

1 **Findings and insights from the genetic investigation of age of first reported occurrence for**  
2 **complex disorders in the UK Biobank and FinnGen**

3

4 Yen-Chen A. Feng<sup>1,2,3</sup>, Tian Ge<sup>1,3,4</sup>, Mattia Cordioli<sup>6</sup>, FinnGen, Andrea Ganna<sup>2,5,6</sup>, Jordan W.  
5 Smoller<sup>1,3,4\*</sup>, Benjamin M. Neale<sup>2,3,4,5\*</sup>

6

7 <sup>1</sup>Psychiatric and Neurodevelopmental Genetics Unit, Department of Psychiatry, Massachusetts General  
8 Hospital, Boston, Massachusetts, USA

9 <sup>2</sup>Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital,  
10 Boston, Massachusetts, USA

11 <sup>3</sup>Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge,  
12 Massachusetts, USA

13 <sup>4</sup>Harvard Medical School, Boston, Massachusetts, USA

14 <sup>5</sup>Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge,  
15 Massachusetts, USA

16 <sup>6</sup>Institute for Molecular Medicine Finland, University of Helsinki, Helsinki, Finland

17

18

19

20 **Correspondence:**

21 Jordan W. Smoller ([jsmoller@partners.org](mailto:jsmoller@partners.org))

22 Benjamin M. Neale ([bneale@broadinstitute.org](mailto:bneale@broadinstitute.org))

23

24 **Abstract**

25 Age of onset contains information on the timing of events relevant to disease etiology, but there  
26 has not been a systematic investigation of its heritability from GWAS data. Here, we characterize  
27 the genetic architecture of age of first occurrence and its genomic relationship with disease  
28 susceptibility for a wide range of complex disorders in the UK Biobank. For diseases with a  
29 sufficient sample size, we discover that age of first occurrence has non-trivial genetic contributions,  
30 some with specific genetic risk factors not associated with susceptibility to the disease. Through  
31 genetic correlation analysis, we show that an earlier health-event occurrence is correlated with a  
32 higher polygenic risk of disease susceptibility. An independent genetic investigation of the  
33 FinnGen cohort replicates the pattern of heritability and genetic correlation estimates. We then  
34 demonstrate that incorporating disease onset age with susceptibility may improve genetic risk  
35 prediction and stratification.

36

## 37 **Introduction**

38 Genome-wide association studies (GWAS) have revealed that most complex traits and diseases  
39 have an underlying polygenic component<sup>1,2</sup>. To date, disease-based GWAS have predominantly  
40 examined the genetic risk associated with *whether* an individual has ever been affected with a  
41 disease of interest (i.e., disease susceptibility) using a case-control design. Such design, however,  
42 typically uses lifetime risk to model the association of phenotypic variation with genetic variation  
43 and ignores the time component of *when* a disease occurs for an individual.

44 Studies have shown that age of onset of a disease itself plays a pertinent role in  
45 understanding the genetic etiology of disease development. Using a genome-wide approach,  
46 efforts have been made to identify genetic modifiers of disease onset age<sup>3-5</sup>, as well as genetic  
47 risk factors distinctive to different age-of-onset groups<sup>6-8</sup>. As shown from twin and single-  
48 nucleotide-polymorphism (SNP) heritability analyses, the phenotypic variance explained by  
49 genetic variation can change across the lifespan for many traits<sup>9-11</sup>. Furthermore, the polygenic  
50 model has suggested that there may exist a correlation between age of onset and susceptibility  
51 to a disease, such that individuals with a higher genetic loading might develop the disease at an  
52 earlier age<sup>12-14</sup>. Despite these efforts, however, there has been comparatively less investigation  
53 of the heritability of age of onset for disease phenotypes.

54 In recent years, large-scale biobank datasets that include broad phenotypic information  
55 (e.g., UK Biobank<sup>15,16</sup>, FinnGen<sup>17,18</sup>) have provided an unprecedented resource to study the  
56 causes of disease at scale with a linkage to genetic data. Related to age of onset, dates of health-  
57 outcome events, such as diagnosis, treatment, and death, have been made available through  
58 questionnaire-based or hospital records within such biobanks. Collectively, this provides an  
59 opportunity to systematically characterize the genetic construct of disease onset age, the specific  
60 genetic risk factors that may alter it, and how these findings correlate or differ from the genetic  
61 basis of susceptibility.

62           Here, we present a deep investigation into the genetic architecture of age of first  
63           occurrence and its genetic overlap with susceptibility across a wide range of complex disorders  
64           in the UK Biobank (UKBB). We leveraged measurements from three different phenotypic datasets  
65           that either directly or approximately capture age of onset: *self-reported* age of diagnosis of a  
66           medical condition (SR), age of first in-patient ICD-10 diagnosis or hospitalization episode from  
67           *hospital in-patient* records (HIP), and age of the earliest event occurrence *combining* self-report,  
68           in-patient, primary care, and death records (COMB; Methods). We refer to these definitions more  
69           generally as “age of first occurrence” of a particular disease or medical condition. For diseases  
70           with a sufficient sample size, we show that age of first occurrence is moderately heritable, some  
71           with specific genetic risk factors not associated with susceptibility. Across disease domains, there  
72           is an overall inverse genetic correlation between age of first occurrence and susceptibility.  
73           Independent of the UKBB cohort, we then show that a similar pattern of heritability and genetic  
74           correlations exists in the FinnGen study, which has a longer follow-up. Finally, we demonstrate  
75           that information on age of first occurrence has the potential to improve polygenic risk prediction  
76           for disease susceptibility and patient stratification.



## 77 **Results**

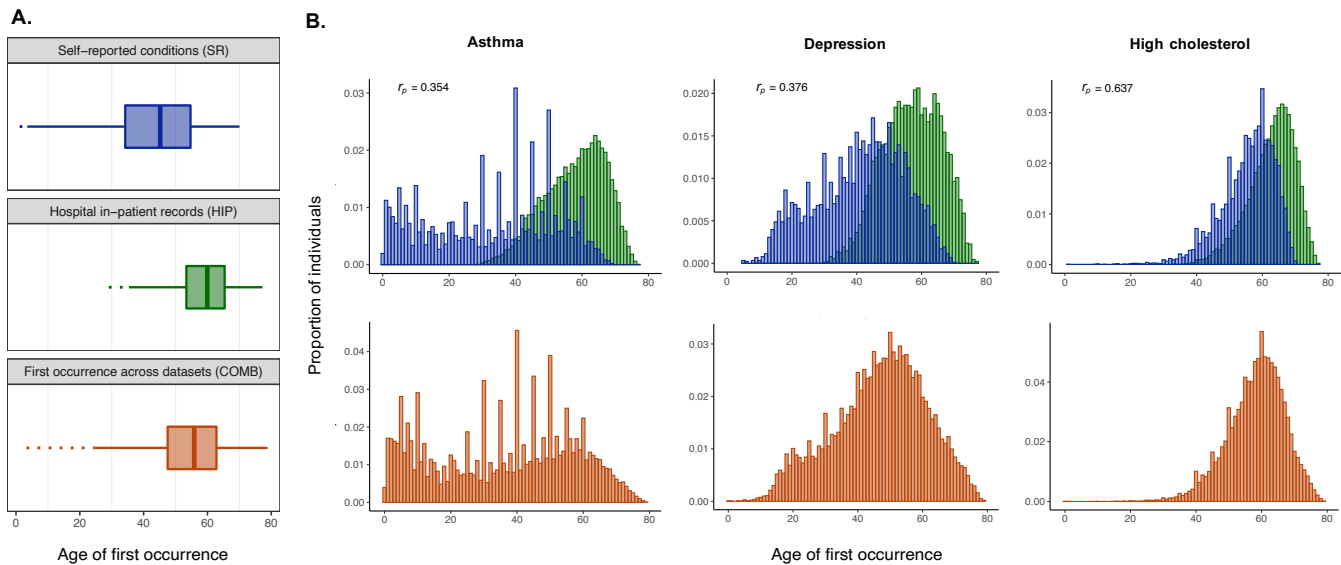
### 78 ***Characterizing age of first occurrence in the UK Biobank***

79 We estimated age of first occurrence for medical conditions aggregated from the SR, HIP, and  
80 COMB datasets in the UKBB White-British subset (N=361,140). We extracted every possible  
81 disease definition following hierarchical disease classifications where appropriate, resulting in  
82 over 3,000 clinical terms (Methods). Among them, 70 SR, 224 HIP, and 164 COMB terms had at  
83 least 5000 affected individuals (prevalence ~1.4%; Tables S1-3) and were retained for GWAS  
84 analysis.

85 Age of first occurrence ranged from 0 to 70 years of age in the SR dataset (median: 17-  
86 59), 30 to 80 in the HIP dataset (median: 36-66), and 0 to 80 in the COMB dataset (median: 6-68;  
87 Figures 1&S1; Tables S4-6). We then compared the distribution of age of first occurrence by data  
88 source for 26 disease phenotypes where definitions were comparable between SR and HIP, of  
89 which 15 mapped across all three datasets (Table S7). Among these diseases, SR age of first  
90 occurrence was consistently younger than that based on the HIP dataset. Some diseases had a  
91 higher prevalence in SR (e.g., hypertension, asthma, and arthrosis) while others (e.g., hernia,  
92 gallstones, and cancers) were more prevalent in the HIP dataset. Distribution of age of first  
93 occurrence in SR varied extensively by trait and exhibited a spiky behavior by quartiles (0.25, 0.5,  
94 0.75, or 1.00), reflecting that age of the reported diagnosis was recorded as integers in the  
95 questionnaire. Comparatively, age of first occurrence in HIP was estimated directly from the  
96 recorded dates and was smoothly and consistently distributed across diseases with a bell-like  
97 shape. The COMB dataset had the largest number of cases per definition among all three sources  
98 and showed an overall “merged” distribution of age of first occurrence in SR and HIP (Figures  
99 1B&S1).

100 Considering diagnosis of the same medical condition from both SR and HIP, the  
101 distribution of age of first occurrence between the two datasets showed little (e.g, migraine,  
102 median difference  $m_{diff} = 31.6$ , phenotypic correlation  $r_p = 0.12$ ; asthma,  $m_{diff} = 27.2$ ,  $r_p = 0.35$ ) to

103 moderate overlap (e.g, high cholesterol,  $m_{diff} = 6.4$ ,  $r_p = 0.64$ ; gallstones,  $m_{diff} = 8.6$ ,  $r_p = 0.87$ ),  
104 with values in HIP shifting toward the right ( $m_{diff}$ : 6-32;  $r_p$ : 0.1-0.9; avg.  $r_p = 0.46$ ; Figures 1B&S2;  
105 Table S7). The extent of overlap generally increased with age of first occurrence of the condition  
106 (Figure S3) but overall suggested that the actual disease “onset” age for most diseases was left-  
107 truncated in the hospitalization records, reflecting the fact that ICD-10 was implemented and  
108 integrated in the UK hospital episode statistics in the 1990s. As expected, age of first occurrences  
109 in the COMB dataset overlapped substantially with both the SR (avg.  $r_p = 0.91$ ) and the HIP (avg.  
110  $r_p = 0.69$ ) datasets (Table S7).  
111



**Figure 1. Distribution of age of first occurrence of disease phenotypes from three phenotypic datasets in UKBB**

- A.** An averaged distribution of age of first occurrence is shown across 70 SR (blue), 224 HIP (green), and 164 COMB (orange) disease definitions in each dataset. Dotted line indicates the outlying range of values. Age of first occurrence ranges from 0-70 in the SR dataset, 30-80 in the HIP dataset, and 0-80 in the COMB dataset. Spikes in the SR phenotypes reflect that the values are recorded in quartiles (0.25, 0.5, 0.75, or 1.00).
- B.** Distribution of age of first occurrence differs by trait and data source. Shown here are three selected disease phenotypes with matching definitions across datasets. SR and HIP conditions show little to moderate overlap in age of first occurrence, as measured by Pearson's correlation coefficient ( $r_p$ ; top), while the COMB conditions exhibit a merged distribution of SR and HIP (bottom).

