

An efficient and accurate frailty model approach for genome-wide survival association analysis controlling for population structure and relatedness in large-scale biobanks

Rounak Dey^{1‡}, Wei Zhou^{2,3,4‡}, Tuomo Kiiskinen^{5,6}, Aki Havulinna^{5,6}, Amanda Elliott^{1,2,3}, Juha Karjalainen^{2, 3,4,5}, Mitja Kurki^{2,3,4,5}, Ashley Qin¹, FinnGen, Seunggeun Lee⁷, Aarno Palotie^{2,3,4,5}, Benjamin Neale^{2,3,4*}, Mark Daly^{2,3,4,5*}, Xihong Lin^{1,3,8*}

¹Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA 02115, USA

²Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, Massachusetts, USA;

³Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge, Massachusetts, USA;

⁴Stanley Center for Psychiatric Research, Broad Institute of Harvard and MIT, Cambridge, Massachusetts, USA;

⁵Institute for Molecular Medicine Finland, Helsinki Institute of Life Sciences, University of Helsinki, Helsinki, Finland

⁶ Finnish Institute for Health and Welfare, Helsinki, Finland.

⁷Graduate School of Data Science, Seoul National University, Seoul, Korea

⁸Department of Statistics, Harvard University, Cambridge, Massachusetts, USA

‡These authors contributed equally to this work

*These authors jointly supervised this work

Abstract

With decades of electronic health records linked to genetic data, large biobanks provide unprecedented opportunities for systematically understanding the genetics of the natural history of complex diseases. Genome-wide survival association analysis can identify genetic variants associated with ages of onset, disease progression and lifespan. We developed an efficient and accurate frailty (random effects) model approach for genome-wide survival association analysis of censored time-to-event (TTE) phenotypes in large biobanks by accounting for both population structure and relatedness. Our method utilizes state-of-the-art optimization strategies to reduce the computational cost. The saddlepoint approximation is used to allow for analysis of heavily censored phenotypes (>90%) and low frequency variants (down to minor allele count 20). We demonstrated the performance of our method through extensive simulation studies and analysis of five TTE phenotypes, including lifespan, with heavy censoring rates (90.9% to 99.8%) on ~400,000 UK Biobank participants with white British ancestry and ~180,000 samples in FinnGen, respectively. We further performed genome-wide association analysis for 871 TTE phenotypes in UK Biobank and presented the genome-wide scale phenome-wide association (PheWAS) results with the PheWeb browser.

Introduction

Survival models, especially the Cox proportional hazard model¹, have been widely used to analyze time-to-event (TTE) outcomes, both in biomedical research²⁻⁴, and in genome-wide association studies (GWAS)⁵⁻¹¹. It has been shown that the proportional hazard model can increase the power to detect genetic variants associated with the age-of-onset of TTE phenotypes in cohort studies compared to modelling the disease status using a logistic regression model¹²⁻¹⁴. With the availability of detailed time-stamped diagnosis data from Electronic Health Records (EHR), large biobanks, such as UK Biobank (UKBB)¹⁵ (> 400,000 individuals) and FinnGen (<https://www.finnngen.fi/en>) (> 200,000 individuals), provide unprecedented opportunities to analyze TTE phenotypes to unravel the complex genetic architectures of disease onset, progression, and lifespan. Genome-wide scans of TTE phenotypes in large biobanks can potentially identify novel genetic variants associated with the onset of human diseases by leveraging both the disease status and the age-of-onset information.

In GWAS analysis, population structure and sample relatedness are often key factors that need to be controlled for. Biobank cohorts often have substantial population structure and relatedness. For example, in the UK Biobank, 91,392 out of 408,582 subjects with White British ancestry have at least one relative (up to 3rd degree) in the data. Several linear¹⁶⁻¹⁸ and logistic^{19,20} mixed effects models have been developed to account for relatedness in GWASs for quantitative and binary phenotypes. To account for related subjects in the proportional hazard model, frailty models, which are mixed effects survival models, have been proposed^{21,22}, where event times are assumed to be independent conditional on unobserved random effects called “frailties”. The frailties are modeled based on the dependence and clustering structure of the observations. Previous research has extensively studied shared frailty models with Gamma-distributed frailties^{21,23-27}. However, the shared frailty model is limited in its scope to model more complicated dependency structures that arise in cohort-based association studies. To model complicated dependency structures, such as known familial structures and cryptic relatedness, the multivariate frailty model with Gaussian frailty was proposed^{28,29}, and was later implemented in the R package COXME³⁰, which, however, lacks scalability for GWASs. Recently the COXME method was further

improved in COXMEG³¹, which utilizes several computational optimization strategies to make it applicable in genetic association studies, but COXMEG still cannot handle biobank-scale genome-wide datasets. Based on our performance benchmarking, even for 20,000 subjects, COXMEG requires 3,356 CPU-hours to perform a GWAS of 46 million variants, which means even with perfect parallelization on 30 CPUs, it would take over 4.6 days to complete the GWAS.

In large-scale GWASs, the score test is particularly useful among different asymptotic tests, because it requires fitting the model only once under the null hypothesis of no association²⁰. Score tests have also been implemented in the COXMEG package³². However, score tests can lead to severe type I error inflation for phenotypes with heavy censoring, which is extremely common in biobank-based phenotypes. In the UK Biobank phenome that we built (see **ONLINE METHODS**), 871 TTE phenotypes have at least 500 events (cases), out of which 811 phenotypes have censoring rate more than 95%. The inaccuracies of the score test in unbalanced case-control phenotypes have been previously shown for logistic regression and logistic mixed effects models^{19,33-35}, and a saddlepoint approximation³⁶ (SPA)-based adjustment has been proposed and successfully implemented¹⁹ to accurately calibrate the p-values in such scenarios. Recently, the SPACox¹¹ method also used SPA to calibrate p-values for time-to-event phenotypes in unrelated samples. However, the SPACox method does not account for sample-relatedness. Through simulations, we show similar inaccuracies are also present in score tests in frailty models for analyzing heavily censored phenotypes.

Here we propose a novel method for genome-wide survival analysis of TTE phenotypes, which accounts for both population structure and sample relatedness, controls type I error rates even for phenotypes with extremely heavy censoring, and is scalable for genome-wide scale PheWASs on biobank-scale data. Our method, Genetic Analysis of Time-to-Event phenotypes (GATE), transforms the likelihood of a multivariate Gaussian frailty model to a modified Poisson generalized linear mixed model (GLMM^{20,37}) likelihood, employs several state-of-the-art optimization techniques to fit the modified GLMM under the null hypothesis, and then performs score tests calculated using the null model for each genetic variant. To obtain well-calibrated p-values for heavily censored phenotypes, GATE uses the SPA to estimate the null distribution of the score

statistic instead of the traditionally used normal approximation. Moreover, our method saves the memory requirement substantially by storing the raw genotypes in binary format and calculating the elements of the GRM on the fly instead of storing or inverting a large dimensional GRM.

Through extensive simulations and analysis of TTE phenotypes from the UK Biobank data of 408,582 subjects with White British ancestry and the FinnGen study, we showed that GATE is scalable to biobank-scale GWASs of TTE phenotypes with type I error rates well controlled even for less frequent variants and heavily censored phenotypes. Benchmarking has shown that GATE can analyze 46 million variants in a GWAS with 408,582 subjects in ~ 14.5 hours using 30 CPUs with peak memory usage under 11 GB.

Results

Overview of Methods

GATE consists of two main steps: 1) Fitting the null frailty model to estimate the variance component and other model parameters, and 2) performing a score statistic-based test for association between each genetic variant and the phenotype. Step 1 involves iteratively fitting the null frailty model using similar optimization strategies as described in GMMAT²⁰ and SAIGE¹⁹, such as using the computationally efficient average information restricted maximum likelihood (AI-REML^{20,38}) algorithm for estimating the variance component, and using pre-conditioned gradient descent (PCG³⁹) method to solve linear systems to avoid inverting the $N \times N$ genetic relatedness matrix (GRM). GATE computes the elements of the GRM on-the-fly when needed using binary vectors of raw genotypes, and thus it doesn't require to supply, store, or invert a pre-computed GRM, which can be extremely time and memory-consuming for large sample sizes (N). For example, in UK Biobank data with $M = 93,511$ markers and $N = 408,582$ subjects with White British ancestry, the memory requirement drops from 622 GB for storing a pre-computed GRM in floating point numbers, to only 8.9 GB for storing the raw genotypes in the binary format.

