

1 **Age-dependent topic modelling of comorbidities in UK Biobank identifies** 2 **disease subtypes with differential genetic risk**

3

4 Xilin Jiang^{1,2,3,4,5,6}, Martin Jinye Zhang^{4,7§}, Yidong Zhang^{1,8,9§}, Michael Inouye^{5,6,10,11,12,13}, Chris
5 Holmes^{1,2,14}, Alkes L. Price^{4,7,15*}, Gil McVean^{1*}

6

7 **Affiliations**

8 ¹ Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, University of
9 Oxford, Oxford OX3 7LF, UK

10 ² Department of Statistics, University of Oxford, Oxford OX1 3LB, UK

11 ³ Wellcome Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK

12 ⁴ Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA, USA

13 ⁵ British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public Health and
14 Primary Care, University of Cambridge, Cambridge UK

15 ⁶ Heart and Lung Research Institute, University of Cambridge, Cambridge UK

16 ⁷ Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge,
17 MA, USA

18 ⁸ CAMS China Oxford Institute, Nuffield Department of Medicine, University of Oxford,
19 Oxford OX3 7BN, UK

20 ⁹ Department of Radiation Oncology, Peking Union Medical College Hospital, Chinese
21 Academy of Medical Sciences and Peking Union Medical College, Beijing, China

22 ¹⁰ Cambridge Baker Systems Genomics Initiative, Department of Public Health and Primary
23 Care, University of Cambridge, Cambridge, UK.

24 ¹¹ Health Data Research UK Cambridge, Wellcome Genome Campus and University of
25 Cambridge, Cambridge, UK.

26 ¹² British Heart Foundation Cambridge Centre of Research Excellence, Department of Clinical
27 Medicine, University of Cambridge, Cambridge, UK.

28 ¹³ Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute,
29 Melbourne, VIC, Australia.

30 ¹⁴ The Alan Turing Institute, London NW1 2DB, UK

31 ¹⁵ Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA

32

33 § These authors contributed equally to this work

34 *These authors jointly supervised the work

35 Corresponding authors:

36 xilinjiang@hsph.harvard.edu

37 aprice@hsph.harvard.edu

38 gil.mcvean@bdi.ox.ac.uk

39

40

41 **Abstract**

42 Longitudinal data from electronic health records (EHR) has immense potential to improve
43 clinical diagnoses and personalised medicine, motivating efforts to identify disease subtypes
44 from age-dependent patient comorbidity information. We introduce an age-dependent topic
45 modelling (ATM) method that provides a low-rank representation of longitudinal records of
46 hundreds of distinct diseases in large EHR data sets. The model learns, and assigns to each
47 individual, topic weights for several disease topics, each of which reflects a set of diseases that
48 tend to co-occur as a function of age. Simulations show that ATM attains high accuracy in
49 distinguishing distinct age-dependent comorbidity profiles. We applied ATM to 282,957 UK
50 Biobank samples, analysing 1,726,144 disease diagnoses spanning 348 diseases with $\geq 1,000$
51 incidences. We inferred 10 disease topics optimising model fit. We identified 52 diseases with
52 heterogeneous comorbidity profiles (≥ 500 incidences assigned to each of ≥ 2 topics), including
53 breast cancer, type 2 diabetes (T2D), hypertension, and hypercholesterolemia; for most of these
54 diseases, topic assignments were highly age-dependent, suggesting differences in disease
55 aetiology for early-onset vs. late-onset disease. We defined subtypes of the 52 heterogeneous
56 diseases based on the topic assignments, and compared genetic risk across subtypes using
57 polygenic risk scores (PRS). We identified 18 disease subtypes whose PRS differed significantly
58 from other subtypes of the same disease, including a subtype of T2D characterised by
59 cardiovascular comorbidities and a subtype of asthma characterised by dermatological
60 comorbidities. We further identified specific SNPs underlying these differences. For example,
61 the T2D-associated SNP rs1063192 in the *CDKN2B* locus has a higher odds ratio in the top
62 quartile of cardiovascular topic weight (1.19 ± 0.02) than in the bottom quartile (1.08 ± 0.02)
63 ($P = 4 \times 10^{-5}$ for difference). In conclusion, ATM identifies disease subtypes with differential
64 genome-wide and locus-specific genetic risk profiles.

