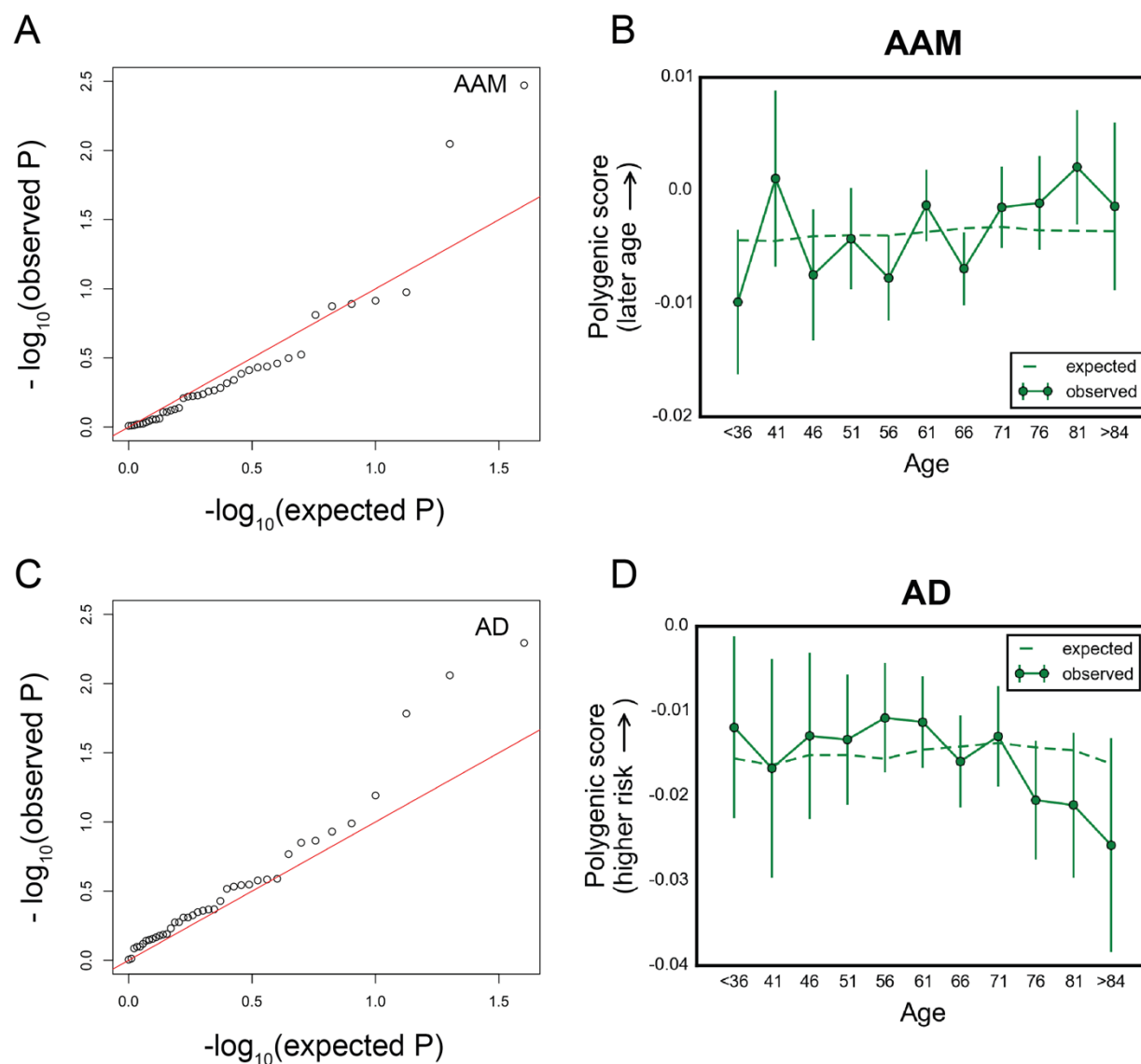


Figure 4. Testing for the influence of set of trait-associated variants on age-specific mortality in the GERA cohort. (A) Quantile-quantile plots for change in the polygenic score of 40 traits (see Table S1) with age treated as a categorical variable. (B) Trajectory of polygenic score of age at menarche with age. (C) Same as (A) but with age treated as an ordinal variable. (D) Trajectory of polygenic score of Alzheimer's disease (excluding the *APOE* locus) with age. The red lines in (A) and (C) indicate distribution of the P values under the null. See Table S2 for P values for all traits. Data points in (B) and (D) are mean polygenic score within 5-year interval age bins (and 95% confidence interval), with the center of the bin indicated on the x-axis. The dashed line shows the expected score based on the null model accounting for confounding batch effects and changes in ancestry.