Step 2 involves scanning the entire genome and testing each variant for association using the score statistic. Since the overall cost of computing the variance of the score statistic for all variants is extremely high because it involves operations on the large-

dimensional GRM, in step 2, GATE uses a variance ratio approximation commonly used in existing LMM and GLMM-based methods such as GRAMMAR-Gamma¹⁷, BOLT-LMM¹⁶, fastGWA¹⁸, and SAIGE¹⁹. The ratio of the variance of the score statistic with and without the random effects (and an attenuation factor due to estimating the baseline hazards) is computed using a subset of genetic markers. Previously, it was shown that this variance ratio remains approximately constant for variants with $MAF \geq 20$ for LMM and GLMMs. Through analytical derivations and simulation examples, we show this observation to hold for frailty models as well (**Supplementary Note section 3 and Supplementary Figure 14**). Therefore, when performing the genome-wide scan, the variance of the score statistic is computed without using the GRM and then calibrated using the variance ratio.

Next, GATE uses the saddlepoint approximation³⁶ (SPA) to approximate the null distribution of score statistics for association tests. SPA-based tests have been successfully used for logistic regression³⁴ and logistic mixed models¹⁹ and provide more accurate p-values than traditional score tests under normal approximation for low-frequency variants when the case-control ratio is unbalanced. In GATE, we have implemented an efficient SPA-based test for frailty models that is similar to the fastSPA method in Dey et al.³⁴. Through simulations and real data analysis, we show that SPA tests provide accurate and calibrated p-values, even for low-frequency variants when the censoring rate is high to 99%.

Both GATE and COXMEG³¹ conduct genetic association tests for TTE phenotypes using the frailty model. Besides the use of SPA-based tests, GATE uses the variance ratio approach to approximate the variances of the score statistics, while COXMEG calculates the variances using the GRM. Using simulation studies, we have shown that GATE provides consistent association p-values to COXMEG (R^2 of $-\log_{10}$ P-values > 0.99) for common variants ($MAF > 5\%$) when the censoring rate is 50% (**Supplementary Figure 1A**) and has well controlled type I error rates, even for less frequent variants and phenotypes with heavy censoring rates (**Supplementary Figure 1B**).

Computation and Memory costs

To assess the computational performance of GATE and the score test implemented in the COXMEG package (COXMEG-Score), we randomly sampled subsets of different sample sizes from 408,582 UK Biobank subjects with White British ancestry. We then benchmarked association tests for overall lifespan (16,375 events, 389,721 censored) adjusting for the top four ancestry principal components, birth year and sex using GATE and COXMEG-Score on 200,000 variants randomly selected from 46 million genetic variants with imputation info ≥ 0.3 and MAC ≥ 20 . In Step 1, 93,511 high-quality genotyped markers were used for the GRM. The projected overall computation time (**Figure 1 and Supplementary Table 1**) for GATE to analyze 46 million variants on $N = 408,582$ subjects was 318 CPU-hours, and the actual computation time on a machine with 30 cores was 14.5 hours. Step 2, which accounts for the majority of the computation time (95.4% for $N = 408,582$) requires substantially less memory (peak memory usage 0.85 GB) than Step 1 (peak memory usage 10.6 GB). However, even for 20,000 subjects, the projected computation time and memory usage for COXMEG-Score were 3,356 CPU-hours and 32.75 GB, compared to only 34 CPU-hours and 0.74 GB required by GATE, achieving 99% and 97.7% reductions in computation time and memory, respectively. This means even with perfect parallelization on 30 CPUs, COXMEG-Score would require 4.6 days to complete the GWAS with only 20,000 subjects. The observations also suggest that the computation time and memory requirements increase nearly linearly with the sample size for GATE, whereas they increase quadratically for COXMEG-Score.

Phenome-wide GWAS of time-to-event phenotypes in the UK Biobank data.

We have applied GATE to perform phenome-wide GWAS for 871 UKBB TTE phenotypes with at least 500 events, adjusting for top four PCs, birth year, and sex (except for 93 sex-specific phenotypes). The TTE phenotypes were created based on the International Classification of Disease (ICD) codes version 9 and 10 mapped to the PheWAS code (PheCode⁴⁰) definitions (See **ONLINE METHODS**) as well as their associated diagnosis dates in the UK Biobank electronic medical records. For each phenotype, we analyzed approximately 46 million genetic markers imputed from the

Haplotype Reference Consortium⁴¹ panel and UK10K⁴² with imputation INFO score ≥ 0.3 and MAC ≥ 20 . Among the 408,582 UK Biobank subjects with White British ancestry, 91,392 had at least one relative up to third degree¹⁵. To account for the relatedness among the subjects, we used 93,511 high-quality genotyped markers with MAF ≥ 0.01 to construct the GRM in Step 1. The same set of markers were used by the UK Biobank research group¹⁵ for estimating kinship among the samples because they are only weakly informative of the ancestry and therefore provide more accurate kinship estimates. We also performed a sensitivity analysis using a larger set of markers (245,745) for the four exemplary phenotypes discussed before (See **Supplementary Note Section 7**). We further applied SPA-based adjustment of the score test because to the censoring rates (**Supplementary Figure 2**) were extremely high for most of the TTE phenotypes in the UKBB (for example, 811 out of 871 have censoring rate more than 95%). The summary statistics for all 871 PheCodes analyzed using GATE are available to download from a public repository (see URL) and browsed in the PheWeb⁴³ (see URL).

Here we discuss the association results using four phenotypes with different censoring rates as exemplars: ischemic heart disease (IHD: PheCode 411, N events=36,962, N censored=370,814, censoring rate=90.9%), female breast cancer (FBC, PheCode 174.1, N events=15,396, N censored=192,764, censoring rate=92.6%), glaucoma (PheCode 365, N events=6,046, N censored=392,925, censoring rate=98.5%), and Alzheimer's Disease (AD: PheCode 290.11, N events=822, N censored=342,059, censoring rate=99.8%). The Manhattan and QQ plots for the GWAS of these phenotypes using GATE with and without SPA are presented in **Figure 2** and **Figure 3**, respectively. The results demonstrate that not adjusting for SPA greatly inflates the type I errors, especially for the low frequency variants, whereas the SPA-adjusted method shows well controlled type I error rates. In total, 114 loci have been identified for the four TTE phenotypes: 55 for IHD, 37 for FBC, 19 for glaucoma, and 3 for AD. We also applied GATE to these four phenotypes in the FinnGen study (see **ONLINE METHODS**) and 81 out of the 114 loci were also tested in the FinnGen study, of which 78 had the same effect direction in both UKBB and FinnGen. 69 out of the 81 loci were successfully replicated in FinnGen with p-value < 0.05 . The complete list of all significant loci and the

association results in the UKBB, FinnGen as well as the meta-analysis of the two data sets are reported in **Supplementary Table 2**. Overall, 99 out of the 114 significant loci have been previously reported to be associated with disease risk in case-control studies to the best of our knowledge. Several loci that are previously well known as associated with the risk of the diseases have been identified in our study. For example, the loci *LPA* and *CELSR2* for IHD^{44,45}, *FGFR2*⁴⁶ and *CASC16*⁴⁷ for breast cancer, *MYOC*⁴⁸ and *TMCO1*⁴⁹ for glaucoma, and *APOE* e4 variant for AD⁵⁰. The age-varying predicted risk of disease onset based on the GATE method, and the age-varying disease-free probability by genotypes based on the Kaplan-Meier curve⁵¹ for the exemplary top hits were plotted in **Figure 4** and **Supplementary Figure 3**, respectively.

GWAS of lifespan in the FinnGen Study and the UK Biobank

We have also applied GATE to the overall lifespan in the FinnGen study (N events = 15,152, N censored = 203,244), in which the age of death ranges from 7 years old to 106 years old as shown in **Supplementary Figure 4**. We identified the previously reported *APOE* locus for lifespan⁵² in FinnGen, in which the most significant variant is the *APOE*-e4 missense variant rs429358 (MAF = 18.3%, p-value = 1.01×10^{-14}) and it is well-known to be associated with lifespan, cardiovascular diseases, stroke, and Alzheimer's disease⁵³⁻⁵⁵. This locus has been replicated in UKBB (N events = 16,375 and N censored = 389,721, see **Supplementary Figure 5**) with p-value 1.92×10^{-5} and meta-analysis p-value 4.04×10^{-17} (**Supplementary Table 3** and **Supplementary Figure 6**). The top hit in UKBB (rs157592, MAF = 18.7%, p-value = 1.87×10^{-8}) had LD $r^2 = 0.7$ with rs429358 as presented in the **Supplementary Table 3**. This variant is in the intergenic region and have no in-silico functions according to the FAVOR functional annotation online portal⁵⁶ (See URL).

Simulation Studies

We investigated the type I error rates and power of GATE in presence of sample relatedness using 10,000 simulated samples. Due to computation burden, we used GATE-noSPA instead of COXMEG-Score for type I error evaluation as **Supplementary**

Figure 1C shows the two approaches provide consistent association p-values (R^2 of -log₁₀ p-values > 0.99).