65
66
67
68
69
70

71 **Introduction**

72 Longitudinal electronic health record (EHR) data, encompassing diagnoses across hundreds of
73 distinct diseases, offers immense potential to improve clinical diagnoses and personalised
74 medicine¹. Despite intense interest in both the genetic relationships between distinct diseases^{2–11}
75 and the genetic relationships between biological subtypes of disease^{12–15}, there has been limited
76 progress on classifying disease phenotypes into groups of diseases with frequent co-occurrences
77 (comorbidities) and leveraging comorbidities to identify disease subtypes. Low-rank modelling
78 has appealing theoretical properties^{16,17} and has produced promising applications^{18–24} to infer
79 meaningful representations of high-dimensional data. In particular, low-rank representation is an
80 appealing way to summarise data across hundreds of distinct diseases^{25–27}, providing the
81 potential to identify patient-level comorbidity patterns and distinguish disease subtypes. Disease
82 subtypes inferred from EHR data could be validated by comparing genetic profiles across
83 subtypes, which is possible with emerging data sets that link genetic data with EHR data^{28–31}.

84
85 Previous studies have used low-rank representation to identify shared genetic components^{25–27}
86 across multiple distinct diseases, identifying relationships between diseases and generating
87 valuable biological insights. However, age at diagnosis information in longitudinal EHR data has
88 the potential to improve such efforts. For example, a recent study used longitudinal disease
89 trajectories to identify disease pairs with statistically significant directionality³², suggesting that
90 age information could be leveraged to infer comorbidity profiles that capture temporal
91 information. In addition, patient-level comorbidity information could potentially be leveraged to
92 identify biological subtypes of disease, complementing its application to increase power for
93 identifying genetic associations¹² and to cluster disease-associated variants into biological
94 pathways⁸; disease subtypes are fundamental to disease aetiology^{14,33–36}.

95
96 Here, we propose an age-dependent topic modelling (ATM) method to provide a low-rank
97 representation of longitudinal disease records. ATM learns, and assigns to each individual, topic
98 weights for several disease topics, each of which reflects a set of diseases that tend to co-occur as
99 a function of age. We applied ATM to 1.7 million disease diagnoses spanning 348 diseases in the
100 UK Biobank, inferring 10 disease topics. We identified 52 diseases with heterogeneous
101 comorbidity profiles that enabled us to define disease subtypes. We used genetic data to validate
102 the disease subtypes, showing that they exhibit differential genome-wide and locus-specific
103 genetic risk profiles.

104

105

106 **Results**

107 *Overview of methods*

108 We propose an age-dependent topic modelling (ATM) model, providing a low-rank
109 representation of longitudinal records of hundreds of distinct diseases in large EHR data sets
110 (Figure 1, Methods). The model assigns to each individual *topic weights* for several *disease*
111 *topics*; each disease topic reflects a set of diseases that tend to co-occur as a function of age,
112 quantified by age-dependent *topic loadings* for each disease. The model assumes that for each
113 disease diagnosis, a topic is sampled based on the individual's topic weights (which sum to 1
114 across topics, for a given individual), and a disease is sampled based on the individual's age and
115 the age-dependent topic loadings (which sum to 1 across diseases, for a given topic at a given
116 age). The model generalises the latent dirichlet allocation (LDA) model^{37,38} by allowing topic
117 loadings for each topic to vary with age (Supplementary Note, Supplementary Figure 1).

118

119 We developed a method to fit this model that addresses several challenges inherent to large EHR
120 data sets; the method estimates topic weights for each individual, topic loadings for each disease,
121 and posterior diagnosis-specific topic probabilities for each disease diagnosis. First, we derived a
122 scalable deterministic method that uses numerical approximation approaches to fit the
123 parameters of the model, addressing the challenge of computational cost. Second, we used the
124 prediction odds ratio³⁹ to compare model structures (e.g. number of topics and parametric form
125 of topic loadings as a function of age), addressing the challenge of appropriate model selection;
126 roughly, the prediction odds ratio quantifies the accuracy of correctly predicting disease
127 diagnoses in held-out patients using comorbidity information, compared to a predictor based
128 only on prevalence (see Methods and Supplementary Table 1). Third, we employed collapsed
129 variational inference⁴⁰, addressing the challenge of sparsity in the data (e.g. in UK Biobank data
130 that we analysed, the average patient has diagnoses for 6 of 348 diseases analysed); collapsed
131 variational inference outperformed mean-field variational inference³⁷ in empirical data. Further
132 details are provided in the Methods section and Supplementary Note; we have publicly released
133 open-source software implementing the method (see Code Availability).