## 112 **Common genetic associations with age of first occurrence**

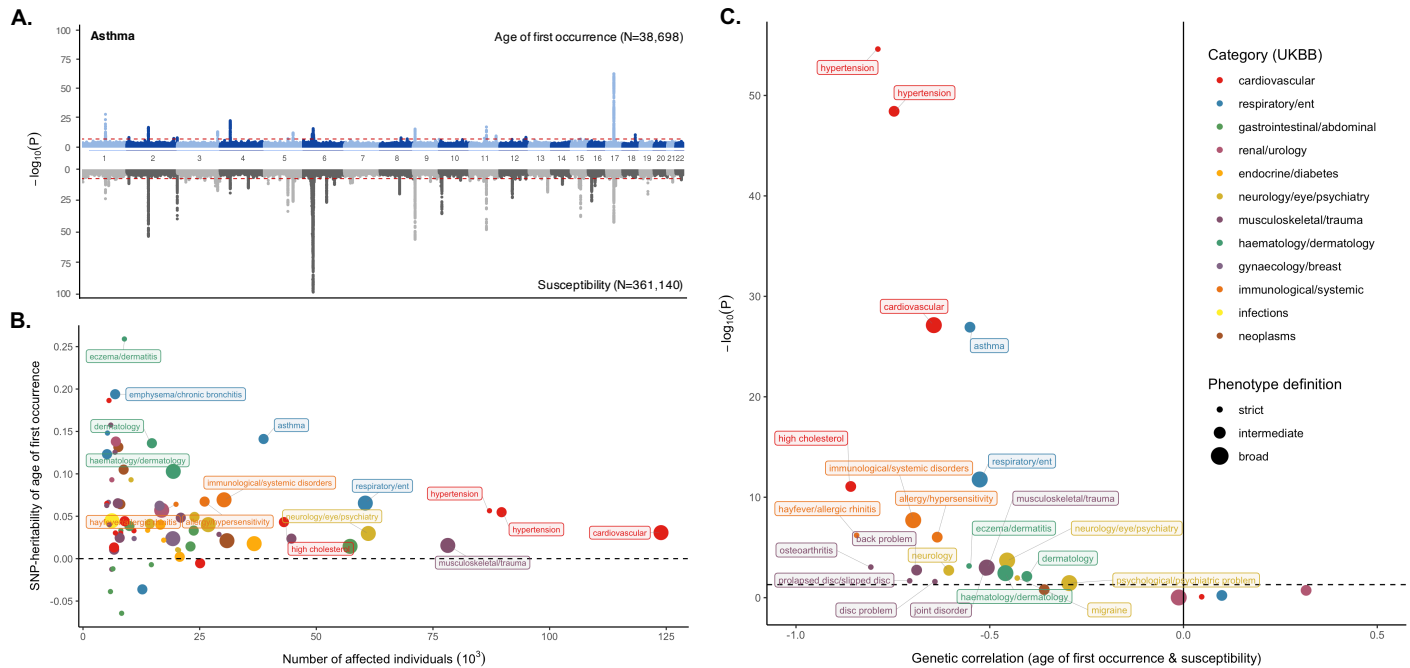
113 We performed a GWAS of age of first occurrence for the 70 SR, 224 HIP, and 164 COMB disease  
114 definitions with >5000 affected individuals. For each condition, a case-control GWAS was also  
115 performed, treating non-diseased individuals as the control group and adjusting for the same sets  
116 of covariates (Methods). Univariate Linkage Disequilibrium Score Regression (LDSR) intercepts  
117 for all GWASes had a value close to 1, suggesting little or no signs of inflation in association  
118 statistics due to population stratification or other confounding factors (Figure S4; Tables S8-10).

119 Genome-wide significant SNPs ( $P < 5 \times 10^{-8}$ ) that may alter age of first occurrence were  
120 identified for 31 SR, 57 HIP, and 42 COMB disease definitions, ranging from 1 to 20 independently  
121 associated loci. Among the three data sources, age of first occurrence defined for hospitalization  
122 events in the HIP dataset had the least number of associated loci (max. 2 loci; Tables S8-10).  
123 While most of the identified signals were a subset of the significant associations in the  
124 corresponding susceptibility GWAS (Figures 2A&S5-7), some of these GWASes contained loci  
125 significantly associated with age of first occurrence but not with susceptibility, suggesting a role  
126 in modifying disease development that precedes disease onset (Tables S8-11).

127 Disease phenotypes showing the largest number of independent associations with age of  
128 first occurrence included asthma (total number of significant hits,  $n_{sig} = 19$ ; total number of unique  
129 hits not seen in susceptibility GWAS,  $n_{uniq} = 5$ ), hypertension ( $n_{sig} = 9$ ), high cholesterol level ( $n_{sig}$   
130  $= 6$ ,  $n_{uniq} = 1$ ), and eczema ( $n_{sig} = 5$ ,  $n_{uniq} = 2$ ) in the SR dataset. Top results in the HIP dataset  
131 included mental disorders ( $n_{sig} = 2$ ,  $n_{uniq} = 2$ ), monoarthritis ( $n_{sig} = 2$ ,  $n_{uniq} = 2$ ), substance-related  
132 disorders ( $n_{sig} = 2$ ,  $n_{uniq} = 2$ ), and type 2 diabetes ( $n_{sig} = 2$ ,  $n_{uniq} = 1$ ); in the COMB dataset those  
133 included asthma ( $n_{sig} = 20$ ,  $n_{uniq} = 7$ ), disorders of lipoprotein metabolism ( $n_{sig} = 13$ ), measles ( $n_{sig}$   
134  $= 6$ ,  $n_{uniq} = 6$ ), and dermatitis ( $n_{sig} = 5$ ,  $n_{uniq} = 1$ ) (Tables S8-10).

135 Comparing the overlapping loci between age of first occurrence and susceptibility revealed  
136 that their effects were often in opposite directions and although in close proximity, the lead SNPs  
137 were mostly different. A few loci showed an even stronger association with age of first occurrence

138 than with susceptibility. For instance, the significant locus on chromosome 1 for SR asthma was  
 139 associated with a reduced age of onset by 4.8 years (lead SNP 1:152285861:G:A,  $P = 4.6 \times 10^{-29}$ )  
 140 but an increased risk of asthma by 1.2 fold (lead SNP 1:152179152:C:T,  $P = 7.4 \times 10^{-24}$ ; Table  
 141 S11).  
 142



**Figure 2. Genetic characterization of age of first occurrence and its relationship with susceptibility in the UKBB SR dataset**

- A.** A Miami plot of GWAS results reveal overlapping and distinct genetic associations between age of first occurrence (top) and case-control status (bottom) of SR asthma. Each dot represents a single SNP. P-values are shown on the  $-\log_{10}$  scale on the y-axis, plotted against chromosome positions on the x-axis. The red dashed lines denote the genome-wide significance threshold at  $P = 5 \times 10^{-8}$ .
- B.** SNP-heritability estimates for age of first occurrence ( $h^2_{aof}$ ) across 70 SR disease definitions suggest non-trivial common genetic contributions.  $h^2_{aof}$  was estimated from univariate LDSR. Each dot represents an individual disease, colored by disease categories used in the SR dataset; a larger dot corresponds to a broader disease definition. Labeled are conditions with a significant  $h^2_{aof}$  at FDR < 0.05. Heritability analysis of HIP and COMB diseases reveal a similar pattern in Figure S9.
- C.** Genetic correlation ( $r_g$ ) analysis suggests an inverse genomic relationship between age of first occurrence and susceptibility for diseases with a significant heritability for both traits.  $r_g$  was estimated using bivariate LDSR. The dashed line denotes nominal significance at  $P = 0.05$ ; labeled are conditions with a significant  $r_g$  at FDR < 0.05. Analysis of HIP and COMB diseases show a similar pattern in Figure S13.

### 143 **Moderate SNP-heritability of age of first occurrence**

144 Heritability analysis using LDSR showed that age of first occurrence for complex diseases was  
145 moderately heritable: 27/70 SR endpoints, 49/224 HIP endpoints, and 30/164 COMB endpoints  
146 had a significantly non-zero SNP-heritability for the age at which an individual first developed a  
147 given condition ( $h_{aof0}^2$ ; nominal p-value < 0.05), ranging from 1 to 25% with an average of 7-9%  
148 (Figures 2B&S8-9; Tables S12-14).  $h_{aof0}^2$  estimates in the HIP dataset were slightly lower than in  
149 the other two datasets.

150 Diseases with a heritable age of first occurrence in the SR dataset included cardiovascular  
151 (e.g., hypertension,  $h_{aof0}^2 = 0.055$ ), respiratory (e.g., asthma,  $h_{aof0}^2 = 0.141$ ), dermatological (e.g.,  
152 eczema/dermatitis,  $h_{aof0}^2 = 0.259$ ) and immunological (e.g., allergy/hypersensitivity/anaphylaxis,  
153  $h_{aof0}^2 = 0.067$ ) related traits (Figure 2B; see p-values in Table S12). Top diseases in the HIP  
154 dataset were seen among circulatory system (e.g., hypertension,  $h_{aof0}^2 = 0.049$ ), genitourinary  
155 (e.g., irregular menstrual cycle/bleeding,  $h_{aof0}^2 = 0.127$ ), digestive (e.g., cholelithiasis and  
156 cholecystitis,  $h_{aof0}^2 = 0.096$ ), mental disorders (e.g., psychological disorders,  $h_{aof0}^2 = 0.044$ ), and  
157 neoplasms (e.g., benign or malignant skin cancer,  $h_{aof0}^2 > 0.10$ ; Figure S9A; Table S13).  
158 Significant  $h_{aof0}^2$  in the COMB dataset encompassed some of the top non-cancer conditions from  
159 SR and HIP, as well as other traits like myocardial infarction ( $h_{aof0}^2 = 0.114$ ), non-insulin-  
160 dependent diabetes ( $h_{aof0}^2 = 0.067$ ), hypothyroidism ( $h_{aof0}^2 = 0.055$ ), and diaphragmatic hernia  
161 ( $h_{aof0}^2 = 0.036$ ) (Figure S9B; Table S14).

162 Among the 26 mapped disease definitions,  $h_{aof0}^2$  estimates showed variability across  
163 datasets, particularly between SR and HIP, with some diseases more consistently estimated than  
164 others (e.g., hypertension,  $h_{aof0}^2 \sim 5\%$ ; Table S7; Figure S10A). Conditions with a significant  $h_{aof0}^2$   
165 in SR but not in HIP—often also larger in magnitude—tended to be chronic that could start early  
166 in life or are mild in presentations (e.g., asthma, disc problem, high cholesterol, and

167 dermatological conditions). Conversely, a significant  $h_{aof0}^2$  in HIP but not in SR was observed  
168 among conditions that appeared more likely to be acute or have an adult-onset that required in-  
169 hospital treatments (e.g., gallstones, diabetes, and cancers; Figure S11). Cross-dataset genetic  
170 correlation estimates ( $r_g$ ) of the mapped age-of-first-occurrence endpoints differed widely and  
171 many did not reach statistical significance. On average, these  $r_g$ 's were highest between SR and  
172 COMB, lowest between SR and HIP, which increased slightly with age of first occurrence,  
173 consistent with the pattern of phenotypic similarity (Table S7; Figure S12 vs. S3).

174 In contrast, SNP-heritability estimates for the corresponding susceptibility endpoints ( $h_{sus}^2$ )  
175 were all significant and showed a lesser degree of heterogeneity across datasets (Figure S10B),  
176 with most of the traits having a cross-dataset  $r_g$  closer to 1 (Tables S7).