The type I error rates of GATE were evaluated based on association tests of 9.4×10^8 simulated genetic markers on 10,000 samples, which contain 500 families and 5,000 independent samples. Each family has 10 members, simulated based on the pedigree shown in **Supplementary Figure 7**. The variance component parameter τ is set to be 0.1 and 0.25 (see **ONLINE METHODS**). The empirical type I error rates at the significance level $\alpha = 1 \times 10^{-6}$ and 5×10^{-8} are shown in the **Supplementary Table 4 and Supplementary Figure 8A**. Our simulation results suggest that GATE has well controlled type I error rates even for low frequency variants (down to MAC = 20) when the phenotype is heavily censored (90%). However, without SPA, the score tests in GATE suffer from inflated type I error rates as the case-control ratios become more unbalanced and the frequency of variants decreases. We also evaluated type I error rates of GATE in a setting with cryptic sample relatedness by randomly selecting 10,000 UKBB participants with white British ancestry. Phenotypes were simulated using the real genotypes to mimic the sample relatedness of a real-world dataset, and association tests were conducted on the imputed genetic markers in the UKBB (see **ONLINE METHODS**). Similarly, we observed that the type I error rates were well controlled in GATE in presence of cryptic sample relatedness with different censoring rates (**Supplementary Table 5, Supplementary Figure 8B and 9**).

Next, we evaluated empirical power of GATE at $\alpha = 5 \times 10^{-8}$ and compared to the power of COXMEG-Score. **Supplementary Figure 10** shows the power curve by hazard ratios for variants with MAF 0.05 and 0.2 when $\tau = 0.25$ and the censoring rate = 50%. Both methods have nearly identical power in all simulation settings. We do not compare their powers in the presence of heavy censoring, in view of the inflated type I error rate of COXMEG-Score.

Overall simulation studies show that GATE can control type I error rates even when censoring rate is high and has similar power for common variants as COXMEG-Score. In contrast, same as GATE-noSPA, COXMEG suffers type I error inflation and the

inflation is especially severe with low MAF and heavy censoring (**Supplementary Figure 1B, 1C, 8 and 9**).

Discussion

In this paper, we have proposed a novel method to perform scalable genome-wide survival association analysis of censored TTE phenotypes in large biobanks using an efficient implementation of the frailty model. Our method can adjust for population structure and sample relatedness and provide accurate p-values even in extreme cases of very low frequency variants and heavily censored phenotypes (incidence rate < 0.1%). Applying this approach to the UK Biobank and the FinnGen study, we demonstrated that our method is scalable to the analysis of large biobank-scale datasets with > 400,000 subjects.

Biobanks with genetic data linked to EHR records/survey questionnaires provide unprecedented opportunities for genetic association studies on TTE phenotypes to identify genetic risk factors that affect the onset and progression of diseases. However, biobanks pose challenges to such analysis because of the high computational and memory cost required to handle large data sets with extensive population structure and relatedness. Moreover, current methods artificially inflate associations when heavily censored phenotypes (e.g., censoring rate > 75%) and low frequency variants (MAF < 1%) are involved. The proposed method, GATE performs a frailty model-based association analysis to account for both population structure and relatedness using score tests with SPA adjustment, which provides accurate p-values under heavy censoring. In addition, it implements several optimization techniques that were previously used in the context of linear and logistic mixed models in BOLT-LMM and SAIGE to make it computationally feasible to analyze large biobank cohorts. We have applied GATE to 871 TTE phenotypes in the UK Biobank data with White British ancestry, which were constructed based on PheCodes mapped to ICD codes and have at least 500 events. The genome-side summary statistics are available for public to download. We have also created a PheWeb⁴³ for users to explore and visualize the PheWAS results.

TTE phenotypes are particularly suited not only for studying disease onsets, but also for exploring other progression phenotypes such as times of surgery, recurrence, times of onset of secondary phenotypes after an initial diagnosis etc. Previously, the lack of scalable GWAS methods for TTE outcomes has hindered such investigations in massive scales. By facilitating large-scale GWAS of TTE phenotypes, GATE opens the door to such deeper investigations.

One consideration while analyzing TTE phenotypes is the appropriate choice of the unit of time. To assess the impact of time-units on the GWAS results, we performed sensitivity analysis using the event and censoring times rounded to the nearest 1 month, 3 months, 6 months and 12 month time-units for the four exemplary UK Biobank phenotypes presented in this paper, and compared the p-values across different time-units (**Supplementary Figure 11**). The p-values were very similar across the four time-units for all phenotypes, with more detailed time-units resulting in slightly more significant p-values.

For the selection of number of markers to construct the GRM, there is a trade-off between computation cost and the accuracy of adjusting the sample relatedness. Increasing the number of markers (M) included in the GRM linearly increases the computation time and memory requirement of step 1, whereas using too few markers may not be sufficient to capture the detailed familial and cryptic relatedness among the samples properly⁵⁷. For the UK Biobank data analysis, we used $M = 93,511$ LD pruned high-quality genotyped markers which were used by the UK Biobank research group for estimating kinship among the samples¹⁵. We performed a sensitivity analysis (see **Supplementary Note Section 7**) by increasing the number of markers to $M = 245,975$ pruned markers with $MAF \geq 0.01$. The results (**Supplementary Figure 12 and 13**) showed that the p-values were generally concordant, and the p-values using $M = 245,975$ markers were slightly larger than the p-values using $M = 93,511$ markers.

There are several limitations to GATE. First, similar to other mixed model methods for genetic association tests, the computation time required for the algorithms to converge in step 1 can vary among different phenotypes and study samples because of the difference in heritability and the extent of sample relatedness. Second, GATE uses a score statistic-based test without fitting the model under the alternate hypothesis, which

can be computationally inefficient. Therefore, it does not provide accurate estimates of hazard ratios for the genetic variants. Following a similar approach as in several other mixed model-based methods^{16,17,19,58}, GATE provides a hazard ratio estimate using the null model parameter estimates (see **Supplementary Note Section 5**). Third, the current implementation of GATE is targeted to perform single-variant association analysis, which can suffer from low power to detect associations in extremely rare variants. With whole genome and whole exome sequencing data available, a possible future extension of this method can allow for mask-based or region-based association tests to improve power for the rare variants^{56,59}. Finally, the current version of GATE does not incorporate left-truncated data, which may not be valid for early-onset phenotypes in biobanks with relatively older participants. For example, the median age of UK Biobank's participants is 59 years old and the earliest dates of health data available are around late 1990s, and assuming no left-censoring can reduce association power for early-onset diseases. The next work will extend GATE to allow for left-truncated phenotypes. In summary, we have proposed a scalable and accurate method, GATE, to perform genome-wide PheWAS of TTE phenotypes on large biobank cohorts accounting for population structure, sample relatedness and heavy censoring. We demonstrated that it is possible to efficiently analyze the current largest biobank (UK Biobank) of > 400,000 subjects using GATE. Our method facilitates biobank-based PheWAS of TTE phenotypes which ultimately contributes towards identifying genetic components that affect the onset and progression of complex diseases.

URLs

GATE is implemented as an open-source R package available at <https://github.com/weizhou0/GATE>. The GWAS results for 871 time-to-event phenotypes in UK Biobank using GATE are currently available for public download at <http://gate.genohub.org/>. Manhattan plots, Q-Q plots, and regional association plots for each TTE phenotype as well as the PheWAS plots can be browsed at <http://phewas.genohub.org/>. The FAVOR⁵⁶ portal is accessed through favor.genohub.org.