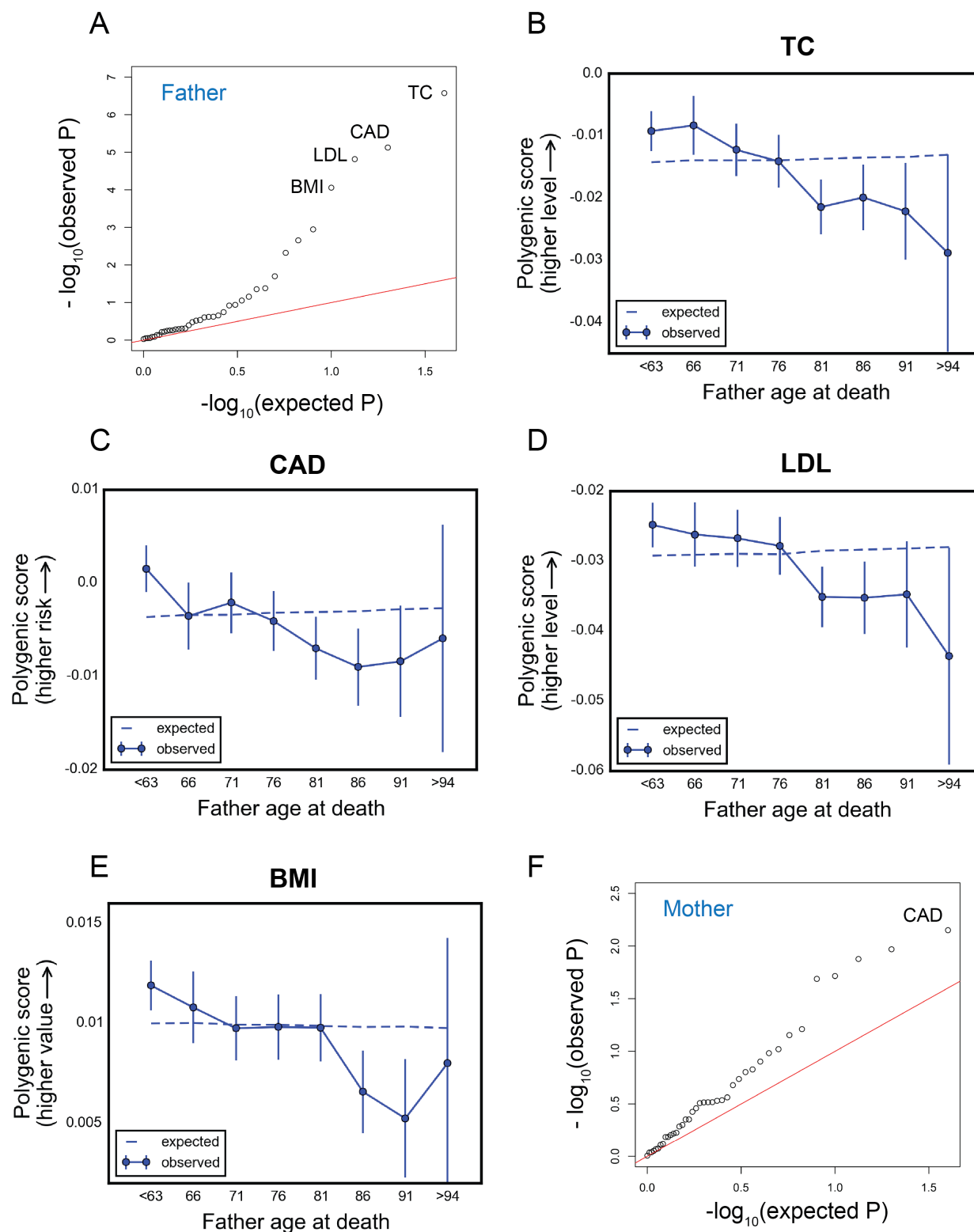


Figure 5. Testing for the influence of set of trait-associated variants on age-specific mortality in the UK Biobank. (A) Quantile-quantile plots for change in the polygenic score of 40 traits (see Table S1) with age at death of fathers of the UK Biobank participants. (B)-(E) Trajectory of

polygenic score of traits showing significant change with age at death of fathers: total cholesterol (B), coronary artery disease (C), low-density lipoproteins (D), and body mass index (E). (F) Quantile-quantile plots for change in the polygenic score of traits in (A) with age at death of mothers of the UK Biobank participants. The red lines in (A) and (F) indicate the distribution of the P values under the null. See Table S2 for P values for all traits. Data points in (B)-(E) are mean polygenic score within 5-year interval age bins (and 95% confidence interval), with the center of the bin indicated on the x-axis. The dashed line shows the expected score based on the null model, accounting for confounding batch effects and changes in ancestry.

Acknowledgments

We thank Guy Sella for helpful discussions. This research has used the UK Biobank Resource (application number 11138), and was funded in part by Columbia University (a Research Initiative in Science and Engineering grant to MP and JKP) and the National Institutes of Health (grant R01MH106842 to JKP and R01GM115889 to Guy Sella). These data analyses were approved by the Columbia University Institutional Review Board, protocols AAAQ2700 and AAAP0478.