134

135 We applied ATM to longitudinal records of 282,957 individuals from the UK Biobank²⁹,
136 containing a total of 1,726,144 disease diagnoses spanning 348 diseases (see Data Availability).
137 Each disease diagnosis had an associated age at diagnosis, defined as the earliest age of reported
138 diagnosis of the disease in that individual; we caution that age at diagnosis may differ from age
139 at disease onset (see Discussion). ATM does not use genetic data, but we used genetic data to
140 validate the inferred topics (Methods).

141

142 *Simulations*

143 We performed simulations to compare ATM with latent dirichlet allocation (LDA)^{37,38}, a simpler
144 topic modelling approach that does not model age. We simulated 61,000 disease diagnoses
145 spanning 20 diseases in 10,000 individuals, using the ATM generative model; the average

146 number of disease diagnoses per individual (6.1), ratio of #individuals/#diseases (500), topic
147 loadings, and standard deviation in age at diagnosis (8.5 years for each disease) were chosen to
148 match empirical UK Biobank data. We assigned each disease diagnosis to one of two subtypes
149 for the underlying disease based on age and other subtype differences, considering high,
150 medium, or low age-dependent effects by specifying an average difference of 20, 10, or 5 years
151 respectively in age at diagnosis for the two subtypes. For each level of age-dependent effects,
152 we varied the proportion of diagnoses belonging to the first subtype (*subtype sample size*
153 *proportion*) from 10-50%. Further details of the simulation framework are provided in the
154 Methods section. Our primary metric for evaluating the LDA and ATM methods was area under
155 the precision-recall curve (AUPRC)⁴¹, where precision is defined as the proportion of disease
156 diagnoses that a given method assigned to the first subtype that were assigned correctly and
157 recall is defined as the proportion of disease diagnoses truly belonging to the first subtype that
158 were assigned correctly. We discretized the subtype assigned to each disease diagnosis by a
159 given method by assigning the subtype with higher inferred probability. We note that AUPRC is
160 larger when classifying the smaller subtype; results using the second subtype as the classification
161 target are also provided. We used AUPRC (instead of prediction odds ratio) in our simulations
162 because the underlying truth is known. Further details and justifications of metrics used in this
163 study are provided in the Methods section and Supplementary Table 1.

164
165 In simulations with high age-dependent effects, ATM attained much higher AUPRC than LDA
166 across all values of subtype sample size proportion (AUPRC difference: 24%-42%), with both
167 methods performing better at more balanced ratios (Figure 2, Supplementary Table 2).
168 Accordingly, ATM attained both higher precision and higher recall than LDA (Supplementary
169 Figure 2). Results were qualitatively similar when using the second subtype as the classification
170 target (Supplementary Figure 3). In simulations with medium or low age-dependent effects,
171 ATM continued to outperform LDA but with smaller differences between the methods. In
172 simulations without age-dependent effects, ATM slightly underperformed LDA (Supplementary
173 Figure 4A).

174
175 We performed three secondary analyses. First, we varied the number of individuals, number of
176 diseases, or number of disease diagnoses per individual. ATM continued to outperform LDA in
177 each case, although increasing the number of individuals or the number of disease diagnoses per
178 individual did not always increase AUPRC (Supplementary Figure 4B). Second, we performed
179 simulations in which we increased the number of subtypes from two to five and changed the
180 number of diseases to 50, and compared ATM models trained using different numbers of topics
181 (in 80% training data) by computing the prediction odds ratio; we used the prediction odds ratio
182 (instead of AUPRC) in this analysis both because it is a better metric to evaluate the overall
183 model fit to the data, and because it is unclear how to compare AUPRC across scenarios of
184 varying topic numbers (see Supplementary Table 1). We confirmed that the prediction odds
185 ratio was maximised using five topics, validating the use of the prediction odds ratio for model

186 selection (Supplementary Figure 5A). Third, we computed the accuracy of inferred topic
187 loadings, topic weights, and grouping accuracy (defined as proportion of pairs of diseases truly
188 belonging to the same topic that ATM correctly assigned to the same topic), varying the number
189 of individuals and number of diseases diagnoses per individual. We determined that ATM also
190 performed well under these metrics (Supplementary Figure 5B-E).