Acknowledgments

The FinnGen project is funded by two grants from Business Finland (HUS 4685/31/2016 and UH 4386/31/2016) and eleven industry partners (AbbVie Inc, AstraZeneca UK Ltd, Biogen MA Inc, Celgene Corporation, Celgene International II Sàrl, Genentech Inc, Merck Sharp & Dohme Corp, Pfizer Inc., GlaxoSmithKline, Sanofi, Maze Therapeutics Inc., Janssen Biotech Inc). Following biobanks are acknowledged for collecting the FinnGen project samples: Auria Biobank (www.auria.fi/biopankki), THL Biobank (www.thl.fi/biobank), Helsinki Biobank (www.helsinginbiopankki.fi), Biobank Borealis of Northern Finland (<https://www.ppshep.fi/Tutkimus-ja-opetus/Biopankki/Pages/Biobank-Borealis-briefly-in-English.aspx>), Finnish Clinical Biobank Tampere ([www.tays.fi/en-US/Research and development/Finnish Clinical Biobank Tampere](http://www.tays.fi/en-US/Research%20and%20development/Finnish%20Clinical%20Biobank%20Tampere)), Biobank of Eastern Finland (www.ita-suomenbiopankki.fi/en), Central Finland Biobank (www.ksshp.fi/fi-FI/Potilaalle/Biopankki), Finnish Red Cross Blood Service Biobank (www.veripalvelu.fi/verenluovutus/biopankkitoiminta) and Terveystalo Biobank (www.terveystalo.com/fi/Yritystietoa/Terveystalo-Biopankki/Biopankki/). All Finnish Biobanks are members of BBMRI.fi infrastructure (www.bbMRI.fi). This research has been conducted using the UK Biobank Resource under application number 52008. X.L. was supported by NCI R35-CA197449, P01-CA134294, U19-CA203654 and NHLBI R01-HL113338. B.M.N. was supported by NHGRI U01-HG009088-04S3 and NIMH R37-MH107649-06. R.D. was supported by NCI R35-CA197449. W.Z. was supported by an NIH T32 fellowship (Grant number: 1T32HG010464-01). A.P. was supported by the Academy of Finland Centre of Excellence in Complex Disease Genetics (Grant No. 312074)

Author Contributions

R.D., W.Z. X.L., B.M.N., and M.J.D. designed experiments. R.D. and W.Z. performed experiments. R.D. and W.Z. implemented the software with input from X.L., B.M.N., and

M.J.D. R.D. constructed phenotypes for UK Biobank data. R.D. and X.L. analyzed UK Biobank data. A.Q., R.D., and W.Z. created the PheWeb browser for UK Biobank results. W.Z., T.K., A.H., A.E., J.K., M.K., and A.P. analyzed data for the FinnGen study. Helpful advice was provided by S.L. R.D. and W.Z. wrote the manuscript with input from all co-authors.

Competing Financial Interests Statement

B.M.N. is on the scientific advisory board of Deep Genomics, and is a consultant for CAMP4 Therapeutics, Takeda and Biogen. X.L. is a consultant to AbbVie Pharmaceuticals and Verily Life Sciences. M.J.D. is a founder of Maze Therapeutics and on the scientific advisory board of BC Platforms.

Figures

Figure 1: Projected computation time (A) and memory usage (B) for GATE and COXMEG-Score as a function of sample size (N). The numerical data are provided in Supplementary Table 1.

Benchmarking was performed for the GWAS of lifespan based on randomly subsampled data from UK Biobank White British ancestry subjects. Association tests were performed on 200,000 randomly selected markers with imputation INFO ≥ 0.3 , with the filtering criteria of MAC ≥ 20 .

The computation times were projected for testing 46 million variants with INFO ≥ 0.3 and MAC ≥ 20 . The reported run times are medians of five runs, each with randomly sampled subjects with different randomization seeds. The x and y axes are plotted in log10 scale.

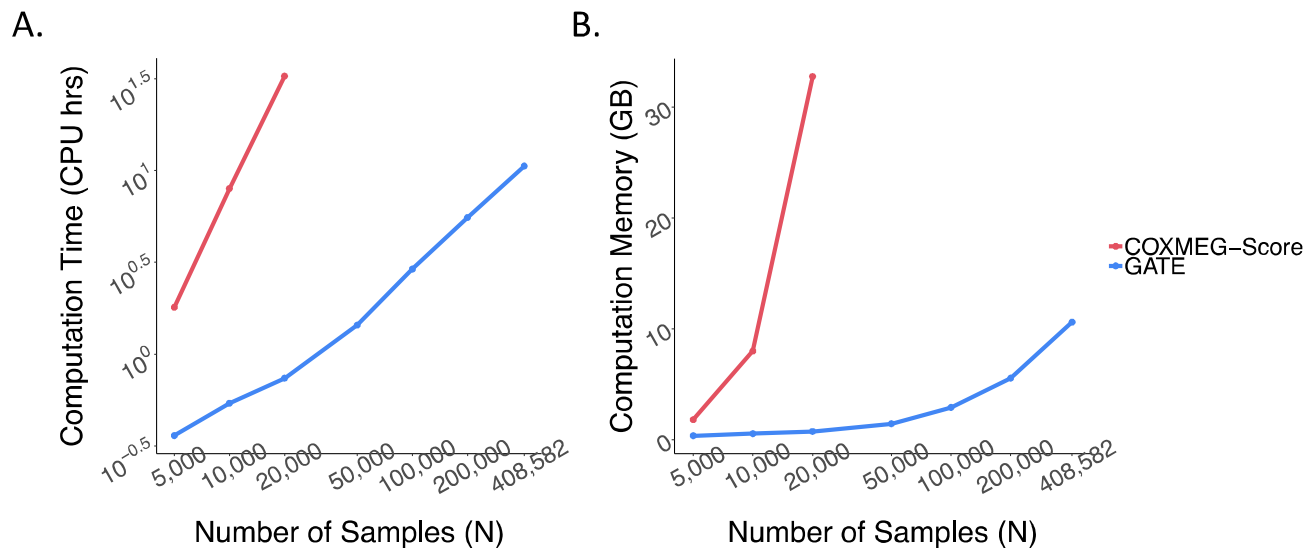
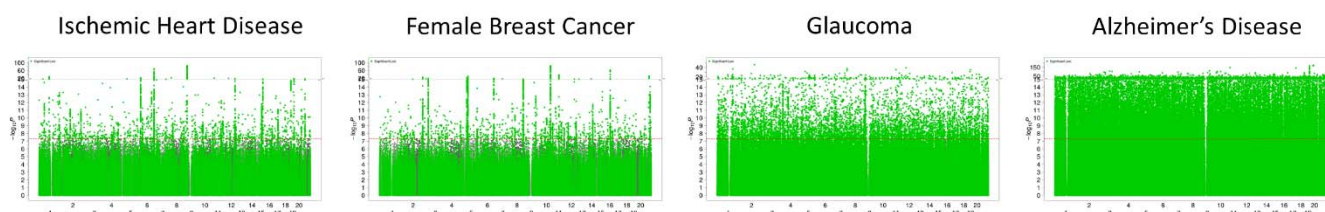


Figure 2: Manhattan plots for GWAS of four time-to-event phenotypes with different censoring rates in the UK Biobank data with White British ancestry: GWAS results using GATE-noSPA (A) and GATE (B) are shown for ischemic heart disease (PheCode 411, N=407776, censoring rate=90.9%), female breast Cancer (PheCode 174.1, N=208160, censoring rate=92.6%), glaucoma (PheCode 365, N=398971, censoring rate=98.5%), and Alzheimer's Disease (PheCode 290.11, N=342881, censoring rate=99.8%).

A. GATE noSPA



B. GATE

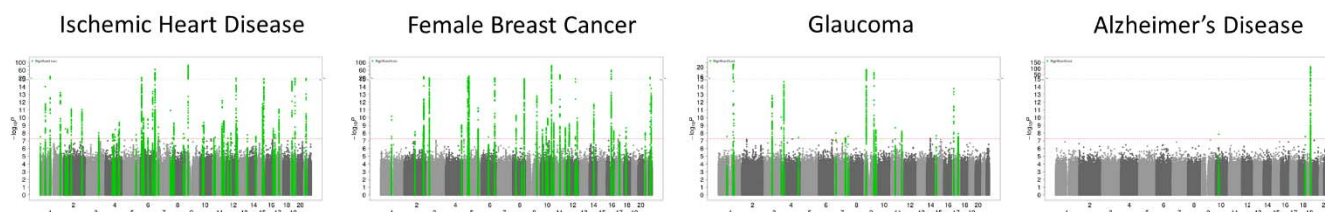
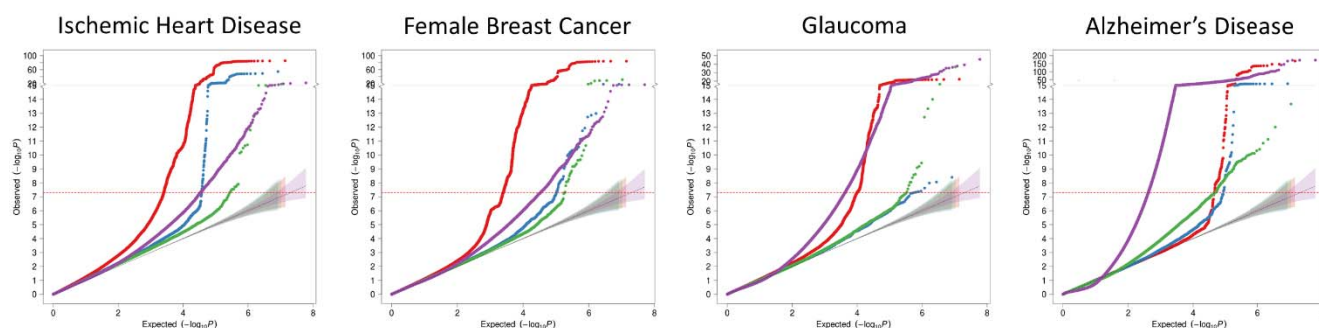
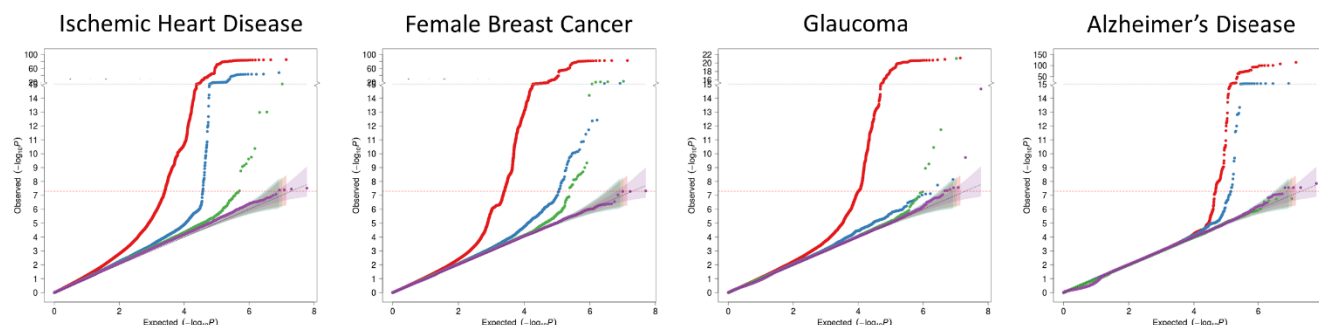