References

- Allison, Anthony Clifford. 1964. "Polymorphism and natural selection in human populations." In *Cold Spring Harbor Symposia on Quantitative Biology*, 137-49. Cold Spring Harbor Laboratory Press.
- Banda, Yambazi, Mark N Kvale, Thomas J Hoffmann, Stephanie E Hesselson, Dilrini Ranatunga, Hua Tang, Chiara Sabatti, Lisa A Croen, Brad P Dispensa, and Mary Henderson. 2015. 'Characterizing race/ethnicity and genetic ancestry for 100,000 subjects in the Genetic Epidemiology Research on Adult Health and Aging (GERA) cohort', *Genetics*, 200: 1285-95.
- Beauchamp, Jonathan P. 2016. 'Genetic evidence for natural selection in humans in the contemporary United States', *Proceedings of the National Academy of Sciences*, 113: 7774-79.
- Bennet, Anna M, Emanuele Di Angelantonio, Zheng Ye, Frances Wensley, Anette Dahlin, Anders Ahlbom, Bernard Keavney, Rory Collins, Björn Wiman, and Ulf de Faire. 2007. 'Association of apolipoprotein E genotypes with lipid levels and coronary risk', *JAMA*, 298: 1300-11.
- Berg, Jeremy J, and Graham Coop. 2014. 'A population genetic signal of polygenic adaptation', *PLoS Genetics*, 10: e1004412.
- Bergman, Aviv, Gil Atzmon, Kenny Ye, Thomas MacCarthy, and Nir Barzilai. 2007. 'Buffering mechanisms in aging: a systems approach toward uncovering the genetic component of aging', *PLoS Computational Biology*, 3: e170.
- Bersaglieri, Todd, Pardis C Sabeti, Nick Patterson, Trisha Vanderploeg, Steve F Schaffner, Jared A Drake, Matthew Rhodes, David E Reich, and Joel N Hirschhorn. 2004. 'Genetic signatures