177

### 178 ***Inverse genetic correlation between age of first occurrence and susceptibility***

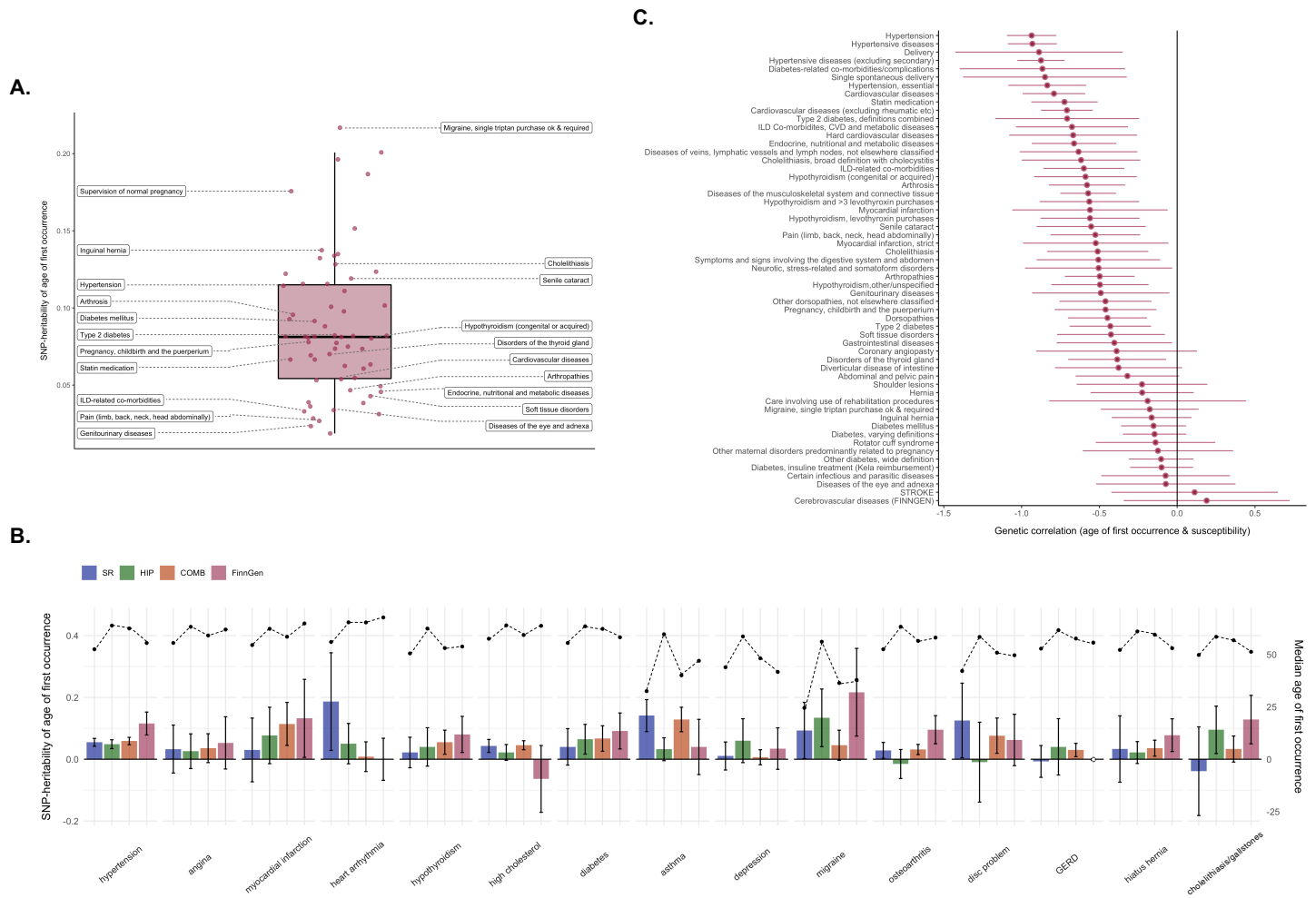
179 For diseases with a significant  $h_{aof0}^2$  and  $h_{sus}^2$ , we estimated  $r_g$  between age of first occurrence  
180 and susceptibility using bivariate LDSR (Methods). Interestingly, more than half of the tested traits  
181 showed a significant, negative  $r_g$  between the risk of developing the disease and the age of  
182 developing the disease, ranging from approximately -0.2 to -0.9 (Figures 2C&S13-14; Tables S12-  
183 14). This inverse genomic relationship was observed across disease categories in all three  
184 phenotypic datasets. A few traits had a non-significant, positive  $r_g$  that could result from a smaller  
185 sample size (Figure S13).

186 As with the pattern of  $h_{aof0}^2$ ,  $r_g$  results showed a different profile in each dataset.  
187 Hypertension was consistently the most significant phenotype with an  $r_g$  for age at first occurrence  
188 and susceptibility around -0.7 in all three datasets ( $P < 10^{-28}$ ). Other top diseases in the SR dataset  
189 with an  $r_g$  at FDR  $< 0.05$  included asthma ( $r_g = -0.55$ ), high cholesterol ( $r_g = -0.86$ ), hay  
190 fever/allergy ( $r_g < -0.6$ ), eczema ( $r_g = -0.55$ ), osteoarthritis ( $r_g = -0.81$ ), and migraine ( $r_g = -0.43$ ,)  
191 (Figure 2C; see  $p$ -values in Table S12). In the HIP dataset additional traits with the strongest  $r_g$   
192 included substance-related/psychological disorders ( $r_g < -0.3$ ), diseases of esophagus and

193 dysphagia ( $r_g = -0.79$ ), type 2 diabetes ( $r_g = -0.51$ ), COPD/bronchiectasis/asthma ( $r_g = -0.40$ ), non-  
194 specific chest pain ( $r_g = -0.58$ ), and cholelithiasis and cholecystitis ( $r_g = -0.37$ ) (Figure S14A; Table  
195 S13). In the COMB dataset other top results were predominantly a combination of traits from SR  
196 and HIP, such as disorders of lipoprotein metabolism ( $r_g = -0.99$ ), asthma ( $r_g = -0.56$ ), diverticular  
197 disease of intestine ( $r_g = -0.82$ ), other arthrosis ( $r_g = -0.72$ ), other hypothyroidism ( $r_g = -0.72$ ), and  
198 non-insulin-dependent diabetes mellitus ( $r_g = -0.59$ ) (Figure S14B; Table S14). Together, these  
199 observations are consistent with predictions of the polygenic model, in which an earlier onset or  
200 occurrence of a disease may correlate with a higher polygenic liability<sup>12</sup>.

201         Across pairs of diseases,  $r_g$ 's for susceptibility were similar to  $r_g$ 's for age of first occurrence:  
202 that is, diseases with extensive sharing of genetic risk appeared to also have genetically  
203 correlated age of first occurrence, especially for phenotypes measured by SR. (Figure S15).

204



**Figure 3. Genetic analysis of age of first occurrence in FinnGen and its comparison with UKBB results**

- A.** Distribution of heritability estimates from FinnGen for 64 diseases that have a significant  $h_{aof}^2$ . Labeled are selected conditions with a significant  $h_{aof}^2$  at FDR < 0.05.
- B.**  $h_{aof}^2$  estimates for 15 comparable disease definitions in UKBB and FinnGen show variable degree of similarity. Left axis denotes SNP-heritability shown in bar plots and the corresponding 95% confidence intervals (95% CI). Right axis shows the median age of first occurrence for each condition indicated in dotted lines. The full comparison of all 26 matched phenotypes is available in Table S18 and Figure S19.
- C.** A negative  $r_g$  between age of first occurrence and disease susceptibility is observed for many of the tested diseases, consistent with the findings in UKBB. Shown are  $r_g$  estimates and its 95% C.I. for diseases with a significant heritability for both traits.

205

206 **Similar genetic architecture of age of first occurrence in FinnGen**

207 Next, we sought replication for UKBB findings, specifically the pattern of heritability and genetic

208 correlation, in FinnGen, based on its v4 release of 130,423 unrelated individuals. FinnGen is a



209 registry-based cohort that follows health events across a lifetime for an individual, including  
210 medication history (Figure S16; Supplementary Notes). A fraction of the FinnGen participants  
211 were ascertained in hospitals or disease-based cohorts<sup>14</sup>, making it a case-enriched cohort  
212 relative to UKBB (Table S15). Because registers data were established since the 1960s, FinnGen  
213 has the advantage of a prospective cohort that contains a longer follow-up time and a wider age  
214 range (0.08 to 98.98 years old at recruitment) for participants compared to UKBB. We analyzed  
215 280 disease phenotypes with a sufficient sample size following the same analytical pipeline in  
216 UKBB where possible (Methods; Table S15). Age of first occurrence was defined as the earliest  
217 age of an event in the registries (range: 0-100; median: 9-72; Table S16; Figure S17).

218 As in the UKBB, some diseases were found to have genetic risk factors associated with  
219 age of first occurrence but not with susceptibility (e.g., diabetes; Table S17). 63 of the 280 medical  
220 conditions had a significant heritability for age of first occurrence, ranging from 2 to 22%, many of  
221 which were among top results in UKBB (Figure 3A; Table S18). These included conditions  
222 identical to or close proxies for phenotypes in UKBB, particular those defined in the HIP dataset,  
223 such as hypertension ( $h_{aof}^2 = 0.116$ ), arthropathies ( $h_{aof}^2 = 0.049$ ), statin medication ( $h_{aof}^2 =$   
224  $0.067$ ), diabetes ( $h_{aof}^2 = 0.093$ ), cholelithiasis ( $h_{aof}^2 = 0.128$ ), and migraine (triptan purchase  
225 and/or ICD diagnosis;  $h_{aof}^2 = 0.2168$ ). Conditions with a significant  $h_{aof}^2$  in FinnGen but not in  
226 UKBB involved pregnancy and childbirth ( $h_{aof}^2 = 0.077$ ) and eye diseases (e.g., senile cataract,  
227  $h_{aof}^2 = 0.119$ ; Figure 3A; see  $p$ -values in Table S18).

228 To make a closer comparison between UKBB and FinnGen, we focused on the 26 mapped  
229 disease definitions. Most of these conditions had a higher in-sample prevalence in FinnGen  
230 compared to UKBB, except for high cholesterol (Table S19). Median age of first occurrence of  
231 FinnGen endpoints generally fell between UKBB-SR and the other two UKBB datasets. Locus-  
232 level comparison of age-of-first-occurrence GWAS revealed variable concordance between the  
233 two cohorts, with a few FinnGen traits showing a consistent direction of effects with UKBB with a  
234 sign-test  $p$ -value  $< 0.05$  (e.g., hypertension, diabetes; Table S20; Figure S18). Heritability

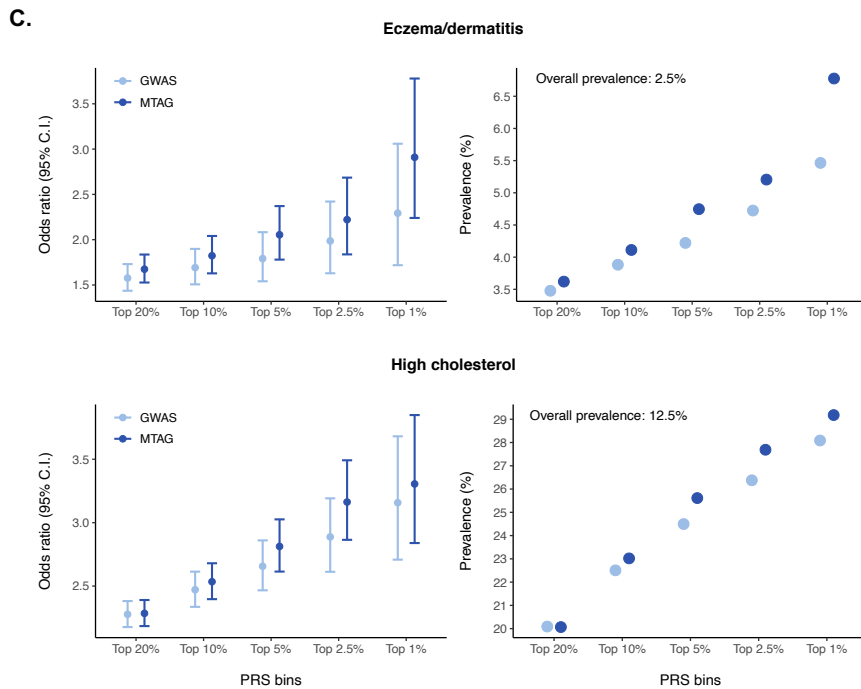
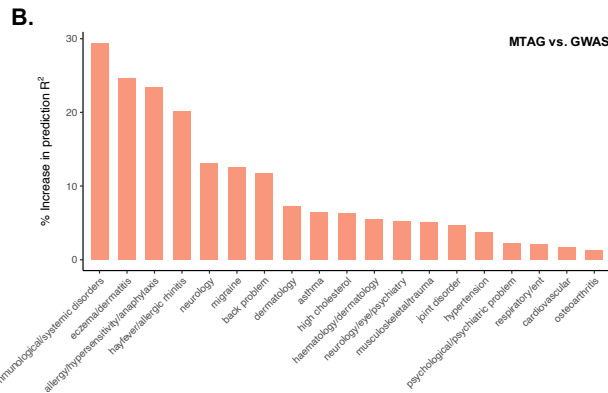
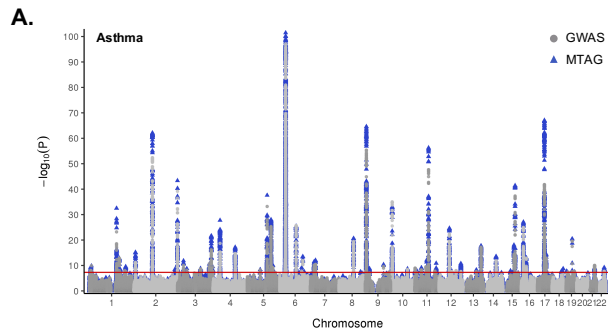
235 estimates showed that several of these diseases had a higher  $h_{aof0}^2$  in FinnGen than in UKBB  
236 (e.g., hypothyroidism, UKBB <5.5%, FinnGen 8.0%; musculoskeletal disorders, UKBB ~1.8%,  
237 FinnGen 4.7%; Figures 3B&S18; Table S19). Hypertension remained as the most significant trait  
238 with a non-zero  $h_{aof0}^2$ , whereas angina and depression both showed little evidence of  $h_{aof0}^2$  in  
239 either cohort. The significant  $h_{aof0}^2$  for asthma and high cholesterol in UKBB (SR and COMB) was  
240 not seen in FinnGen. These heterogeneous results for individual diseases might result from  
241 differences in sample ascertainment, criteria to define event occurrence age, as well as the  
242 sample size limitation of a case-only analysis.