Figure 3: Quantile-quantile (QQ) plots for GWAS of four time-to-event phenotypes with different censoring rates in the UK Biobank data with White British ancestry: GWAS results using GATE-noSPA (A) and GATE (B) are shown for ischemic heart disease (PheCode 411, N=407776, censoring rate=90.9%), female breast Cancer (PheCode 174.1, N=208160, censoring rate=92.6%), glaucoma (PheCode 365, N=398971, censoring rate=98.5%), and Alzheimer's Disease (PheCode 290.11, N=342881, censoring rate=99.8%). QQ plots are color-coded based on different minor allele frequency categories.

A. GATE noSPA

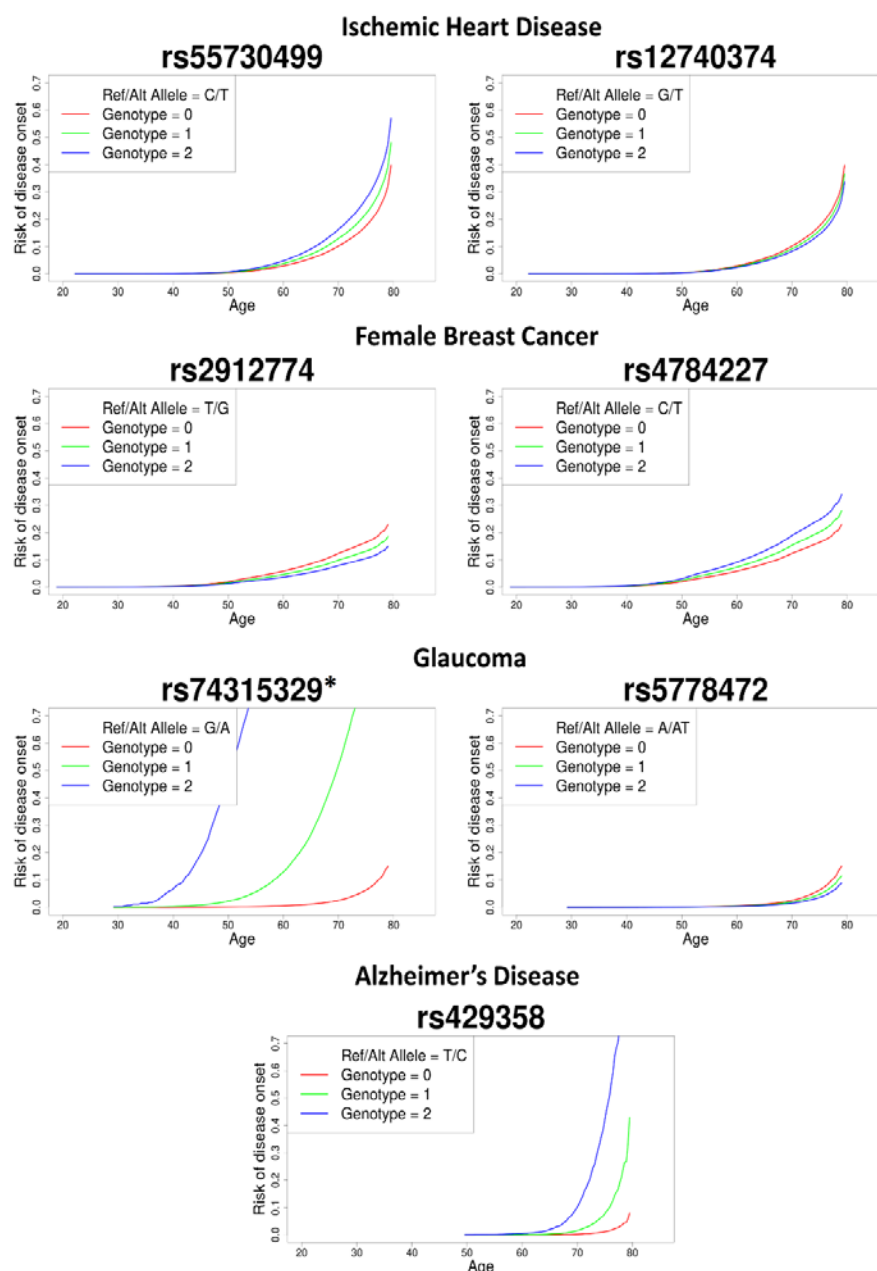


B. GATE



■ MAF=[0.05,0.5] ■ MAF=(0.005,0.05] ■ MAF=(0.001,0.005] ■ MAF=[4.8e-05,0.001]

Figure 4: Predicted risk of disease onset over-time by genotypes for loci LPA and CELSR2 for ischemic heart disease, FGFR2 and CASC16 for female breast cancer, MYOC and TMCO1 for glaucoma, and APOE e4 variant for AD. The red, green and blue lines represent the risk of disease onset for alternate allele counts zero, one and two, respectively for a female subject born in 1950 (median birth year in the UKBB data) with the top four PC coordinates each set at the mean level across the UK Biobank subjects with white British ancestry.



* No homozygous alternate subject was present among the study subjects for rs74315329 (Alternate allele frequency = 0.0013)

Online Methods

Frailty model for Time-to-event phenotypes.

Consider a study of N subjects, where for the i -th subject, we observe the data pair (δ_i, t_i) , where δ_i is a censoring indicator, with $\delta_i = 1$ if the i -th subject experiences an event during the study period, and $\delta_i = 0$ otherwise, i.e., censored. Let t_i denote the observed event or censoring time. For the i -th subject, let the $p \times 1$ vector X_i denote the covariates, and $G_i = 0, 1, 2$ denote the minor allele counts for the genetic variant of interest. Then, in a frailty model^{25,28,60}, the conditional hazard function of subject i at time t given the covariates, genotype and random effect/frailty b_i is modeled as

$$\lambda_i(t|b_i) = \lambda_0(t) \exp(X_i^T \beta + G_i \gamma + b_i),$$

where β and γ are the regression coefficients of the covariates X_i and the genotype G_i respectively, and $\lambda_0(t)$ is the baseline hazard function at time t , the frailty $b = (b_1, \dots, b_N)$ follows a multivariate normal distribution $N(0, \tau V)$, with V being the Genetic Related Matrix (GRM). Unlike standard generalized linear mixed models, the covariate vector X_i in a frailty model does not include the intercept term, instead the baseline hazard $\lambda_0(t)$ works as the intercept in a frailty model. We test the null hypothesis of no genetic association $H_0: \gamma = 0$ vs $H_1: \gamma \neq 0$.

Estimating the variance component and other null model parameters (step 1).

First, the likelihood for the observed event status-time pairs (δ_i, t_i) under the frailty model is derived and expressed as a modified Poisson mixed effects model likelihood, with the mean function weighted by the cumulative baseline hazard (CBH) function $\Lambda_0(t) = \int_0^t \lambda_0(u) du$. The CBH function is estimated by the Breslow's estimator $\hat{\Lambda}_0(t)$ as a step function. Breslow⁶¹ showed that the maximum likelihood approach for the proportional hazard model (for unrelated subjects) that leads to the estimator $\hat{\Lambda}_0(t)$, is equivalent to maximizing the partial likelihood proposed by Cox¹. In the **Supplementary Note Section 6**, we have shown that the same maximum likelihood approach holds for frailty models (related subjects) as well given the random effects. Then, using the penalized quasi-likelihood (PQL³⁷) method and the AI-REML³⁸ algorithm, the model

parameters under H_0 are estimated iteratively. To avoid storing large $N \times N$ GRMs, GATE only calculates the elements of the GRM when they are needed using raw binary format genotypes. For scalable computation of quantities of the form $A^{-1}x$ that arises in the model fitting steps, where A is a large matrix and x is a vector, GATE uses the PCG algorithm³⁹, which has been previously used in BOLT-LMM¹⁶ and SAIGE¹⁹ to accurately compute quantities like $y = A^{-1}x$ by solving the linear system of equations $Ay = x$, instead of explicitly inverting the large matrix A .