- of strong recent positive selection at the lactase gene', *The American Journal of Human Genetics*, 74: 1111-20.
- Byars, Sean G, Douglas Ewbank, Diddahally R Govindaraju, and Stephen C Stearns. 2010. 'Natural selection in a contemporary human population', *Proceedings of the National Academy of Sciences*, 107: 1787-92.
- Cantor, Rita M, Kenneth Lange, and Janet S Sinsheimer. 2010. 'Prioritizing GWAS results: a review of statistical methods and recommendations for their application', *The American Journal of Human Genetics*, 86: 6-22.
- Chang, Christopher C, Carson C Chow, Laurent CAM Tellier, Shashaank Vattikuti, Shaun M Purcell, and James J Lee. 2015. 'Second-generation PLINK: rising to the challenge of larger and richer datasets', *Gigascience*, 4: 1.
- Charlesworth, Brian. 1994. *Evolution in age-structured populations* (Cambridge University Press Cambridge).
- Christensen, Kaare, Thomas E Johnson, and James W Vaupel. 2006. 'The quest for genetic determinants of human longevity: challenges and insights', *Nature Reviews Genetics*, 7: 436-48.
- Collins, Francis S, and Harold Varmus. 2015. 'A new initiative on precision medicine', *New England Journal of Medicine*, 372: 793-95.
- Coop, Graham, Joseph K Pickrell, John Novembre, Sridhar Kudaravalli, Jun Li, Devin Absher, Richard M Myers, Luigi Luca Cavalli-Sforza, Marcus W Feldman, and Jonathan K Pritchard. 2009. 'The role of geography in human adaptation', *PLoS Genetics*, 5: e1000500.
- Corder, EH, AM Saunders, WJ Strittmatter, DE Schmechel, PC Gaskell, GWet Small, AD Roses, JL Haines, and Margaret A Pericak-Vance. 1993. 'Gene dose of apolipoprotein E type 4 allele and the risk of Alzheimer's disease in late onset families', *Science*, 261: 921-23.
- Day, Felix R, Brendan Bulik-Sullivan, David A Hinds, Hilary K Finucane, Joanne M Murabito, Joyce Y Tung, Ken K Ong, and John RB Perry. 2015. 'Shared genetic aetiology of puberty timing between sexes and with health-related outcomes', *Nature communications*, 6.
- Day, Felix R, Cathy E Elks, Anna Murray, Ken K Ong, and John RB Perry. 2015. 'Puberty timing associated with diabetes, cardiovascular disease and also diverse health outcomes in men and women: the UK Biobank study', *Scientific Reports*, 5.
- Field, Yair, Evan A Boyle, Natalie Telis, Ziyue Gao, Kyle J Gaulton, David Golan, Loic Yengo, Ghislain Rocheleau, Philippe Froguel, and Mark I McCarthy. 2016. 'Detection of human adaptation during the past 2,000 years', *bioRxiv*: 052084.
- Fu, Wenqing, and Joshua M Akey. 2013. 'Selection and adaptation in the human genome', *Annual review of genomics and human genetics*, 14: 467-89.
- Galinsky, Kevin J, Gaurav Bhatia, Po-Ru Loh, Stoyan Georgiev, Sayan Mukherjee, Nick J Patterson, and Alkes L Price. 2016. 'Fast principal-component analysis reveals convergent evolution of ADH1B in Europe and East Asia', *The American Journal of Human Genetics*, 98: 456-72.
- Galinsky, Kevin, Po-Ru Loh, Swapn Mallick, Nick J Patterson, and Alkes L Price. 2016. 'Population structure of UK Biobank and ancient Eurasians reveals adaptation at genes influencing blood pressure', *bioRxiv*: 055855.