243         Aside from these differences, genetic correlation analysis found that the majority of the  
244 heritable traits in FinnGen had a negative  $r_g$  between age of first occurrence and susceptibility  
245 (Figure 3C; Table S18), corroborating the UKBB finding of an inverse genomic relationship  
246 between event onset age and polygenic burden of the disease.

247

#### 248 ***MTAG of susceptibility and age of first occurrence improved polygenic risk prediction***

249 Given the strong genetic overlap between susceptibility and age of first occurrence, we then  
250 performed Multi-trait Analysis of GWAS (MTAG)<sup>19</sup> in UKBB to jointly analyze the two traits that  
251 shared a significant  $r_g$  for each of the phenotypic datasets (Methods). For many of the tested  
252 disease endpoints, MTAG analysis identified additional significant loci not found in the original  
253 GWAS of susceptibility, showing an enhanced power for loci discovery equivalent to 0.5% to 30%  
254 increase in sample size. HIP endpoints had the least number of additional loci among the three  
255 datasets (Tables S21-23). The increase in association strength did not appear to occur uniformly  
256 across the genome, but rather, at loci associated with age of first occurrence or associated with  
257 specific life periods pertinent to a disease. For example, MTAG of asthma in the SR dataset re-  
258 captured some of the previously reported childhood-onset-specific loci not seen in GWAS,  
259 including a damaging missense variant on *TESPA1*<sup>7</sup> (12:55368291:C:T, MTAG  $P = 4.4 \times 10^{-10}$ ;  
260 Figure 4A).



**Figure 4. MTAG and PRS analysis in the UKBB SR dataset**

- A.** A Manhattan plot of asthma MTAG that incorporates age of first occurrence information (blue triangle) compared to the original case-control GWAS (grey circle). The solid red line denotes  $P = 5 \times 10^{-8}$ .
- B.** Improvement in disease risk prediction using MTAG-PRS versus GWAS-PRS. MTAG was performed for diseases with a significant  $r_g$  between age of first occurrence and susceptibility. PRS and the proportion increase in prediction  $R^2$  of MTAG relative to GWAS (y-axis) were computed in an independent sample of 91K EUR individuals. Results for the HIP and COMB datasets are shown in Figure S20.
- C.** The application of MTAG-PRS and GWAS-PRS in risk stratification for two selected disease phenotypes. The left panel shows the adjusted odds ratio and its 95% CI (y-axis) comparing individuals in each of the top PRS percentiles (x-axis) to the rest of the population. Showing on the right is the corresponding disease prevalence in the top PRS percentiles computed using either GWAS or MTAG summary statistics. Full results for all three UKBB datasets can be found in Tables S24-26.

262

263 To explore how age of first occurrence could aid in risk prediction, we constructed

264 polygenic risk score (PRS) using GWAS and MTAG summary statistics in the left-out set of 91,436

265 ancestry-matched individuals in UKBB. We calculated the relative change in prediction accuracy

266 of MTAG versus GWAS on top of a covariate-only model (Methods). For diseases with an  
267 adequate sample size (>2000 cases in the target sample), the MTAG model that incorporated  
268 age of first occurrence outperformed the susceptibility-only GWAS in predicting disease risk for  
269 many traits. For SR conditions the increase in predictive accuracy can be as high as 30% and on  
270 average around 10%, while the improvement was overall less obvious for HIP endpoints. (Figures  
271 4B&S20; Tables S24-26). Notably, diseases with a larger increase in prediction  $R^2$  seemed to  
272 involve several childhood-onset, allergy-related conditions (Figures 4B).

273 Finally, we evaluated the clinical utility of PRS in patient stratification using MTAG-PRS  
274 compared to GWAS-PRS (Methods). In both models, individuals in the top PRS percentiles (1%,  
275 2.5%, 5%, 10%, and 20%) had a significantly elevated disease risk versus those who were not.  
276 Furthermore, for many traits, MTAG-PRS consistently identified individuals at a higher risk than  
277 predicted by GWAS, corresponding to a larger proportion of cases across top PRS percentiles  
278 (e.g., eczema, high cholesterol; Figures 4C&S21; Tables S24-S26). For instance, individuals in  
279 the top 1% GWAS-PRS were at a 2.2-fold increased risk for SR eczema compared to the average  
280 PRS group ( $P = 4.0 \times 10^{-7}$ ) and a 2.3-fold risk than the rest of the sample ( $P = 1.7 \times 10^{-8}$ ), while  
281 those defined by the top 1% MTAG-PRS had an even higher risk of OR = 2.5 ( $P = 1.9 \times 10^{-11}$ ) and  
282 OR = 2.9 ( $P = 1.2 \times 10^{-15}$ ), respectively (Figures 4C&S21). Overall, risk prediction and stratification  
283 incorporating age of first occurrence information showed more gains for SR phenotypes than  
284 defined in the other two datasets, particularly for conditions that tend to have a pediatric population.

285

## 286 Discussion

287 Age of onset has long been an important phenotype of interest in epidemiological studies due to  
288 its rich information on timing of events relevant to disease etiology, but a systematic  
289 characterization of its genetic underpinnings has been lacking. Here, we provide a deep  
290 investigation into the genetic architecture of age of first occurrence of a wide range of complex  
291 disorders among White-British individuals in UKBB using three different phenotypic criteria (SR,  
292 HIP, COMB). We discovered that age of first occurrence has non-trivial genetic contributions  
293 across many diseases, some with unique genetic risk factors not associated with simply disease  
294 susceptibility. Although the extent of sharing varies by disease, genetic correlation analysis  
295 suggests that an earlier health-event occurrence indexes a heavier polygenic burden of the  
296 disease. Independent of UKBB, genetic investigation of the FinnGen cohort yields a similar  
297 pattern of heritability and genetic correlation estimates. Further, we demonstrate that  
298 incorporating onset age with susceptibility could improve risk loci discovery and patient  
299 stratification based on genetic risk prediction.

300 The scan of SNP-heritability of age of first occurrence ( $h_{aof}^2$ ) revealed a similar landscape  
301 of what has been observed with SNP-heritability of susceptibility to complex disorders ( $h_{sus}^2$ )<sup>20</sup>.  
302 We note that  $h_{aof}^2$  estimates the degree to which genetic variation explains age at first event  
303 among affected individuals whereas  $h_{sus}^2$  is defined at the *population* level and estimates the  
304 genetic contribution to a continuous disease liability. Notably, our use of different criteria to define  
305 age of onset of a disease and the inherent variability in the manifestation or diagnosis of a disease  
306 led to heterogeneous  $h_{aof}^2$  estimates between datasets, more so than variability among  $h_{sus}^2$   
307 estimates (Figure S10). Aside from the sample size constraint, the extent to which  $h_{aof}^2$  varied  
308 by source of definition is disease-dependent and may relate to health-reporting and health-  
309 seeking behaviors as well as the time span of each data type. The SR dataset serves as a good  
310 source for documenting the approximate age of onset for health events, including those starting  
311 early in life, but might suffer from recall bias, particularly for diseases with a mild manifestation.

312 On the other hand, for conditions that require in-patient stays, such as an acute onset (e.g., heart  
313 failure, diabetes complications), HIP phenotypes may be a more valid source. The downside of  
314 HIP-defined event occurrence is that it might not capture the exact age of onset for many diseases,  
315 including occurrences prior to when the ICD-10 system came into use and incorporated into the  
316 Health Episode Statistics database in the UK (1990s), and ICD codes can be given to patients for  
317 administrative purposes. The replication of our results in FinnGen, where registry data covers an  
318 older time period, partially addresses this limitation. Our analysis showed that SR and HIP  
319 endpoints often exhibit a differential pattern of prevalence and  $h_{aof}^2$  (Figure S11). The third  
320 approach of combining different sources to define the earliest event (e.g., UKBB-COMB and  
321 FinnGen) is most effective in increasing sample size, which is ideal for case-control GWAS but  
322 should be interpreted with caveats for age of onset GWAS. As different data sources can start at  
323 a different point in time and not all individuals are included in all datasets (e.g., Figure S16), this  
324 approach could result in a synthetic or multi-modal distribution of disease onset age not reflective  
325 of the actual timing of events. In short, there is unlikely to be one simple measure to define age  
326 of onset, and the heterogeneity in case ascertainment methods can affect the findings and  
327 interpretation of genetic analysis with respect to time.

328 Our findings suggest that the genetic basis of age of first occurrence and susceptibility—  
329 two seemingly orthogonal phenotypic components of a disease—have a correlated nature for  
330 many complex disorders. The shared and distinct loci between age of first occurrence and  
331 susceptibility indicate that some variants affect both disease risk and disease onset age while  
332 others are specifically associated with delaying or accelerating disease occurrence. The growth  
333 of biobank-scale datasets may reveal more genetic risk factors that can modify age of onset and  
334 may represent new therapeutic targets. On the genome-wide scale, our observation that earlier  
335 onset is associated with a greater polygenic loading for disease is consistent with what has long  
336 been hypothesized for polygenic disorders and has been increasingly observed in individual  
337 GWAS studies<sup>13,14,21</sup>. For Mendelian disorders, studies have also shown that its polygenic

338 background can modify the age of disease occurrence among those with rare monogenic  
339 mutations<sup>22,23</sup>. Through genetic correlation analysis in the UKBB and FinnGen cohorts, we  
340 demonstrated that such an inverse genomic relationship between when and whether a disease  
341 occurs is in fact widespread among complex disorders. While the interpretation will vary by  
342 disease, a negative  $r_g$  might imply a heterogeneous genetic architecture at different ages of onset,  
343 suggesting age-related disease subtypes or a continuum whereby earlier onset of the disease is  
344 more genetically driven while later onset might reflect a greater contribution of non-genetic life  
345 events.

346 As an alternative to our case-only approach, Cox proportional hazard modeling, a common  
347 method to study time-to-event endpoints that models both disease status (whether) and onset  
348 time (when) accounting for censoring events, can be a powerful approach for detecting SNP-  
349 disease associations<sup>24,25</sup>. However, such models do not allow explicit assessment of genetic risk  
350 factors underlying age of first occurrence separately from disease susceptibility, and hence  
351 estimation of their genetic overlap. In addition, heritability from the Cox model is typically defined  
352 for cumulative hazards<sup>25,26</sup> and does not have the same interpretation as  $h_{afo}^2$ . As more scalable  
353 survival models for GWAS are being developed<sup>25,27</sup>, a comparison of the two approaches will be  
354 informative for distinguishing their genetic findings and implications for the dissecting genetic  
355 architecture.