Once the null model parameters, random effects and cumulative baseline hazard functions $(\hat{\beta}, \hat{b}_i, \hat{\Lambda}_0(t_i))$ have been estimated, GATE estimates the variance ratio from a small number of markers. Denote the fitted means by $\hat{\mu}_i = \hat{\Lambda}_0(t_i) \exp(X_i^T \hat{\beta} + \hat{b}_i)$, and the weight matrix $\hat{W} = \text{diag}(\hat{\mu}_1, \dots, \hat{\mu}_N)$. Then the score statistic, under $H_0: \gamma = 0$ is $T = G^T(\delta - \hat{\mu}) = \tilde{G}^T(\delta - \hat{\mu})$, where $G = (G_1, \dots, G_N)$, $\delta = (\delta_1, \dots, \delta_N)$, $\hat{\mu} = (\hat{\mu}_1, \dots, \hat{\mu}_N)$. The covariate-and-intercept-adjusted genotypes are denoted by $\tilde{G} = G - \tilde{X}(\tilde{X}^T \hat{W} \tilde{X})^{-1} \tilde{X}^T G$, where $\tilde{X} = [1 \ X]$ is the augmented covariate matrix. Then, the variance of the score statistic under H_0 is given by $V_T = G^T \hat{Q} G = \tilde{G}^T \hat{Q} \tilde{G}$, where $\hat{Q} = \hat{S}^{-1} - \hat{S}^{-1} X (X^T \hat{S}^{-1} X)^{-1} X^T \hat{S}^{-1}$, $\hat{S} = (\hat{W} - \hat{U})^{-1} + \hat{t}V$. The expression of \hat{U} is described in detail in the **Supplementary Note Section 1.3**. Unlike in the GLMMs, the term \hat{U} appears in the variance of the score statistic due to the attenuation of information (additional variability) for estimating $\Lambda_0(t_i)$ s. The variance ratio is then calculated as $\hat{r} = \frac{\tilde{G}^T \hat{Q} \tilde{G}}{\tilde{G}^T \hat{W} \tilde{G}}$. GATE calculates the variance ratio based on 30 randomly selected genotyped markers with $\text{MAC} \geq 20$ and computes the coefficient of variation (CV). If the CV of the variance ratios is smaller than 0.001, then the mean of the variance ratios is selected as \hat{r} , otherwise more markers are selected at an increment of 10 markers, and the CV is recalculated until the CV becomes smaller than 0.001.

Score test using SPA.

Using the estimated variance ratio \hat{r} , the variance-adjusted test statistic can be calculated as $T_{adj} = \tilde{G}^T(\delta - \hat{\mu}) / \sqrt{\hat{r} \tilde{G}^T \hat{W} \tilde{G}}$, which under the null hypothesis has mean zero and variance unity. The traditional score test then assumes asymptotic normality of the score statistic T (and thus T_{adj} as well) under H_0 , to calculate the p-value. However,

observations have been made before in the context of logistic mixed models that the asymptotic normality assumption of the score test statistic leads to severe Type I error inflation for low-frequency and rare variants when the case-control ratio is unbalanced¹⁹. We make the same observations in frailty models as well when the censoring rate is high. In order to provide well calibrated p-values in such situations, we used saddle point approximation (SPA) to approximate the null distribution of the score statistic, which has been shown to have better approximation error bounds compared to the normal approximation^{34,36,62,63}, especially at the extremely small tail probability region of $\alpha = 5 \times 10^{-8}$. Contrary to the normal approximation which only utilizes the first two moments only to approximate, SPA utilizes the entire moment generating function (MGF). In fact, it uses the cumulant generating function (CGF), i.e., is the logarithm of the MGF, which for the frailty model, based on the modified Poisson mixed model likelihood, can be derived as $K(\xi) = \sum_{i=1}^N \hat{\mu}_i (e^{\tilde{G}_i c \xi} - \tilde{G}_i c \xi - 1)$, where $c = (\hat{r} \tilde{G}^\top \hat{W} \tilde{G})^{-1/2}$. Then, the distribution of T_{adj} can be calculated based on the SPA by $Pr(T_{adj} < s) \approx \Phi \left\{ w + \frac{1}{w} \log \left(\frac{v}{w} \right) \right\}$, and the p-value is given by $p = Pr(T_{adj} < -|s|) + Pr(T_{adj} > |s|)$, where $T_{adj} = s$ is the observed adjusted score statistic, $w = \text{sign}(\hat{\xi}) \sqrt{2(\hat{\xi}s - K(\hat{\xi}))}$, $v = \hat{t} \sqrt{K''(\hat{\xi})}$, $\hat{\xi}$ is the solution to the equation $K'(\hat{\xi}) = s$, and $K'(\xi)$ and $K''(\xi)$ are the first and second derivatives of the CGF $K(\xi)$, respectively.

Since the normal approximation works well around the mean, we use the normal approximation when T_{adj} is less than two standard deviations away from the mean for faster computation. In addition, a faster version of the SPA similar to Dey et al.³⁴ is also implemented which reduces the computation time even further, from $O(N)$ to $O(N_c)$, where N_c is the number of minor allele carriers.

Data Simulation.

We carried out a series of simulations to evaluate the performance of GATE, including the type I error rates and power. To evaluate whether GATE can control type I error rates in presence of sample relatedness, we randomly simulated a set of 1,000,000 base-pair “pseudo” sequences, in which variants are independent to each other. Alleles

for each variant were randomly drawn from Binomial($n = 2$, $p = \text{MAF}$). Then we performed the gene-dropping⁶⁴ simulation using these sequences as founder haplotypes that were propagated through the pedigree of 10 family members shown in **Supplementary Figure 7**. We simulated genotypes of 150,000 genetic variants with $\text{MAF} \geq 1\%$ for 5,000 independent samples and 500 families based on the pedigree to estimate the GRM on-the-fly in Step 1 of GATE and genotypes of 1.9 million genetic variants with $\text{MAC} \geq 20$ for association tests in Step 2. MAFs were randomly sampled from the MAF spectrum in UK Biobank imputation data as shown in **Supplementary Figure 9**. For each subject i , the censoring time T_{ci} was randomly selected from exponential distribution with mean $1/\lambda_c$ and the underlying failure time T_{fi} was generated from a frailty model with the underlying exponential hazard function $T_{fi} = \frac{-\log(U_i)}{\lambda \exp(\eta_i)}$, where $U_i \sim \text{uniform}(0,1)$ and η_i is the linear predictor. Under the null hypothesis of no genetic effects, $\eta_i = X_1^T \alpha + b_i$, where X_1 is a covariate that was randomly drawn from $N(0, 1)$, α is the coefficient and is 0.5 and b_i is the random effect simulated from $N(0, \tau \psi)$ with $\tau = 0.1$ and 0.25, respectively, which is the variance component parameter. The time for subject i is $t_i = \min(T_{ci}, T_{fi})$ and $\delta_i = I(T_{fi} \leq T_{ci})$. We selected λ , the mean of the exponential hazard function, corresponding to different censoring rates $\sum_{i=1}^N \delta_i / N = 50\%, 75\%$ and 90% . We repeated the simulation for 500 times. For each phenotype set, a null frailty model was fitted in Step 1 with the covariate X_1 . In Step 2, we conducted single variant association tests on 1.9 million simulated genetic markers. In totally, about 9.4×10^8 association tests were conducted. We evaluated the empirical type I error rates at the type I error rate $\alpha = 1 \times 10^{-6}$ and 5×10^{-8} as shown in **Supplementary Table 4** and **Supplementary Figure 8A**. These results have indicated that GATE can produce well calibrated type I error rates in the presence of sample relatedness at the significance level, while GATE-no SPA (similar to COXMEG) has inflated type I error rates and inflation gets larger than censoring rates is higher (**Supplementary Table 4**). For example, GATE-no SPA has type I error rate 8.9×10^{-6} at $\alpha = 5 \times 10^{-8}$ when censoring rate is 75% and 2.8×10^{-5} when censoring rate is 90% with $\tau = 0.1$.