191
192 We conclude that ATM (which models age) assigns disease diagnoses to subtypes with higher
193 accuracy than LDA (which does not model age) in simulations with age-dependent effects. We
194 caution that our simulations largely represent a best-case scenario for ATM given that the
195 generative model and inference model are very similar (although there are some differences, e.g.
196 topic loadings were generated using a model different from the inference model), thus it is important
197 to analyse empirical data to validate the method.

198

199 *Age-dependent disease topic loadings capture comorbidity profiles in the UK Biobank*

200 We applied ATM to longitudinal records of 282,957 individuals from the UK Biobank²⁹. We
201 used Phecode⁴² to define 1,726,144 disease diagnoses spanning 348 diseases with at least 1,000
202 diagnoses each; the average individual had 6.1 disease diagnoses, and the average disease had a
203 standard deviation of 8.5 years in age at diagnosis. The optimal ATM model structure included
204 10 topics and modelled age-dependent topic loadings for each disease as a spline function with
205 one knot (see below). We assigned names (and corresponding acronyms) to each of the 10
206 inferred topics based on the Phecode systems⁴² assigned to diseases with high topic loadings
207 (aggregated across ages) for that topic (Table 1, Supplementary Table 3).

208

209 Age-dependent topic loadings across all 10 topics and 348 diseases (stratified into Phecode
210 systems), summarised as averages across age<60 and age≥60, are reported in Figure 3,
211 Supplementary Figure 6, and Supplementary Table 4. Some topics such as NRI span diseases
212 across the majority of Phecode systems, while other topics such as ARP are concentrated in a
213 single Phecode system. Conversely, a single Phecode system may be split across multiple topics,
214 e.g. the digestive system is split across UGI, LGI, and MDS. We note that topic loadings in
215 diseases that span multiple topics are heavily age-dependent. For example, type 2 diabetes
216 patients assigned to the CVD topic are associated with early onset of type 2 diabetes whereas
217 type 2 diabetes patients assigned to MGND topic are associated with late onset of type 2
218 diabetes.

219

220 We performed seven secondary analyses to validate the inferred comorbidity topics. First, we fit
221 ATM models with different model structures using 80% training data, and computed their
222 prediction odds ratios using 20% testing data. The ATM model structure with 10 topics and age-
223 dependent topic loadings modelled as a spline function performed optimally (Supplementary
224 Figure 7; see Methods). Second, we confirmed that ATM attained higher prediction odds ratios
225 than LDA across different values of the number of topics (Supplementary Figure 8). Third, we

226 reached similar conclusions using evidence lower bounds³⁹ (ELBO; see Supplementary Table 1)
227 when fitting the model without splitting training and testing data (Supplementary Figure 9).
228 Fourth, we confirmed that collapsed variational inference⁴⁰ outperformed mean-field variational
229 inference³⁷ (Supplementary Figure 10). Fifth, we computed a co-occurrence odds ratio
230 evaluating whether diseases grouped into the same topic by ATM in the training data have higher
231 than random probability of co-occurring in the testing data (Supplementary Table 1). The co-
232 occurrence odds ratio is consistently above one and increases with the number of comorbid
233 diseases, for each inferred topic (Supplementary Figure 11). Sixth, we compared the topic
234 loadings by repeating the inference on female-only or male-only populations and observed no
235 major discrepancies, except for genitourinary topics MGND and FGND (topic loading R^2
236 (female vs. all) = 0.788, topic loading R^2 (male vs. all) = 0.773, Supplementary Figure 12).
237 Lastly, we verified that BMI, sex, Townsend deprivation index, and birth year explained very
238 little of the information in the inferred topics (Supplementary Table 3).