356 Through PRS analysis, we showed that genetic variation underlying age of first occurrence  
357 of health events can improve disease risk prediction and patient stratification, especially for  
358 diseases that tend to have an early onset. This demonstrates how axes of information correlated  
359 with case-control status can be useful in planning prevention and intervention strategies tailored  
360 toward individual genetic predisposition. The rich spectrum of phenotypic information in large  
361 biobank cohorts provides an unparalleled opportunity for epidemiological and genetic research to  
362 study clinical features beyond simply disease susceptibility through longitudinal healthcare  
363 records. Genome-wide analysis of such quantitative disease dimensions is critical and will

364 continue to generate valuable insights into the genetic basis of disease development, severity,  
365 and progression.  
366



367 **Online Methods**

368 ***UK Biobank, genotyping, and quality control***

369 UK Biobank (UKBB) is a population-based cohort with extensive phenotype and genomic data on  
370 more than 500,000 individuals in the United Kingdom aged 40 to 69 years at recruitment<sup>15,16</sup>. The  
371 study started recruitment between 2006 and 2010 and has been followed up prospectively. The  
372 genetic data in the UK Biobank is available for a subset of 488,377 participants, with 49,950  
373 individuals genotyped using the UK BiLEVE Axiom Array and the other 438,427 participants  
374 genotyped using the UK Biobank Axiom Array in GRCh37 coordinates. The Haplotype Reference  
375 Consortium (HRC)<sup>28</sup> data and the merged UK10K and 1000 Genomes phase 3 data<sup>29</sup> were used  
376 as reference panels for imputation.

377 We obtained the lists of post-QC samples and variants from the Neale Lab GWAS for  
378 analysis, which comprised 361,194 unrelated individuals of predominantly White-British descent  
379 and 13.7 million genetic markers. Quality control procedures of the genotype data were detailed  
380 in the Neale Lab blog posts and GitHub repository<sup>20,30</sup>. In brief, the initial UKBB cohort was filtered  
381 to those who were unrelated and did not have sex chromosome aneuploidy using the provided  
382 UKBB sample QC metrics. Among them, individuals who were self-reported “White-British”, “Irish”,  
383 or “White”, and fell within seven standard deviations of the first six principal components (PCs)  
384 were retained. Variants were filtered to those with an imputation INFO score > 0.8, a minor allele  
385 frequency (MAF) > 0.1% and a p-value for the Hardy-Weinberg equilibrium test >  $1 \times 10^{-10}$ .

386 Here, we further removed participants who have since withdrawn consent and variants  
387 that were located on the sex chromosomes or had a MAF < 1%. The final dataset consisted of  
388 361,140 individuals and 9.4 million autosomal, common variants. We then converted the data  
389 from .bgen to .pgen format using PLINK 2.0<sup>31</sup>, which preserved dosage information for association  
390 analysis.

391

392 ***Age of first occurrence: phenotype definition and selection***

393 Age of first occurrence of a disease or medical condition was estimated from three different  
394 phenotypic datasets in the UKBB for the 360K White-British subset, including *self-reported*  
395 medical conditions (SR), *hospital in-patient* records (HIP), and the *combined* first occurrence data-  
396 fields of diagnostic codes that mapped across different phenotypic datasets (COMB). Although  
397 the classification schemes are different, both SR and HIP datasets follow a tree-structure topology  
398 to define disease endpoints.

399 For the SR dataset, we aggregated all parent nodes and their children nodes from medical  
400 conditions ascertained through touchscreen questionnaires and verbal interviews (data-fields  
401 20001 and 20002) into a total of 562 clinical terms across 11 non-cancer disease classes and  
402 106 terms across 9 cancer categories. For each SR medical condition, participants were given  
403 the option to report either year or age when first diagnosed by a doctor in integers and the value  
404 was rounded to the nearest quarter age (data-fields 20007 and 20009). We identified affected  
405 individuals for each term and extracted their interpolated *age of diagnosis*. The earliest age of  
406 diagnosis among all children nodes was considered when the trait of interest was a parent node.  
407 Individuals who had a missing age of diagnosis—either uncertain/unknown or preferred not to  
408 answer—were excluded from the analysis.

409 The HIP dataset was based on the Health Episode Statistics database in the UK and  
410 contained hospitalization episodes for each participant in the form of International Classification  
411 of Diseases (ICD) classifications, predominantly in ICD-10. Each in-patient episode had a  
412 corresponding primary diagnosis and might be associated with one or more secondary diagnoses;  
413 dates of admission, episode-start, episode-end, and discharge were also recorded. The earliest  
414 dates in HIP can be traced back to 1990s, roughly when the ICD-10 system came into place  
415 ([https://biobank.ctsu.ox.ac.uk/showcase/exinfo.cgi?src=Data\\_providers\\_and\\_dates](https://biobank.ctsu.ox.ac.uk/showcase/exinfo.cgi?src=Data_providers_and_dates)). We  
416 mapped the available ICD-10 codes to PheWAS Codes (PheCodes)<sup>32,33</sup> and defined all possible  
417 PheCode-based phenotypes based on its hierarchical structure, ranging from individual

418 PheCodes, comorbid medical conditions with related PheCodes, to the entire PheCode category.  
419 Together, this resulted in 1,843 cleanly defined clinical terms across 18 PheCode categories. For  
420 each term, individuals who had at least one associated ICD-10 codes in the primary or secondary  
421 diagnoses were considered a case, and their *age of first hospitalization episode, or in-patient*  
422 *diagnosis*, was calculated as the interval between the earliest episode-start date and the mid-  
423 month of their birth year. Where episode-start date was missing, admission date was used instead.

424 The COMB dataset is a special data type that contains the first-reported dates of a  
425 diagnostic code for a range of different non-cancer health outcomes, generated by mapping  
426 across SR, HIP, primary care, and death records. Note that the current release of primary care  
427 data includes only 45% of the UKBB cohort. For data types not based on the ICD-10 classification  
428 system, health outcomes were first mapped to a related 3-character ICD-10 code where  
429 appropriate by the UKBB team, and its date of first occurrence was recorded as the earliest among  
430 the four different datasets. The mapping mechanism however has not been externally validated  
431 and many of the SR endpoints did not have a corresponding ICD-10 code. We extracted all the  
432 available first occurrence fields, totaling 1,165 terms across 16 disease categories (data-fields  
433 2401-2417), set to missing any improbable dates for each term (e.g., an event date before birth,  
434 at birth, or in the future), and estimated the number of affected individuals as those with a non-  
435 missing first occurrence date. Due to the composite nature of COMB endpoints, we did not map  
436 these 3-character ICD-10 codes to PheCodes. Age when a given condition was first reported was  
437 estimated as described for the HIP dataset.

438 We noted that age first occurrence estimated in the SR and the COMB datasets could be  
439 as early as <1 year of age, and the small values seemed unlikely for some diseases. However,  
440 such instances only accounted for a tiny fraction of the total affected individuals, particularly  
441 among adult-onset conditions (Tables S4&S6). Imposing any threshold to truncate the distribution  
442 would seem arbitrary given that the biologically plausible age of onset differs widely across

443 diseases. We instead acknowledged that measurement error or misclassification may exist for  
444 these small values but their impact on analysis should be limited as a whole.

445 To compare the distribution of age of first occurrence across data types, we selected 26  
446 phenotypes whose definitions could be well mapped between SR and HIP, 15 of which were  
447 comparable across all three datasets (hypertension, angina, myocardial infarction, heart  
448 arrhythmia, hypothyroidism, high cholesterol, diabetes, asthma, depression, migraine,  
449 osteoarthritis, disc problem, gastroesophageal reflux disease/GERD, hiatus hernia,  
450 cholelithiasis/gallstones, breast cancer, skin cancer, endocrine/metabolic disorders, psychiatric  
451 disorders, neurological diseases, cardiovascular diseases, respiratory diseases,  
452 gastrointestinal/digestive diseases, dermatologic diseases, musculoskeletal disorders, and  
453 neoplasms) (Table S7).

454

#### 455 ***Association analysis***

456 We performed a GWAS of age of first occurrence for diseases with at least 5000 affected  
457 individuals, a cutoff we considered sufficiently powered for genetic analysis. This led to a total of  
458 70 SR, 224 HIP, and 164 COMB medical conditions (Tables S1-3). Age of first occurrence was  
459 analyzed as a quantitative outcome in a linear regression to estimate its association with imputed  
460 SNP dosages in PLINK 2.0 (PLINK v2.00a2LM), adjusting for sex, genotyping array, and the first  
461 20 PCs. For comparison, a GWAS of susceptibility to the same disease definition was also  
462 performed, which treated “ever affected with the condition” as the endpoint and considered all  
463 non-case individuals as controls; the same sets of covariates were included in the regression  
464 model. We deliberately did not adjust for current age in the model to avoid double counting the  
465 age information. Nonetheless, when current age was included as an additional covariate, we  
466 observed comparable GWAS results and heritability estimates.

467 To obtain independently associated loci, we performed linkage disequilibrium (LD)  
468 clumping on GWAS summary statistics. The clumping procedure started by identifying genome-

469 wide significant SNPs ( $p$ -value  $< 5 \times 10^{-8}$ ) and then selecting any other SNPs that had a  $r^2 > 0.1$   
470 with the index SNP within a 500kb window to form a clump (`--clump-p1 5e-08 --clump-p2 0.05 --`  
471 `clump-r2 0.1 --clump-kb 500`). The procedure stopped when all genome-wide significant SNPs  
472 were assigned to a locus. Any overlapping associated loci were merged using  
473 BEDTools/bedtools<sup>34</sup> (v2.27.1). A randomly selected sample of 10,000 White British individuals  
474 in UKBB was used as the reference panel to compute LD. Due to its strong and extensive LD  
475 structure, the major histocompatibility complex (MHC) region (chr6: 25Mb-35Mb) was treated as  
476 one genomic locus. Finally, associated genomic loci of age of first occurrence and susceptibility  
477 from their respective GWAS were compared using “bedtools intersect” to identify shared or  
478 distinct loci.

479

#### 480 ***Inflation evaluation, SNP-heritability, and genetic correlation***

481 We performed LD Score regression (LDSR)<sup>35</sup> analysis on GWAS summary statistics to assess  
482 the extent of residual confounding and to estimate SNP-heritability. LDSR can distinguish inflation  
483 in GWAS association  $\chi^2$  statistics due to confounding such as population stratification from true  
484 polygenicity. An LDSR intercept close to one, or a ratio of  $(\text{intercept}-1)/(\text{mean } \chi^2-1)$  close to zero,  
485 would indicate the contribution of confounding biases is well-controlled. The analysis was done  
486 using the pre-computed LD scores of 1.2 million high-quality HapMap3 SNPs (excluding the MHC  
487 region) from the European samples in the 1000 Genomes Project<sup>36</sup> and GWAS summary-level  
488 results from the tested UKBB disease phenotypes. SNP-heritability for age of first occurrence was  
489 estimated based on the slope of LDSR to measure the degree to which the phenotypic variation  
490 is explained by common genetic variation. SNP-heritability for susceptibility was estimated on the  
491 liability-scale, assuming that population prevalence equals sample prevalence in the dataset. For  
492 conditions with a significant heritability of its age of first occurrence and susceptibility endpoints,  
493 we calculated a genetic correlation ( $r_g$ ) between the two traits using bivariate LDSR<sup>37</sup>. The genetic  
494 covariance is estimated using the slope from the regression of the product of z-scores from two

495 GWAS studies against the LD score. The estimate obtained from this method represents the  
496 genetic correlation between the two traits attributable to all polygenic effects captured by common  
497 SNPs. In addition, we estimated pairwise  $r_g$ 's separately for age-of-first-occurrence and  
498 susceptibility phenotypes to examine the pattern of genetic sharing across diseases.