To evaluate whether GATE can control type I error rates in presence of cryptic sample relatedness, we have randomly selected $N = 10,000$ samples with white British ancestry from UK Biobank and simulated TTE phenotypes based on the observed genotypes of these subjects in the approach described above for pedigree-based data sets, except that under the null hypothesis of no genetic effects, $\eta_i = X_{1i}^\top \alpha + \sum_{j=1}^L \hat{G}_{ij} \beta$ and was simulated based on real genotypes of randomly selected $L = 30,000$ LD-pruned ($r^2 < 0.2$) markers from the odd chromosomes with $MAF \geq 1\%$. The real genotypes were used for simulating real sample relatedness in the null model. In particular, X_1 is a covariate that was randomly drawn from $N(0, 1)$, α is the coefficient and is 1, \hat{G}_{ij} is the standardized genotype value for the j th marker of i th subject and β is the genetic effect size following $N(0, \tau/L)$, where $\tau = 0.25$, which is the variance component parameter. The time for subject i is $t_i = \min(T_{ci}, T_{fi})$ and $\delta_i = I(T_{fi} \leq T_{ci})$. We selected λ , the mean of the exponential hazard function, corresponding to different censoring rates $\sum_{i=1}^N \delta_i / N = 50\%, 75\%$ and 90% . We repeated the simulation for 100 times. For each phenotype set, a null frailty model was fitted in Step 1 with covariates including the first 4 genetic principal components, which were estimated for all White-British participants in the UK Biobank, and X_1 . In Step 2, we conducted single variant association tests on genetic markers on the even chromosome. In total, 8.3×10^8 were conducted. We evaluated the empirical type I error rates at the type I error rate $\alpha = 1 \times 10^{-6}$ and 5×10^{-8} as shown in **Supplementary Table 5** and **Supplementary Figure 8B**, which suggests that GATE produces well calibrated type I error rates in the presence of cryptic relatedness at the corresponding significance levels.

To evaluate the empirical power of GATE and compare the power to COXMEG, phenotypes were generated under the alternative hypothesis for 10,000 samples, which contain 500 families and 5,000 independent samples. The family pedigree is shown in the **Supplementary Figure 7**. We simulated 100 datasets with 10 genetic markers with different hazard ratios. Power was evaluated at $\alpha = 5 \times 10^{-8}$ with the censoring rate 50% for MAF 0.05 and 0.2 as presented in the **Supplementary Figure 10**.

Building the UK Biobank TTE Phenome.

The time-to-event phenotypes for the UK Biobank were constructed as the disease phenotypes defined based on the hierarchical PheCodes⁴⁰ that represent different disease groups. The ICD9 and ICD10 codes were mapped to PheCodes using a combination of available maps through the Unified Medical Language System (see URLs) and other sources, string matching, and manual review^{19,40}. For each PheCode, the subjects who had the PheCode were regarded as having events, and the subjects who did not have the PheCode were regarded as censored. For each failed subject, the TTE (failure time) was calculated by subtracting the birth year from the earliest time of diagnosis of any of the PheCode-specific ICD codes, rounded to the nearest full month. To obtain the TTE (censoring time) for each censored subject, the birth year was subtracted from the time of the last non-imaging visit to any of the UK Biobank ascertainment centers, or the last time any ICD code was recorded for that subject, or the time of death if death was recorded during the course of the study, whichever is latest, rounded to the nearest full month. For lifespan, the subjects who had their death recorded, were assigned the failed status with the ages at death as the corresponding TTE, and the subjects who did not have their death recorded were assigned the censored status with the TTE defined as before.

FinnGen

FinnGen is a public-private partnership project combining genotype data from Finnish biobanks and digital health record data from Finnish health registries (<https://www.finnngen.fi/en>). Release 5 analysis contains 218,792 samples after quality control with population outliers excluded via principal component analysis based on genetic data. TTE phenotypes were constructed from population registries and ICD10 codes, and harmonizing definitions over ICD8 and ICD9, including ischemic heart disease (N events=30,952, N censored=187838, censoring rate=85.8%), female breast cancer (N events=8,401, N censored=114,878, censoring rate=93.2%), glaucoma (N events=8,591, N censored=210199, censoring rate=96.1%) and Alzheimer's disease (N events=3,899, N censored = 207,324, censoring rate=98.2%). We conducted genome-

wide survival analysis using GATE with the first ten genetic PCs, sex, genotyping batch and birth year as covariates and 240,000 pruned genetic markers for GRM estimation. Patients and control subjects in FinnGen provided informed consent for biobank research, based on the Finnish Biobank Act. Alternatively, older research cohorts, collected prior the start of FinnGen (in August 2017), were collected based on study-specific consents and later transferred to the Finnish biobanks after approval by Fimea, the National Supervisory Authority for Welfare and Health. Recruitment protocols followed the biobank protocols approved by Fimea. The Coordinating Ethics Committee of the Hospital District of Helsinki and Uusimaa (HUS) approved the FinnGen study protocol Nr HUS/990/2017.

The FinnGen study is approved by Finnish Institute for Health and Welfare (THL), approval number THL/2031/6.02.00/2017, amendments THL/1101/5.05.00/2017, THL/341/6.02.00/2018, THL/2222/6.02.00/2018, THL/283/6.02.00/2019, THL/1721/5.05.00/2019, Digital and population data service agency VRK43431/2017-3, VRK/6909/2018-3, VRK/4415/2019-3 the Social Insurance Institution (KELA) KELA 58/522/2017, KELA 131/522/2018, KELA 70/522/2019, KELA 98/522/2019, and Statistics Finland TK-53-1041-17. The Biobank Access Decisions for FinnGen samples and data utilized in FinnGen Data Freeze 5 include: THL Biobank BB2017_55, BB2017_111, BB2018_19, BB_2018_34, BB_2018_67, BB2018_71, BB2019_7, BB2019_8, BB2019_26, Finnish Red Cross Blood Service Biobank 7.12.2017, Helsinki Biobank HUS/359/2017, Auria Biobank AB17-5154, Biobank Borealis of Northern Finland_2017_1013, Biobank of Eastern Finland 1186/2018, Finnish Clinical Biobank Tampere MH0004, Central Finland Biobank 1-2017, and Terveystalo Biobank STB 2018001.

Genome build.

The genomic coordinates reported in this paper were based on NCBI Build 37/UCSC hg19.

References

1. Cox, D.R. Regression Models and Life-Tables. *Journal of the Royal Statistical Society. Series B (Methodological)* **34**, 187-220 (1972).
2. Lee, E. & Go, O. Survival analysis in public health research. *Annual Review of Public Health* **18**, 105-34 (1997).
3. Dg, A., BI De, S., Sb, L. & Ka, S. Review of survival analyses published in cancer journals. *British Journal of Cancer* **72**, 511 (1995).
4. Kasza, J., Wraith, D., Lamb, K. & Wolfe, R. Survival analysis of time-to-event data in respiratory health research studies. Vol. 19 483-492 (2014).
5. He, L. *et al.* Genome-wide time-to-event analysis on smoking progression stages in a family-based study. *Brain and Behavior* **6**, n/a-n/a (2016).
6. Phipps, A.I. *et al.* Common genetic variation and survival after colorectal cancer diagnosis: a genome-wide analysis. *Carcinogenesis* **37**, 87-95 (2016).
7. Johnson, D.C. *et al.* Genome-wide association study identifies variation at 6q25.1 associated with survival in multiple myeloma. *Nature Communications* **7**(2016).
8. Kulminski, A.M. *et al.* Pleiotropic Associations of Allelic Variants in a 2q22 Region with Risks of Major Human Diseases and Mortality.(Research Article)(Report). *PLoS Genetics* **12**, e1006314 (2016).
9. Wu, C. *et al.* Genome-wide association study of survival in patients with pancreatic adenocarcinoma. *Gut* **63**, 152 (2014).
10. Lee, S. & Lim, H. Review of statistical methods for survival analysis using genomic data. *Genomics & informatics* **17**, e41-e41 (2019).
11. Bi, W., Fritsche, L.G., Mukherjee, B., Kim, S. & Lee, S. A Fast and Accurate Method for Genome-Wide Time-to-Event Data Analysis and Its Application to UK Biobank. *Am J Hum Genet* **107**, 222-233 (2020).
12. Green, M.S. & Symons, M.J. A comparison of the logistic risk function and the proportional hazards model in prospective epidemiologic studies. *Journal of Chronic Diseases* **36**, 715-723 (1983).
13. Callas, P., Pastides, H. & Hosmer, D. Empirical comparisons of proportional hazards, Poisson, and logistic regression modeling of occupational cohort data. *American Journal of Industrial Medicine* **33**, 33-47 (1998).
14. Staley, J.R. *et al.* A comparison of Cox and logistic regression for use in genome-wide association studies of cohort and case-cohort design. *European journal of human genetics : EJHG* **25**, 854-862 (2017).
15. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203-209 (2018).
16. Loh, P.R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet* **47**, 284-90 (2015).
17. Svishcheva, G.R., Axenovich, T.I., Belonogova, N.M., van Duijn, C.M. & Aulchenko, Y.S. Rapid variance components-based method for whole-genome association analysis. *Nat Genet* **44**, 1166-70 (2012).
18. Jiang, L. *et al.* A resource-efficient tool for mixed model association analysis of large-scale data. *Nature Genetics* **51**, 1749-2 (2019).

19. Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat Genet* **50**, 1335-1341 (2018).
20. Chen, H. *et al.* Control for Population Structure and Relatedness for Binary Traits in Genetic Association Studies via Logistic Mixed Models. *The American Journal of Human Genetics* **98**, 653-666 (2016).
21. Vaupel, J., Manton, K. & Stallard, E. The impact of heterogeneity in individual frailty on the dynamics of mortality. *Demography* **16**, 439-454 (1979).
22. Hougaard, P. Frailty models for survival data. *Lifetime data analysis* **1**, 255-273 (1995).
23. Clayton, D. & Cuzick, J. Multivariate Generalizations of the Proportional Hazards Model. *Journal of the Royal Statistical Society: Series A (General)* **148**, 82-108 (1985).
24. Klein, J.P. Semiparametric estimation of random effects using the Cox model based on the EM algorithm. *Biometrics* **48**, 795-806 (1992).
25. McGilchrist, C.A. REML estimation for survival models with frailty. *Biometrics* **49**, 221-5 (1993).
26. Petersen, J.H., Andersen, P.K. & Gill, R.D. Variance components models for survival data. *Statistica Neerlandica* **50**, 193-211 (1996).
27. Korsgaard, I.R. & Andersen, A.H. The Additive Genetic Gamma Frailty Model. *Scandinavian Journal of Statistics* **25**, 225-269 (1998).
28. Ripatti, S. & Palmgren, J. Estimation of multivariate frailty models using penalized partial likelihood. *Biometrics* **56**, 1016-22 (2000).
29. Therneau, T.M., Grambsch, P.M. & Pankratz, V.S. Penalized Survival Models and Frailty. *Journal of computational and graphical statistics* **12**, 156-175 (2003).
30. Therneau, T.M. *coxme: Mixed Effects Cox Models*. (2019).
31. He, L. & Kulminski, A.M. Fast Algorithms for Conducting Large-Scale GWAS of Age-at-Onset Traits Using Cox Mixed-Effects Models. *Genetics* **215**, 41-58 (2020).
32. He, L. *coxme: Cox Mixed-Effects Models for Genome-Wide Association Studies*. (2020).
33. Ma, C., Blackwell, T., Boehnke, M., Scott, L.J. & Go, T.D.i. Recommended joint and meta-analysis strategies for case-control association testing of single low-count variants. *Genet Epidemiol* **37**, 539-50 (2013).
34. Dey, R., Schmidt, E.M., Abecasis, G.R. & Lee, S. A Fast and Accurate Algorithm to Test for Binary Phenotypes and Its Application to PheWAS. *Am J Hum Genet* **101**, 37-49 (2017).
35. Dey, R. *et al.* Robust meta-analysis of biobank-based genome-wide association studies with unbalanced binary phenotypes. *Genet Epidemiol* **43**, 462-476 (2019).
36. Daniels, H.E. Saddlepoint Approximations in Statistics. *Ann. Math. Statist.* **25**, 631-650 (1954).
37. Breslow, N.E. & Clayton, D.G. Approximate Inference in Generalized Linear Mixed Models. *Journal of the American Statistical Association* **88**, 9-25 (1993).
38. Gilmour, A.R., Thompson, R. & Cullis, B.R. Average Information REML: An Efficient Algorithm for Variance Parameter Estimation in Linear Mixed Models. *Biometrics* **51**, 1440-1450 (1995).
39. Tsuruta, S., Misztal, I. & Strandén, I. Use of the preconditioned conjugate gradient algorithm as a generic solver for mixed-model equations in animal breeding applications. *J Anim Sci* **79**, 1166-72 (2001).
40. Denny, J.C. *et al.* Systematic comparison of phenome-wide association study of electronic medical record data and genome-wide association study data. *Nat Biotechnol* **31**, 1102-10 (2013).
41. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nature Genetics* **48**, 1279-1283 (2016).
42. Walter, K. *et al.* The UK10K project identifies rare variants in health and disease. *Nature* **526**, 82-90 (2015).

43. Gagliano Taliun, S.A. *et al.* Exploring and visualizing large-scale genetic associations by using PheWeb. *Nature Genetics* **52**, 550-552 (2020).
44. Nelson, C.P. *et al.* Association analyses based on false discovery rate implicate new loci for coronary artery disease. *Nature Genetics* **49**, 1385-1391 (2017).
45. Deloukas, P. *et al.* Large-scale association analysis identifies new risk loci for coronary artery disease. *Nature genetics* **45**, 25-33 (2012).
46. Meyer, Kerstin B. *et al.* Fine-Scale Mapping of the FGFR2 Breast Cancer Risk Locus: Putative Functional Variants Differentially Bind FOXA1 and E2F1. *American journal of human genetics* **93**, 1046-1060 (2013).
47. Udler, M.S. *et al.* Fine scale mapping of the breast cancer 16q12 locus. *Human molecular genetics* **19**, 2507-2515 (2010).
48. Stone, E.M. Identification of a Gene That Causes Primary Open Angle Glaucoma. *Science (American Association for the Advancement of Science)* **275**, 668-670 (1997).
49. Burdon, K.P. *et al.* Genome-wide association study identifies susceptibility loci for open angle glaucoma at TMCO1 and CDKN2B-AS1. *Nature genetics* **43**, 574-578 (2011).
50. Moreno-Grau, S. *et al.* Genome-wide association analysis of dementia and its clinical endophenotypes reveal novel loci associated with Alzheimer's disease and three causality networks: The GR@ACE project. *Alzheimers Dement* **15**, 1333-1347 (2019).
51. Kaplan, E.L. & Meier, P. *Nonparametric Estimation from Incomplete Observations*, (Springer New York, 1992).
52. Wolters, F. *et al.* The impact of APOE genotype on survival: Results of 38,537 participants from six population-based cohorts (E2-CHARGE). *PLoS ONE* **14**, e0219668 (2019).
53. Rovio, S. *et al.* Leisure-time physical activity at midlife and the risk of dementia and Alzheimer's disease. *Lancet Neurol* **4**, 705-11 (2005).
54. Schuit, A.J., Feskens, E.J., Launer, L.J. & Kromhout, D. Physical activity and cognitive decline, the role of the apolipoprotein e4 allele. *Med Sci Sports Exerc* **33**, 772-7 (2001).
55. Smith, J.C., Nielson, K.A., Woodard, J.L., Seidenberg, M. & Rao, S.M. Physical activity and brain function in older adults at increased risk for Alzheimer's disease. *Brain Sci* **3**, 54-83 (2013).
56. Li, X. *et al.* Dynamic incorporation of multiple in silico functional annotations empowers rare variant association analysis of large whole-genome sequencing studies at scale. *Nature Genetics* **52**, 969-983 (2020).
57. Yang, J., Zaitlen, N.A., Goddard, M.E., Visscher, P.M. & Price, A.L. Advantages and pitfalls in the application of mixed-model association methods. *Nature genetics* **46**, 100-106 (2014).
58. Kang, H.M. *et al.* Variance component model to account for sample structure in genome-wide association studies. *Nature genetics* **42**, 348-354 (2010).
59. Wu, Michael C. *et al.* Rare-Variant Association Testing for Sequencing Data with the Sequence Kernel Association Test. *American journal of human genetics* **89**, 82-93 (2011).
60. Therneau, T.M., Grambsch, P.M. & SpringerLink (Online service). *Modeling Survival Data: Extending the Cox Model*. (Springer New York : Imprint: Springer, New York, NY, 2000).
61. Breslow, N.E. Discussion of the paper by D. R. Cox. *Journal of the Royal Statistical Society. Series B (Methodological)* **34**, 216-217 (1972).
62. Barndorff-Nielsen, O.E. Approximate Interval Probabilities. *Journal of the Royal Statistical Society. Series B (Methodological)* **52**, 485-496 (1990).
63. Kuonen, D. Saddlepoint Approximations for Distributions of Quadratic Forms in Normal Variables. *Biometrika* **86**, 929-935 (1999).
64. Abecasis, G.R., Cherny, S.S., Cookson, W.O. & Cardon, L.R. Merlin—rapid analysis of dense genetic maps using sparse gene flow trees. *Nature genetics* **30**, 97-101 (2001).