- Gazave, Elodie, Li Ma, Diana Chang, Alex Coventry, Feng Gao, Donna Muzny, Eric Boerwinkle, Richard A Gibbs, Charles F Sing, and Andrew G Clark. 2014. 'Neutral genomic regions refine models of recent rapid human population growth', *Proceedings of the National Academy of Sciences*, 111: 757-62.
- Gogarten, Stephanie M, Tushar Bhangale, Matthew P Conomos, Cecelia A Laurie, Caitlin P McHugh, Ian Painter, Xiuwen Zheng, David R Crosslin, David Levine, and Thomas Lumley. 2012. 'GWASTools: an R/Bioconductor package for quality control and analysis of Genome-Wide Association Studies', *Bioinformatics*, 28: 3329-31.
- He, Chunyan, and Joanne M Murabito. 2014. 'Genome-wide association studies of age at menarche and age at natural menopause', *Molecular and Cellular Endocrinology*, 382: 767-79.
- Howie, Bryan N, Peter Donnelly, and Jonathan Marchini. 2009. 'A flexible and accurate genotype imputation method for the next generation of genome-wide association studies', *PLoS Genetics*, 5: e1000529.
- Jones, Eric, Travis Oliphant, and Pearu Peterson. 2001. 'SciPy: Open source scientific tools for Python'. <http://www.scipy.org/>.
- Joshi, Peter K, Krista Fischer, Katharina E Schraut, Harry Campbell, Tõnu Esko, and James F Wilson. 2016. 'Variants near CHRNA3/5 and APOE have age-and sex-related effects on human lifespan', *Nature communications*, 7.
- KAar, Pekka, Jukka Jokela, Timo Helle, and Ilpo Kojola. 1996. 'Direct and correlative phenotypic selection on life-history traits in three pre-industrial human populations', *Proceedings of the Royal Society of London B: Biological Sciences*, 263: 1475-80.
- Karn, Mary N, and LS Penrose. 1951. 'Birth weight and gestation time in relation to maternal age, parity and infant survival', *Annals of Eugenics*, 16: 147-64.
- Kvale, Mark N, Stephanie Hesselton, Thomas J Hoffmann, Yang Cao, David Chan, Sheryl Connell, Lisa A Croen, Brad P Dispensa, Jasmin Eshragh, and Andrea Finn. 2015. 'Genotyping informatics and quality control for 100,000 subjects in the Genetic Epidemiology Research on Adult Health and Aging (GERA) Cohort', *Genetics*, 200: 1051-60.
- Liu, Chia-Chan, Takahisa Kanekiyo, Huaxi Xu, and Guojun Bu. 2013. 'Apolipoprotein E and Alzheimer disease: risk, mechanisms and therapy', *Nature Reviews Neurology*, 9: 106-18.
- Liu, Jimmy Z, Yaniv Erlich, and Joseph K Pickrell. 2016. 'Case-control association mapping without cases', *bioRxiv*: 045831.
- Loh, Po-Ru, Pier Francesco Palamara, and Alkes L Price. 2016. 'Fast and accurate long-range phasing in a UK Biobank cohort', *Nature Genetics*.
- Marchini, Jonathan, Bryan Howie, Simon Myers, Gil McVean, and Peter Donnelly. 2007. 'A new multipoint method for genome-wide association studies by imputation of genotypes', *Nature Genetics*, 39: 906-13.
- Marioni, Riccardo E, Stuart J Ritchie, Peter K Joshi, Saskia P Hagenaars, Aysu Okbay, Krista Fischer, Mark J Adams, W David Hill, Gail Davies, and Reka Nagy. 2016. 'Genetic variants linked to education predict longevity', *Proceedings of the National Academy of Sciences*: 201605334.
- Mathieson, Iain, Iosif Lazaridis, Nadin Rohland, Swapna Mallick, Nick Patterson, Songül Alpaslan Roodenberg, Eadaoin Harney, Kristin Stewardson, Daniel Fernandes, and Mario Novak.