499

### 500 ***Replication analysis in FinnGen***

501 To replicate the results of the heritability pattern of age of first occurrence and its genetic  
502 correlation with susceptibility, we analyzed the FinnGen cohort<sup>17</sup> of 130,423 unrelated individuals.  
503 Age at recruitment of FinnGen participants ranged from 0.08 to 98.98 (median: 54.36, IQR: 25.45).  
504 FinnGen is a public-private partnership project combining genotype data from Finnish biobanks  
505 and digital health record data from Finnish health registries (<https://www.finnngen.fi/en>). Six  
506 regional and three country-wide Finnish biobanks participate in FinnGen, which also includes data  
507 from previously established populations and disease-based cohorts.

508 FinnGen disease endpoints are defined using nationwide registries. Data are harmonized  
509 over the ICD revisions 8, 9 and 10, cancer-specific ICD-O-3, (NOMESCO) procedure codes,  
510 Finnish-specific Social Insurance Institute (KELA) drug reimbursement codes, and ATC-codes for  
511 medications. These registries span decades (Figure S16) and are electronically linked to the  
512 cohort baseline data using the unique national personal identification numbers assigned to all  
513 Finnish citizens and residents. We used genotype and phenotype data from FinnGen release v4  
514 of 130,423 unrelated Finnish participants, excluding population outliers via PCA and related  
515 individuals (<3rd degree) using the KING software<sup>38</sup> (Supplementary Notes).

516 We analyzed diseases with at least 5000 cases and a few additional ones with slightly  
517 fewer than 5000 cases that we considered relevant for comparison with UKBB, leading to a total  
518 of 280 medical conditions. For each condition, age of first occurrence was defined as the earliest  
519 age of an event in the registries. Age of first occurrence across the analyzed diseases ranged  
520 from 0 to 105.65 (median: 52.04, IQR: 28.22). fastGWA<sup>39</sup> linear regression model was used for

521 GWAS analysis. Sex, 10 PCs, and genotyping batch were used as covariates. Low quality  
522 variants with missingness rate  $> 0.1$  and variants with  $MAF < 0.0001$  were excluded from the  
523 analysis. Following the same procedures in the UKBB analysis, we identified loci uniquely  
524 associated with age of first occurrence, estimated the SNP-heritability of age of first occurrence  
525 and susceptibility using LDSR, and computed a genetic correlation between them for each  
526 analyzed disease definition. Standard LD scores were used based on the 1000 genomes  
527 reference set, restricting to European populations.

528 In addition to the phenome-wide analysis, we focused on the 26 mapped phenotypes to  
529 compare the results in UKBB and FinnGen (Table S19). To evaluate how the estimated effect  
530 sizes of age of first occurrence in UKBB replicate in FinnGen, we first identified the SNPs present  
531 both in UKBB and FinnGen. For the common SNPs, we then obtained the independently  
532 associated loci in the UKBB performing LD clumping with a p-value threshold of 0.0001. We finally  
533 looked at the correlation of the effect sizes in UKBB and FinnGen for the index SNPs of each  
534 clump and evaluated the concordance of the direction of the effects with a one-sample binomial  
535 test.

536

### 537 ***MTAG analysis of susceptibility and age of first occurrence***

538 For UKBB diseases where age of first occurrence and susceptibility shared a significant genetic  
539 correlation, we meta-analyzed the two traits using Multi-Trait Analysis of GWAS (MTAG)<sup>19</sup>. Built  
540 upon the LDSR framework, MTAG boosts power for loci discovery for a trait by factoring in its  
541 shared genetic architecture with other traits while accounting for sample overlap. MTAG by default  
542 imposes additional SNP filters, which resulted in  $\sim 7.8$  million SNPs, fewer than included in the  
543 GWAS. All the comparisons of MTAG and GWAS were based on this subset of SNPs. The  
544 number of significant loci in MTAG of *susceptibility*—now with the information on age of first  
545 occurrence incorporated—was calculated based on the same clumping and merging procedure

546 described previously and were compared against the original susceptibility GWAS to identify any  
547 additional loci.

548

### 549 ***Polygenic risk prediction of disease susceptibility***

550 To evaluate if genetic sharing with age of first occurrence can improve risk prediction of disease  
551 susceptibility, we constructed polygenic risk scores (PRS) from MTAG association statistics and  
552 compared its performance to that based on the susceptibility GWAS. We used PRS-CS<sup>40</sup>, a  
553 Bayesian polygenic prediction method that imposes a continuous shrinkage prior on SNP effect  
554 sizes and is robust to diverse genetic architectures, to obtain the posterior effect size of each SNP,  
555 using 1000 Genomes European samples as the LD reference panel. PRS was calculated as the  
556 sum of allele dosages weighted by the posterior effect sizes of each SNP in an independent,  
557 ancestry-matched sample of 91,436 UKBB Individuals using PLINK 2.0 (--score). The target  
558 sample was selected as follows: starting with 1000 Genomes variants that overlapped with UKBB  
559 genotyped variants, we filtered to high-quality autosomal SNPs (no strand-ambiguous alleles, not  
560 in long-range LD regions, with a call rate > 0.98 and a MAF > 0.05), and pruned for LD ( $r^2 < 0.2$ )  
561 down to 149,501 nearly independent SNPs. Using the pruned SNPs, we performed PCA on the  
562 2,504 individuals in 1000 Genomes data and then projected the 488,377 UKBB individuals onto  
563 the computed PC space. With the first 6 PCs in 1000 Genomes as the training data, we used the  
564 Random Forest classifier to assign a “super population” label with a prediction probability  $\geq 0.9$   
565 for each UKBB participant (AFR, AMR, EAS, EUR, or SAS). This resulted in 91,436 individuals  
566 who were classified as EUR and were not included in the discovery GWAS. We focused on  
567 diseases with >2,000 cases in the target sample for PRS analysis.

568 To estimate the predictive power of PRS, we calculated an incremental McFadden’s  
569 pseudo- $R^2$  ( $pR^2$ ), comparing a logistic regression model that included PRS and a set of covariates  
570 (age, sex, genotyping array, and the top 10 PCs) to a covariate-only model. The percentage  
571 improvement in the predictive power of MTAG-PRS over GWAS-PRS was computed as



572  $(pR^2_{MTAG}/pR^2_{GWAS} - 1) \cdot 100\%$ . Next, we evaluated the performance of risk stratification of MTAG-  
573 PRS and GWAS-PRS; both PRS were standardized to have a mean of 0 and a standard deviation  
574 of 1. We dichotomized individuals into those who belonged to the top PRS percentile (1%, 2.5%,  
575 5%, 10%, and 20%) versus those who did not, or those among the average percentiles (20-80%).  
576 We then modeled disease risk as a function of the binary PRS indicator as well as age, sex,  
577 genotyping array, and the top 10 PCs as covariates. The proportion of cases (prevalence) was  
578 also identified in each top PRS percentile.  
579

580 **Acknowledgements**

581 We thank Chia-Yen Chen and Raymond Walters for their insightful suggestions and comments.  
582 We would like to acknowledge the participants and investigators of the UK Biobank and the  
583 FinnGen study for their gracious contribution. This research has been conducted using the UK  
584 Biobank Resource under Application Numbers 31063 and 32568. T.G. is supported by the  
585 National Institutes of Health (NIH) grant K99/R00 AG054573. A special thanks to the Neale Lab  
586 UKBB team (Liam Abbott, Sam Bryant, Claire Churchhouse, Andrea Ganna, Daniel Howrigan,  
587 Duncan Palmer, Ben Neale, Raymond Walters, Caitlin Carey, Cotton Seed, Jonathan Bloom, Tim  
588 Poterba, Dan King, Jackie Goldstein, Arcturus Wang, Patrick Schultz, John Compitello, and Jack  
589 Goldsmith) who together contributed to data management, quality control, and analysis of the  
590 UKBB genotype data and made QC metrics and results publicly available to the scientific  
591 community.

592

593 **Author contribution**

594 Y.A.F. conceived the idea of the project through discussions and guidance from T.G., J.W.S. and  
595 B.M.N. Y.A.F. conducted the primary analysis and drafted the manuscript with feedback from all  
596 co-authors. T.G. provided logistic and methodological support for the UK Biobank analysis. M.C.  
597 contributed to FinnGen analysis and writing methods, with A.G. overseeing the analysis. J.W.S.  
598 and B.M.N. co-supervised the work and provided critical interpretation of the results.

599

600 **Competing interests**

601 J.W.S is an unpaid member of the Bipolar/Depression Research Community Advisory Panel of  
602 23andMe, a member of the Leon Levy Foundation Neuroscience Advisory Board and received an  
603 honorarium for an internal seminar at Biogen, Inc. He is PI of a collaborative study of the genetics  
604 of depression and bipolar disorder sponsored by 23andMe for which 23andMe provides analysis  
605 time as in-kind support but no payments. B.M.N. is a member of the Deep Genomics Scientific

606 Advisory Board and serves as a consultant for the Camp4 Therapeutics Corporation, Takeda  
607 Pharmaceutical and Biogen. The remaining authors declare no conflict of interests.

608

609 **Data availability**

610 The UK Biobank data may be obtained through online registration and application  
611 (<https://www.ukbiobank.ac.uk/>). Study details of the FinnGen study cohort may be accessed  
612 through Finnish Biobanks' FinnBB portal ([www.finbb.fi](http://www.finbb.fi)). GWAS summary statistics of age of first  
613 occurrence generated in the study will be made available at a public data repository prior to  
614 publication.