239
240 Disease topics capture known biology as well as the age-dependency of comorbidities for the
241 same diseases. For example, early onset of essential hypertension is associated with the CVD
242 topic⁴³, which captures the established connection between lipid dysfunction
243 (“hypercholesterolemia”) and cardiovascular diseases⁴⁴, while later onset of essential
244 hypertension is associated with the CER topic, which pertains to type 2 diabetes, obesity and
245 COPD (Figure 4A). Continuously varying age-dependent topic loadings for all 10 topics,
246 restricted to diseases with high topic loadings, are reported in Supplementary Figure 13 and
247 Supplementary Table 5. We note that most diseases have their topic loadings concentrated into a
248 single topic (Figure 4B, Supplementary Figure 14A and Supplementary Table 4), and that most
249 individuals have their topic weights concentrated into 1-2 topics (Figure 4C and Supplementary
250 Figure 14B). For diseases spanning multiple topics (Supplementary Figure 6 and Supplementary
251 Table 4), the assignment of type 2 diabetes patients to the CVD topic is consistent with known
252 pathophysiology and epidemiology^{45,46} and has been shown in other comorbidity clustering
253 studies, e.g. with the Beta Cell and Lipodystrophy subtypes described in ref.³⁵ and the severe
254 insulin-deficient diabetes (SIDD) subtype described in ref.¹⁴, which are characterised by early
255 onset of type 2 diabetes and have multiple morbidities including hypercholesterolemia,
256 hyperlipidemia, and cardiovascular diseases⁴⁷. In addition, early-onset breast cancer and late-
257 onset breast cancer are associated with different topics, e.g. NRI and FGND, consistent with
258 known treatment effects for breast cancer patients which increase susceptibility to infections,
259 especially bacterial pneumonias⁴⁸ and hypothyroidism⁴⁹

260
261 We conclude that ATM identifies latent disease topics that robustly compress age-dependent
262 comorbidity profiles and capture disease comorbidities both within and across Phecode systems.

263
264

265 *Disease subtypes defined by distinct topics are genetically heterogeneous*

266 We sought to define disease subtypes based on the diagnosis-specific topic probabilities of each
267 disease diagnosis. We assigned a discrete topic assignment to each disease diagnosis based on its
268 maximum diagnosis-specific topic probability, and defined the disease subtype of each disease
269 diagnosis based on the topic assignment. We restricted our disease subtype analyses to 52
270 diseases with at least 500 diagnoses assigned to each of two distinct subtypes (Methods,
271 Supplementary Figure 6, Supplementary Figure 15, and Supplementary Table 6).

272
273 Age-dependent distributions of subtypes (topics) for four diseases (type 2 diabetes, asthma,
274 hypercholesterolemia, and essential hypertension) are reported in Figure 5A and Supplementary
275 Table 7; results for all 52 diseases are reported in Supplementary Figure 16 and Supplementary
276 Table 7. The number of subtypes can be large, e.g. six subtypes for essential hypertension.
277 Subtypes are often age-dependent, e.g. for the CVD and MGND subtypes of type 2 diabetes^{14,35}
278 (discussed above).

279
280 ATM and the resulting subtype assignments do not make use of genetic data. However, we used
281 genetic data to assess genetic heterogeneity across inferred subtypes of each disease. We first
282 assessed whether polygenic risk scores (PRS) for overall disease risk varied across subtypes of
283 each disease; PRS were computed using BOLT-LMM with five-fold cross validation^{50,51} (see
284 Methods and Code Availability). Results for four diseases (from Figure 5A) are reported in
285 Figure 5B and Supplementary Table 8; results for all 10 well-powered diseases (10 of 52
286 diseases with highest z-scores for nonzero SNP-heritability) are reported in Supplementary
287 Figure 17 and Supplementary Table 8. We identified 18 disease-topic pairs (of 100 disease-topic
288 pairs analysed) for which PRS values in disease cases vary with patient topic weight. For
289 example, for essential hypertension, hypercholesterolemia, and type 2 diabetes, patients assigned
290 to the CVD subtype had significantly higher PRS values than patients assigned to other subtypes.
291 For essential hypertension, patients assigned to the CER subtype had significantly higher PRS
292 values; for type 2 diabetes, patients assigned to the CER subtype had lower PRS values than the
293 CVD subtype, even though the majority of type 2 diabetes diagnoses are assigned to the CER
294 subtype. We further verified that most of the variation in PRS values with disease subtype could
295 not be explained by age⁵² or differences in subtype sample size (Supplementary Figure 18).
296 These associations between subtypes (defined using comorbidity data) and PRS (defined using
297 genetic data) imply that disease subtypes identified through comorbidity are genetically
298 heterogeneous, consistent with differences in disease aetiology.