2015. 'Genome-wide patterns of selection in 230 ancient Eurasians', *Nature*, 528: 499-503.
- Maynard Smith, John, and John Haigh. 1974. 'The hitch-hiking effect of a favourable gene', *Genetical Research*, 23: 23-35.
- Meyer, Wynn K, Barbara Arbeithuber, Carole Ober, Thomas Ebner, Irene Tiemann-Boege, Richard R Hudson, and Molly Przeworski. 2012. 'Evaluating the evidence for transmission distortion in human pedigrees', *Genetics*, 191: 215-32.
- Milot, Emmanuel, Francine M Mayer, Daniel H Nussey, Mireille Boisvert, Fanie Pelletier, and Denis Réale. 2011. 'Evidence for evolution in response to natural selection in a contemporary human population', *Proceedings of the National Academy of Sciences*, 108: 17040-45.
- Murabito, Joanne M, Rong Yuan, and Kathryn L Lunetta. 2012. 'The search for longevity and healthy aging genes: insights from epidemiological studies and samples of long-lived individuals', *The Journals of Gerontology Series A: Biological Sciences and Medical Sciences*, 67: 470-79.
- Nielsen, Rasmus, Carlos Bustamante, Andrew G Clark, Stephen Glanowski, Timothy B Sackton, Melissa J Hubisz, Adi Fladel-Alon, David M Tanenbaum, Daniel Civello, and Thomas J White. 2005. 'A scan for positively selected genes in the genomes of humans and chimpanzees', *PLoS Biology*, 3: e170.
- Nielsen, Rasmus, Ines Hellmann, Melissa Hubisz, Carlos Bustamante, and Andrew G Clark. 2007. 'Recent and ongoing selection in the human genome', *Nature Reviews Genetics*, 8: 857-68.
- Pe'er, Itsik, Roman Yelensky, David Altshuler, and Mark J Daly. 2008. 'Estimation of the multiple testing burden for genomewide association studies of nearly all common variants', *Genetic Epidemiology*, 32: 381-85.
- Pennings, Pleuni S, and Joachim Hermisson. 2006. 'Soft sweeps III: the signature of positive selection from recurrent mutation', *PLoS Genetics*, 2: e186.
- Perry, George H, Nathaniel J Dominy, Katrina G Claw, Arthur S Lee, Heike Fiegler, Richard Redon, John Werner, Fernando A Villanea, Joanna L Mountain, and Rajeev Misra. 2007. 'Diet and the evolution of human amylase gene copy number variation', *Nature Genetics*, 39: 1256-60.
- Peto, Richard, Sarah Darby, Harz Deo, Paul Silcocks, Elise Whitley, and Richard Doll. 2000. 'Smoking, smoking cessation, and lung cancer in the UK since 1950: combination of national statistics with two case-control studies', *BMJ*, 321: 323-29.
- Pickrell, Joseph K, Tomaz Berisa, Jimmy Z Liu, Laure Ségurel, Joyce Y Tung, and David A Hinds. 2016. 'Detection and interpretation of shared genetic influences on 42 human traits', *Nature Genetics*.
- Pierce, John P, Karen Messer, Martha M White, Sheila Kealey, and David W Cowling. 2010. 'Forty years of faster decline in cigarette smoking in California explains current lower lung cancer rates', *Cancer Epidemiology Biomarkers & Prevention*, 19: 2801-10.
- Pilling, Luke C, Janice L Atkins, Kirsty Bowman, Samuel E Jones, Jessica Tyrrell, Robin N Beaumont, Katherine S Ruth, Marcus A Tuke, Hanieh Yaghootkar, and Andrew R Wood. 2016. 'Human longevity is influenced by many genetic variants: evidence from 75,000 UK Biobank participants', *Aging (Albany NY)*, 8: 547.