615

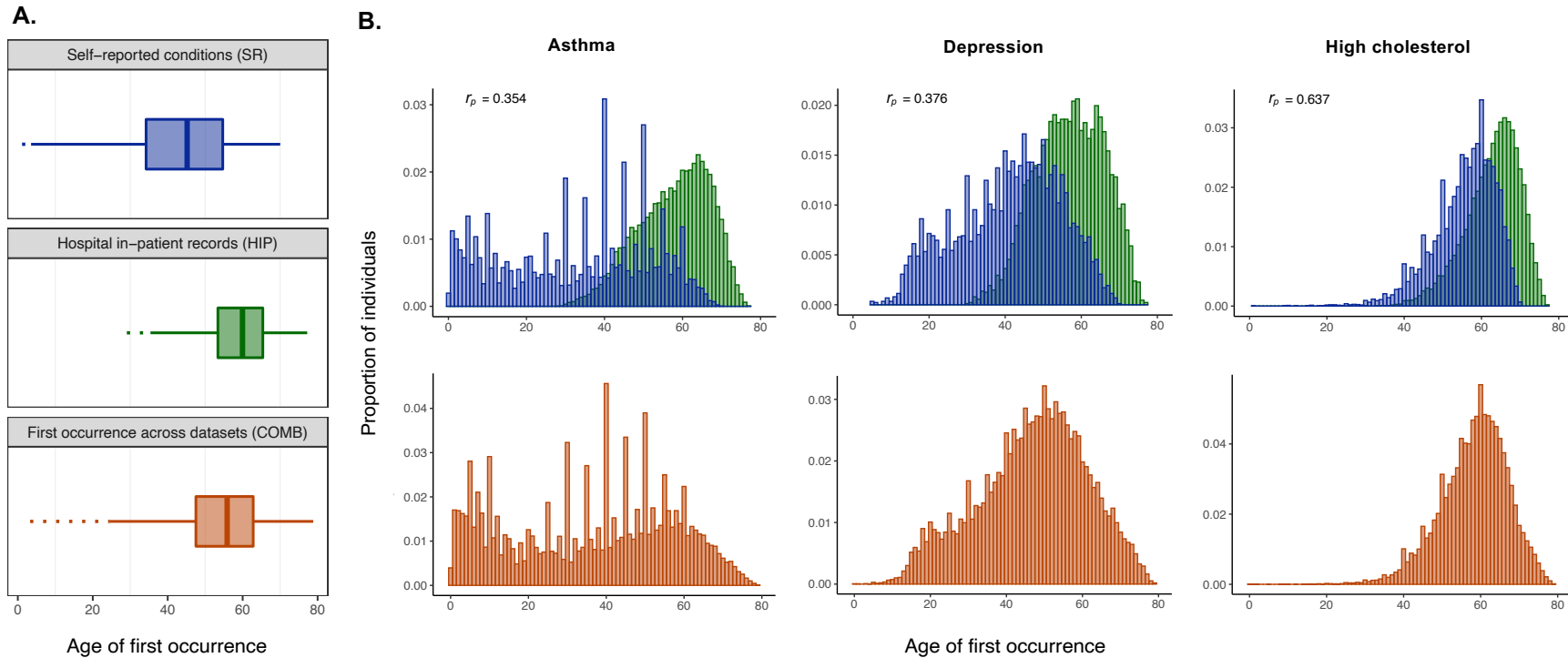
616 **References**

- 617 1. Claussnitzer, M. *et al.* A brief history of human disease genetics. *Nature* **577**, 179-189  
618 (2020).
- 619 2. Visscher, P.M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation.  
620 *Am J Hum Genet* **101**, 5-22 (2017).
- 621 3. Kamboh, M.I. *et al.* Genome-wide association analysis of age-at-onset in Alzheimer's  
622 disease. *Mol Psychiatry* **17**, 1340-6 (2012).
- 623 4. Woolston, A.L. *et al.* Genetic loci associated with an earlier age at onset in multiplex  
624 schizophrenia. *Scientific Reports* **7**, 6486 (2017).
- 625 5. Blauwendraat, C. *et al.* Parkinson's disease age at onset genome-wide association  
626 study: Defining heritability, genetic loci, and  $\alpha$ -synuclein mechanisms. *Mov Disord* **34**,  
627 866-875 (2019).
- 628 6. Ferreira, M.A.R. *et al.* Genetic Architectures of Childhood- and Adult-Onset Asthma Are  
629 Partly Distinct. *Am J Hum Genet* **104**, 665-684 (2019).
- 630 7. Pividori, M., Schoettler, N., Nicolae, D.L., Ober, C. & Im, H.K. Shared and distinct  
631 genetic risk factors for childhood-onset and adult-onset asthma: genome-wide and  
632 transcriptome-wide studies. *Lancet Respir Med* **7**, 509-522 (2019).
- 633 8. Power, R.A. *et al.* Genome-wide Association for Major Depression Through Age at  
634 Onset Stratification: Major Depressive Disorder Working Group of the Psychiatric  
635 Genomics Consortium. *Biol Psychiatry* **81**, 325-335 (2017).
- 636 9. Ge, T., Chen, C.Y., Neale, B.M., Sabuncu, M.R. & Smoller, J.W. Phenome-wide  
637 heritability analysis of the UK Biobank. *PLoS Genet* **13**, e1006711 (2017).
- 638 10. Davis, O.S.P., Haworth, C.M.A. & Plomin, R. Dramatic Increase in Heritability of  
639 Cognitive Development from Early to Middle Childhood: An 8-Year Longitudinal Study of  
640 8,700 Pairs of Twins. *Psychological Science* **20**, 1301-1308 (2009).

- 641 11. Nivard, M.G. *et al.* Stability in symptoms of anxiety and depression as a function of  
642 genotype and environment: a longitudinal twin study from ages 3 to 63 years.  
643 *Psychological Medicine* **45**, 1039-1049 (2015).
- 644 12. Gottesman, II & Shields, J. A polygenic theory of schizophrenia. *Proc Natl Acad Sci U S*  
645 *A* **58**, 199-205 (1967).
- 646 13. Musliner, K.L. *et al.* Association of Polygenic Liabilities for Major Depression, Bipolar  
647 Disorder, and Schizophrenia With Risk for Depression in the Danish Population. *JAMA*  
648 *Psychiatry* **76**, 516-525 (2019).
- 649 14. Mars, N. *et al.* Polygenic and clinical risk scores and their impact on age at onset and  
650 prediction of cardiometabolic diseases and common cancers. *Nat Med* **26**, 549-557  
651 (2020).
- 652 15. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.  
653 *Nature* **562**, 203-209 (2018).
- 654 16. Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a  
655 wide range of complex diseases of middle and old age. *PLoS Med* **12**, e1001779 (2015).
- 656 17. FinnGen. FinnGen Documentation of R3 release  
657 (<https://finngen.gitbook.io/documentation/>). (2020).
- 658 18. Borodulin, K. *et al.* Cohort Profile: The National FINRISK Study. *Int J Epidemiol* **47**, 696-  
659 696i (2018).
- 660 19. Turley, P. *et al.* Multi-trait analysis of genome-wide association summary statistics using  
661 MTAG. *Nat Genet* **50**, 229-237 (2018).
- 662 20. <http://www.nealelab.is/uk-biobank/>.
- 663 21. Ferreira, M.A.R. *et al.* Age-of-onset information helps identify 76 genetic variants  
664 associated with allergic disease. *PLoS Genet* **16**, e1008725 (2020).

- 665 22. Genetic Modifiers of Huntington's Disease (GeM-HD) Consortium. Identification of  
666 Genetic Factors that Modify Clinical Onset of Huntington's Disease. *Cell* **162**, 516-26  
667 (2015).
- 668 23. Fahed, A.C. *et al.* Polygenic background modifies penetrance of monogenic variants for  
669 tier 1 genomic conditions. *Nature Communications* **11**, 3635 (2020).
- 670 24. Staley, J.R. *et al.* A comparison of Cox and logistic regression for use in genome-wide  
671 association studies of cohort and case-cohort design. *Eur J Hum Genet* **25**, 854-862  
672 (2017).
- 673 25. He, L. & Kulminski, A.M. Fast Algorithms for Conducting Large-Scale GWAS of Age-at-  
674 Onset Traits Using Cox Mixed-Effects Models. *Genetics* **215**, 41-58 (2020).
- 675 26. Yazdi, M.H., Visscher, P.M., Ducrocq, V. & Thompson, R. Heritability, reliability of  
676 genetic evaluations and response to selection in proportional hazard models. *J Dairy Sci*  
677 **85**, 1563-77 (2002).
- 678 27. Dey, R. *et al.* An efficient and accurate frailty model approach for genome-wide survival  
679 association analysis controlling for population structure and relatedness in large-scale  
680 biobanks. *bioRxiv* (2020).
- 681 28. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat*  
682 *Genet* **48**, 1279-83 (2016).
- 683 29. Huang, J. *et al.* Improved imputation of low-frequency and rare variants using the UK10K  
684 haplotype reference panel. *Nat Commun* **6**, 8111 (2015).
- 685 30. [https://github.com/Nealelab/UK\\_Biobank\\_GWAS](https://github.com/Nealelab/UK_Biobank_GWAS).
- 686 31. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer  
687 datasets. *Gigascience* **4**, 7 (2015).
- 688 32. Denny, J.C. *et al.* Systematic comparison of phenome-wide association study of  
689 electronic medical record data and genome-wide association study data. *Nat Biotechnol*  
690 **31**, 1102-10 (2013).

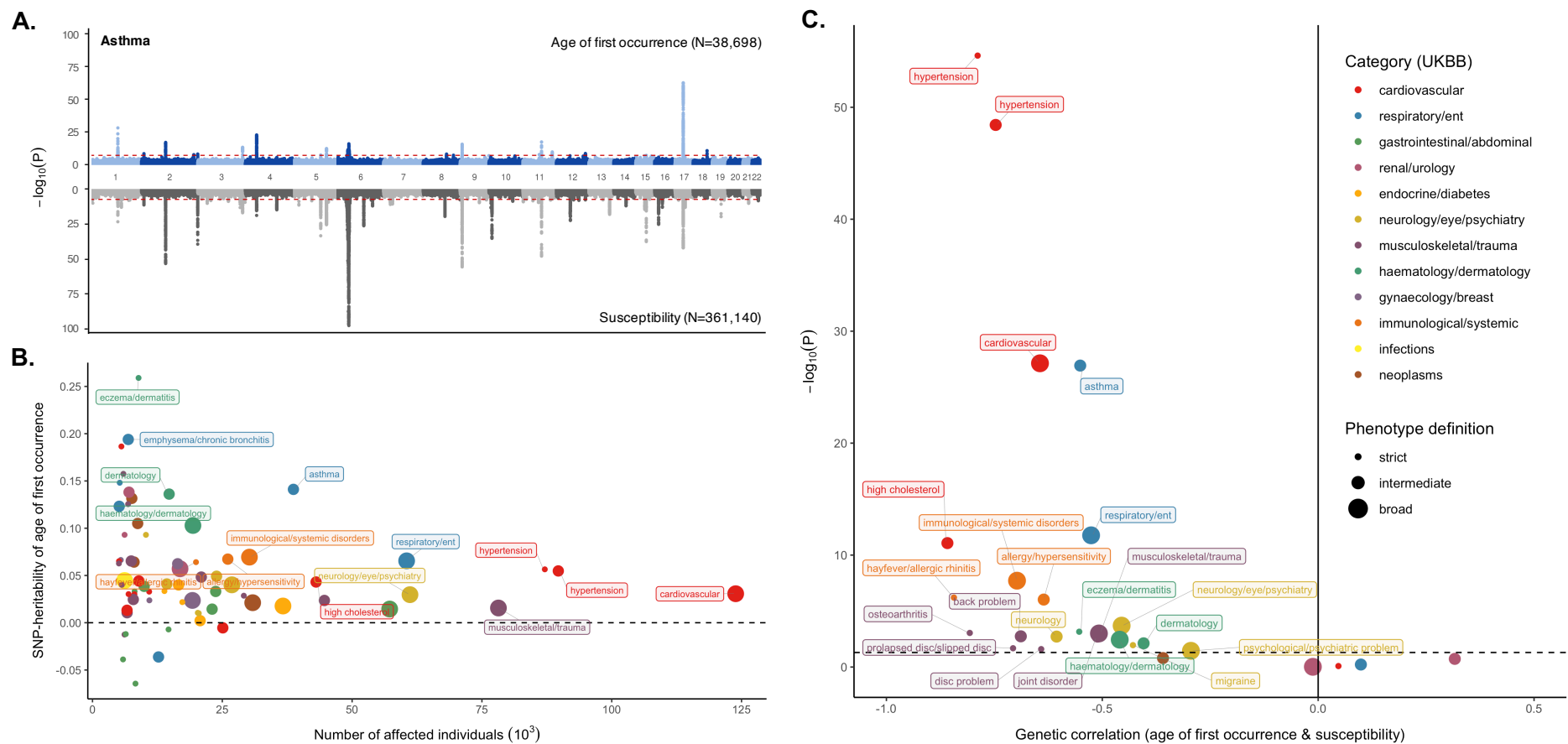
- 691 33. Wu, P. *et al.* Mapping ICD-10 and ICD-10-CM Codes to Phecodes: Workflow  
692 Development and Initial Evaluation. *JMIR Med Inform* **7**, e14325 (2019).
- 693 34. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic  
694 features. *Bioinformatics* **26**, 841-2 (2010).
- 695 35. Bulik-Sullivan, B.K. *et al.* LD Score regression distinguishes confounding from  
696 polygenicity in genome-wide association studies. *Nat Genet* **47**, 291-5 (2015).
- 697 36. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation.  
698 *Nature* **526**, 68-74 (2015).
- 699 37. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and  
700 traits. *Nat Genet* **47**, 1236-41 (2015).
- 701 38. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies.  
702 *Bioinformatics* **26**, 2867-73 (2010).
- 703 39. Jiang, L. *et al.* A resource-efficient tool for mixed model association analysis of large-  
704 scale data. *Nat Genet* **51**, 1749-1755 (2019).
- 705 40. Ge, T., Chen, C.Y., Ni, Y., Feng, Y.A. & Smoller, J.W. Polygenic prediction via Bayesian  
706 regression and continuous shrinkage priors. *Nat Commun* **10**, 1776 (2019).
- 707