299
300 We further investigated whether subtype assignments (defined using comorbidity data) revealed
301 subtype-specific excess genetic correlations. We estimated excess genetic correlations between
302 pairs of disease subtypes (relative to genetic correlations between the underlying diseases).
303 Excess genetic correlations for 15 disease subtypes (spanning 11 diseases and 3 topics: CER,
304 MGND and CVD) are reported in Figure 6A and Supplementary Table 9 (relative to genetic
305 correlations between the underlying diseases; Figure 6B), and excess genetic correlations for all

306 89 well-powered disease subtypes (89 of 378 disease subtypes with z-score > 4 for nonzero SNP-
307 heritability) are reported in Supplementary Figure 19 and Supplementary Table 9. Genetic
308 correlations between pairs of subtypes involving the same disease were significantly less than 1
309 (FDR<0.1) for hypertension (CER vs. CVD: $\rho = 0.86 \pm 0.04$, $P=0.0004$; MGND vs. CVD: $\rho =$
310 0.74 ± 0.05 , $P=3 \times 10^{-8}$) and type 2 diabetes (CER vs. MGND: $\rho = 0.64 \pm 0.09$, $P=8 \times 10^{-5}$)
311 (Figure 6A; Supplementary Table 9). In addition, we observed significant excess genetic
312 correlations (FDR<0.1) for 8 pairs of disease subtypes involving different diseases (Figure 6A;
313 Supplementary Table 9). We verified that the excess genetic correlations could not be explained
314 by non-disease-specific differences in the underlying topics (which are weakly heritable;
315 Supplementary Table 3) by repeating the analysis using disease cases and controls with matched
316 topic weights (Methods, Supplementary Figure 20). We also estimated subtype-specific SNP-
317 heritability and identified some instances of differences between subtypes, albeit with limited
318 power (Supp Table 10).

319
320 Finally, we used the population genetic parameter $F_{ST}^{53,54}$ to quantify genome-wide differences
321 in allele frequency between two subtypes of the same disease; we used F_{ST} on control sets with
322 matched topic weights to assess statistical significance while accounting for non-disease-specific
323 differences in the underlying topics (excess F_{ST} ; Methods). We determined that 63 of 104 pairs
324 of disease subtypes involving the same disease (spanning 29 of 49 diseases, excluding 3 diseases
325 that did not have enough controls with matched topic weights) had significant excess F_{ST}
326 estimates (FDR < 0.1) (Supplementary Figure 21, Supplementary Table 11). For example, the
327 CVD, CER, and MGND subtypes of type 2 diabetes had significant excess F_{ST} estimates (F -
328 statistic=0.0003, $P=0.001$ based on 1,000 matched control sets). This provides further evidence
329 that disease subtypes as determined by comorbidity have different molecular and physiological
330 aetiologies.

331
332 We conclude that disease subtypes defined by distinct topics are genetically heterogeneous.

333
334 *Disease-associated SNPs have subtype-dependent effects*
335 We hypothesised that disease genes and pathways might differentially impact the disease
336 subtypes identified by ATM. We investigated the genetic heterogeneity between disease
337 subtypes at the level of individual disease-associated variants. We employed a statistical test that
338 tests for SNP x topic interaction effects on disease phenotype in the presence of separate SNP
339 and topic effects (Methods). We verified via simulations that this statistical test is well-calibrated
340 under a broad range of scenarios with no true interaction, including direct effect of topic on
341 disease, direct effect of disease on topic, pleiotropic SNP effects on disease and topic, and
342 nonlinear effects (Supplementary Figure 22). We also assessed the power to detect true
343 interactions (Supplementary Figure 23). To limit the number of hypotheses tested, we applied
344 this test to independent SNPs with genome-wide significant main effects on disease (Methods).
345 We thus performed 2,530 statistical tests spanning 888 disease-associated SNPs, 14 diseases, and
346 35 disease subtypes (Supplementary Table 12). We assessed statistical significance using global

347 FDR<0.1 across the 2,530 statistical tests. We also computed main SNP effects specific to each
348 quartile of topic weights across individuals, as an alternative way to represent SNP x topic
349 interactions.