- Price, Alkes L, Nick J Patterson, Robert M Plenge, Michael E Weinblatt, Nancy A Shadick, and David Reich. 2006. 'Principal components analysis corrects for stratification in genome-wide association studies', *Nature Genetics*, 38: 904-09.
- Przeworski, Molly, Graham Coop, and Jeffrey D Wall. 2005. 'The signature of positive selection on standing genetic variation', *Evolution*, 59: 2312-23.
- Raichlen, David A, and Gene E Alexander. 2014. 'Exercise, APOE genotype, and the evolution of the human lifespan', *Trends in Neurosciences*, 37: 247-55.
- Robinson, Matthew R, Gibran Hemani, Carolina Medina-Gomez, Massimo Mezzavilla, Tonu Esko, Konstantin Shakhbazov, Joseph E Powell, Anna Vinkhuyzen, Sonja I Berndt, and Stefan Gustafsson. 2015. 'Population genetic differentiation of height and body mass index across Europe', *Nature Genetics*, 47: 1357-62.
- Roden, Dan M, Jill M Pulley, Melissa A Basford, Gordon R Bernard, Ellen W Clayton, Jeffrey R Balser, and Dan R Masys. 2008. 'Development of a large-scale de-identified DNA biobank to enable personalized medicine', *Clinical Pharmacology and Therapeutics*, 84: 362-69.
- Sabeti, Pardis C, Stephen F Schaffner, Ben Fry, Jason Lohmueller, Patrick Varilly, Oleg Shamovsky, Alejandro Palma, TS Mikkelsen, D Altshuler, and ES Lander. 2006. 'Positive natural selection in the human lineage', *Science*, 312: 1614-20.
- Stearns, Stephen C, Sean G Byars, Diddahally R Govindaraju, and Douglas Ewbank. 2010. 'Measuring selection in contemporary human populations', *Nature Reviews Genetics*, 11: 611-22.
- Sudlow, Cathie, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul Downey, Paul Elliott, Jane Green, and Martin Landray. 2015. 'UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age', *PLoS Medicine*, 12: e1001779.
- Teshima, Kosuke M, Graham Coop, and Molly Przeworski. 2006. 'How reliable are empirical genomic scans for selective sweeps?', *Genome Research*, 16: 702-12.
- The 1000 Genomes Project Consortium. 2015. 'A global reference for human genetic variation', *Nature*, 526: 68-74.
- Tishkoff, Sarah A, Floyd A Reed, Alessia Ranciaro, Benjamin F Voight, Courtney C Babbitt, Jesse S Silverman, Kweli Powell, Holly M Mortensen, Jibril B Hirbo, and Maha Osman. 2007. 'Convergent adaptation of human lactase persistence in Africa and Europe', *Nature Genetics*, 39: 31-40.
- Tobacco and Genetics Consortium. 2010. 'Genome-wide meta-analyses identify multiple loci associated with smoking behavior', *Nature Genetics*, 42: 441-47.
- Tropf, Felix C, Gert Stulp, Nicola Barban, Peter M Visscher, Jian Yang, Harold Snieder, and Melinda C Mills. 2015. 'Human fertility, molecular genetics, and natural selection in modern societies', *PloS One*, 10: e0126821.
- Turchin, Michael C, Charleston WK Chiang, Cameron D Palmer, Sriram Sankararaman, David Reich, Joel N Hirschhorn, and Genetic Investigation of ANthropometric Traits Consortium. 2012. 'Evidence of widespread selection on standing variation in Europe at height-associated SNPs', *Nature Genetics*, 44: 1015-19.
- Turner, Stephen D. 2014. 'qqman: an R package for visualizing GWAS results using QQ and manhattan plots', *bioRxiv*: 005165.
- UK Biobank. <http://www.ukbiobank.ac.uk/>.

- . 2015a. 'Genotype imputation and genetic association studies of UK Biobank, Interim Data Release'. http://www.ukbiobank.ac.uk/wp-content/uploads/2014/04/imputation_documentation_May2015.pdf.
- . 2015b. 'Genotyping and quality control of UK Biobank, a large-scale, extensively phenotyped prospective resource'. http://www.ukbiobank.ac.uk/wp-content/uploads/2014/04/UKBiobank_genotyping_QC_documentation-web.pdf.
- Visscher, Peter M, Matthew A Brown, Mark I McCarthy, and Jian Yang. 2012. 'Five years of GWAS discovery', *The American Journal of Human Genetics*, 90: 7-24.
- Voight, Benjamin F, Sridhar Kudaravalli, Xiaoquan Wen, and Jonathan K Pritchard. 2006. 'A map of recent positive selection in the human genome', *PLoS Biology*, 4: e72.
- Williams, George C. 1957. 'Pleiotropy, natural selection, and the evolution of senescence', *Evolution*, 11: 13.
- Williamson, Scott H, Melissa J Hubisz, Andrew G Clark, Bret A Payseur, Carlos D Bustamante, and Rasmus Nielsen. 2007. 'Localizing recent adaptive evolution in the human genome', *PLoS Genetics*, 3: e90.
- Yang, Ziheng, and Joseph P Bielawski. 2000. 'Statistical methods for detecting molecular adaptation', *Trends in Ecology & Evolution*, 15: 496-503.