**Figure 1. Distribution of age of first occurrence of disease phenotypes from three phenotypic datasets in UKBB**

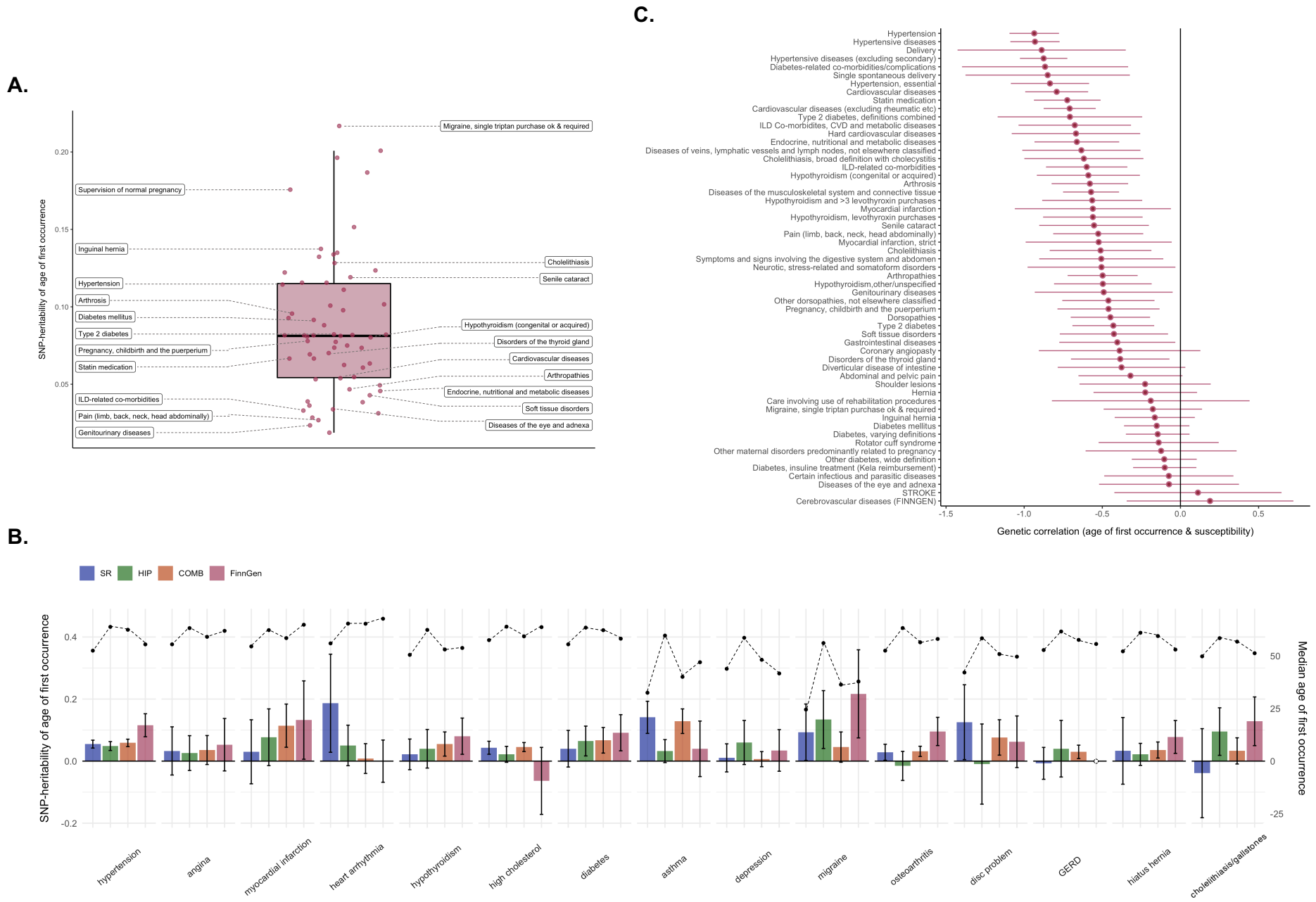
- A.** An averaged distribution of age of first occurrence is shown across 70 SR (blue), 224 HIP (green), and 164 COMB (orange) disease definitions in each dataset. Dotted line indicates the outlying range of values. Age of first occurrence ranges from 0-70 in the SR dataset, 30-80 in the HIP dataset, and 0-80 in the COMB dataset. Spikes in the SR phenotypes reflect that the values are recorded in quartiles (0.25, 0.5, 0.75, or 1.00).
- B.** Distribution of age of first occurrence differs by trait and data source. Shown here are three selected disease phenotypes with matching definitions across datasets. SR and HIP conditions show little to moderate overlap in age of first occurrence, as measured by Pearson's correlation coefficient ( $r_p$ ; top), while the COMB conditions exhibit a merged distribution of SR and HIP (bottom).





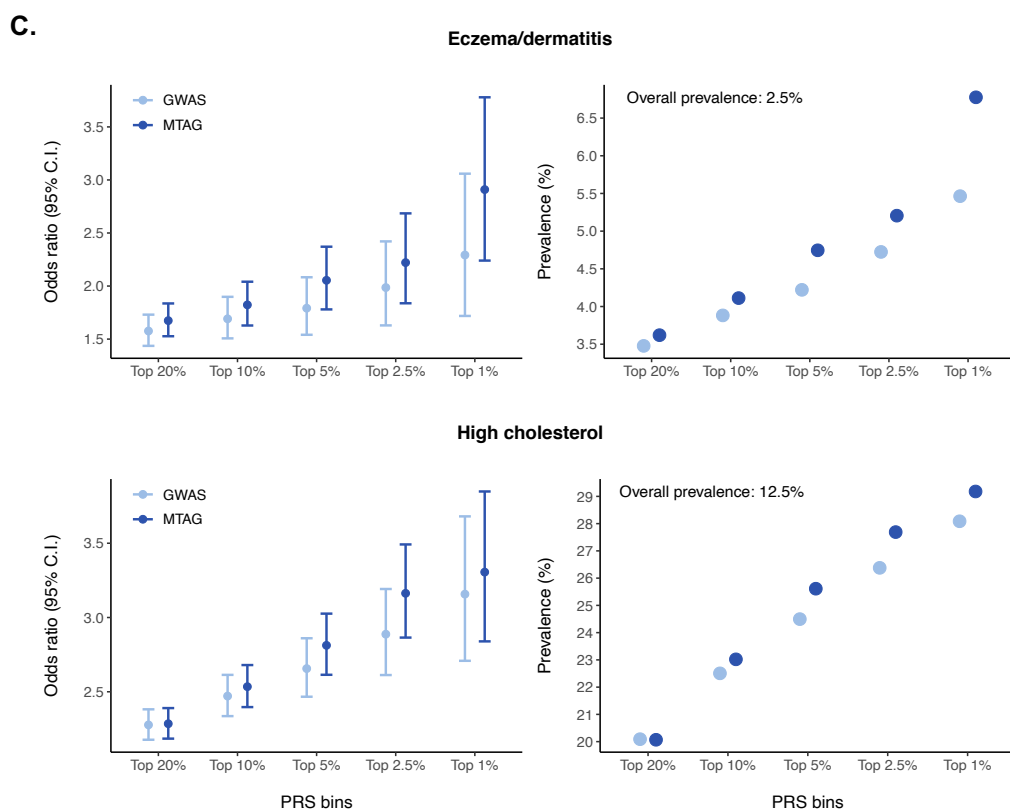
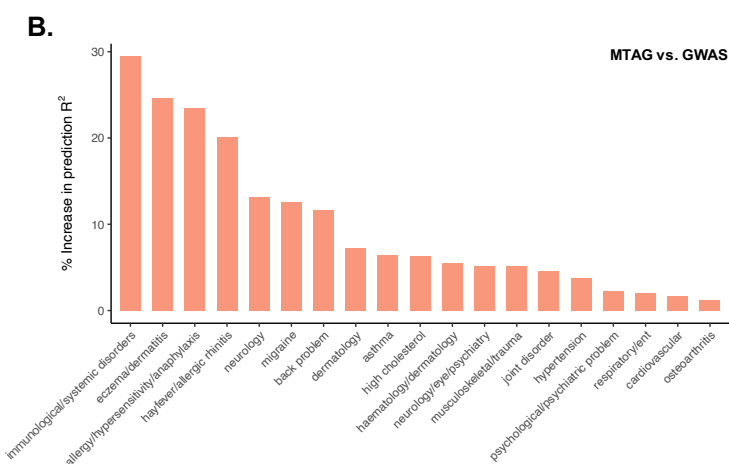
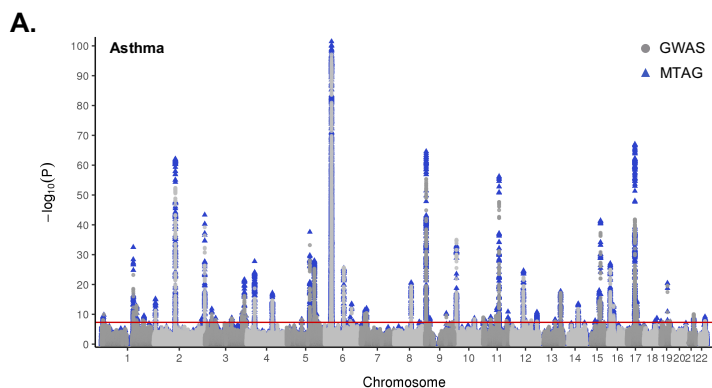
**Figure 2. Genetic characterization of age of first occurrence and its relationship with susceptibility in the UKBB SR dataset**

- A.** A Manhattan plot of GWAS results reveal overlapping and distinct genetic associations between age of first occurrence (top) and case-control status (bottom) of SR asthma. Each dot represents a single SNP. P-values are shown on the  $-\log_{10}$  scale on the y-axis, plotted against chromosome positions on the x-axis. The red dashed lines denote the genome-wide significance threshold at  $P = 5 \times 10^{-8}$ .
- B.** SNP-heritability estimates for age of first occurrence ( $h_{aof0}^2$ ) across 70 SR disease definitions suggest non-trivial common genetic contributions.  $h_{aof0}^2$  was estimated from univariate LDSR. Each dot represents an individual disease, colored by disease categories used in the SR dataset; a larger dot corresponds to a broader disease definition. Labeled are conditions with a significant  $h_{aof0}^2$  at FDR < 0.05. Heritability analysis of HIP and COMB diseases reveal a similar pattern in Figure S9.
- C.** Genetic correlation ( $r_g$ ) analysis suggests an inverse genomic relationship between age of first occurrence and susceptibility for diseases with a significant heritability for both traits.  $r_g$  was estimated using bivariate LDSR. The dashed line denotes nominal significance at  $P = 0.05$ ; labeled are conditions with a significant  $r_g$  at FDR < 0.05. Analysis of HIP and COMB diseases show a similar pattern in Figure S13.



**Figure 3. Genetic analysis of age of first occurrence in FinnGen and its comparison with UKBB results**

- A.** Distribution of heritability estimates from FinnGen for 64 diseases that have a significant  $h_{aof0}^2$ . Labeled are selected conditions with a significant  $h_{aof0}^2$  at FDR < 0.05.
- B.**  $h_{aof0}^2$  estimates for 15 comparable disease definitions in UKBB and FinnGen show variable degree of similarity. Left axis denotes SNP-heritability shown in bar plots and the corresponding 95% confidence intervals (95% CI). Right axis shows the median age of first occurrence for each condition indicated in dotted lines. The full comparison of all 26 matched phenotypes is available in Table S18 and Figure S19.
- C.** A negative  $r_g$  between age of first occurrence and disease susceptibility is observed for many of the tested diseases, consistent with the findings in UKBB. Shown are  $r_g$  estimates and its 95% C.I. for diseases with a significant heritability for both traits.



**Figure 4. MTAG and PRS analysis in the UKBB SR dataset**

**A.** A Manhattan plot of asthma MTAG that incorporates age of first occurrence information (blue triangle) compared to the original case-control GWAS (grey circle). The solid red line denotes  $P = 5 \times 10^{-8}$ .

**B.** Improvement in disease risk prediction using MTAG-PRS versus GWAS-PRS. MTAG was performed for diseases with a significant  $r_g$  between age of first occurrence and susceptibility. PRS and the proportion increase in prediction  $R^2$  of MTAG relative to GWAS (y-axis) were computed in an independent sample of 91K EUR individuals. Results for the HIP and COMB datasets are shown in Figure S20.

**C.** The application of MTAG-PRS and GWAS-PRS in risk stratification for two selected disease phenotypes. The left panel shows the adjusted odds ratio and its 95% CI (y-axis) comparing individuals in each of the top PRS percentiles (x-axis) to the rest of the population. Showing on the right is the corresponding disease prevalence in the top PRS percentiles computed using either GWAS or MTAG summary statistics. Full results for all three UKBB datasets can be found in Tables S24-26.