350

351 We identified 43 SNP x topic interactions at FDR<0.1 (Figure 7, Supplementary Figure 24,
352 Supplementary Table 13 and Supplementary Table 14). Here, we highlight a series of examples.
353 First, the type 2 diabetes-associated SNP rs1063192 in the *CDKN2B* locus has a higher odds
354 ratio in the top quartile of CVD topic weight (1.19 ± 0.02) than in the bottom quartile (1.08 ± 0.02)
355 ($P = 4 \times 10^{-4}$ for difference). *CDKN2B* is associated with both coronary artery disease and type 2
356 diabetes⁵⁵⁻⁵⁹, suggesting that shared pathways underlie the observed SNP x topic interaction.
357 Second, the asthma-associated SNP rs1837253 in the *TSLP* locus has a higher odds ratio in the
358 top quartile of SRD topic weight (1.17 ± 0.02) than in the bottom quartile (1.05 ± 0.02)
359 ($P = 1 \times 10^{-4}$ for difference). *TSLP* plays an important role in promoting Th2 cellular responses
360 and is considered a potential therapeutic target, which is consistent with assignment of asthma
361 and atopic/contact dermatitis⁶⁰ to the SRD topic (Supplementary Table 4). Third, the
362 hypertension-associated SNP rs3735533 within the *HOTTIP* long non-coding RNA has a lower
363 odds ratio in the top quartile of CVD topic weight (1.07 ± 0.02) than in the bottom quartile
364 (1.13 ± 0.02). *HOTTIP* is associated with blood pressure^{27,61} and conotruncal heart
365 malformations⁶². Fourth, the hypothyroidism-associated SNP rs9404989 in the *HCG26* long non-
366 coding RNA has a higher odds ratio in the top quartile of FGND topic weight (1.90 ± 0.24) than in
367 the bottom quartile (1.19 ± 0.13) ($P = 3 \times 10^{-3}$ for difference). Hypothyroidism associations have
368 been reported in the HLA region²⁷, but not to our knowledge in relation to the *HCG26*. To verify
369 correct calibration, we performed control SNP x topic interaction tests using the same 888
370 disease-associated SNPs together with random topics that did not correspond to disease subtypes,
371 and confirmed that these control tests were well-calibrated (Supplementary Figure 24B).

372

373 We conclude that genetic heterogeneity between disease subtypes can be detected at the level of
374 individual disease-associated variants.

375

376

377 Discussion

378 We have introduced an age-dependent topic modelling (ATM) method to provide a low-rank
379 representation of longitudinal disease records, leveraging age-dependent comorbidity profiles to
380 identify and validate biological subtypes of disease. Our study builds on previous studies on
381 topic modelling^{37,38,40,63}, genetic subtype identification^{13–15}, and low-rank modelling of multiple
382 diseases to identify shared genetic components^{25–27}. We highlight three specific contributions of
383 our study. First, we incorporated age at diagnosis information into our low-rank representation,
384 complementing the use of age information in other contexts^{32,52,64}; we showed that age
385 information is highly informative for our inferred comorbidity profiles in both simulated and
386 empirical data, emphasising the importance of accounting for age in efforts to classify disease
387 diagnoses. Second, we identified 52 diseases with heterogeneous comorbidity profiles that we
388 used to define disease subtypes, many of which had not previously been identified
389 (Supplementary Table 15). Third, we used genetic data (including PRS, genetic correlation and
390 F_{ST} analyses) to validate these disease subtypes, confirming that the inferred subtypes reflect true
391 differences in disease aetiology.

392

393 We emphasise three downstream implications of our findings. First, it is of interest to perform
394 disease subtype-specific GWAS on the disease subtypes that we have identified here, analogous
395 to GWAS of previously identified disease subtypes^{13–15}. Second, our findings motivate efforts to
396 understand the functional biology underlying the disease subtypes that we identified; the recent
397 availability of functional data that is linked to EHR is likely to aid this endeavor^{29,65}. Third, it is
398 of interest to apply ATM to identify age-dependent comorbidity profiles and disease subtypes in
399 other EHR data sets^{30,31}; establishing representations of disease topics that are transferable and
400 robust across different healthcare systems and data sources represents a major future challenge.

401

402 Our findings reflect a growing understanding of the importance of context, such as age, sex,
403 socioeconomic status and previous medical history, in genetic risk^{52,66,67}. To maximise power
404 and ensure accurate calibration, context information needs to be integrated into clinical risk
405 prediction tools that combine genetic information (such as polygenic risk scores^{1,68}) and non-
406 genetic risk factors. Our work focuses on age, but motivates further investigation of other
407 contexts. We note that aspects of context are themselves influenced by genetic risk factors, hence
408 there is an open and important challenge in determining how best to combine medical history
409 and/or causal biomarker measurements with genetic risk to predict future events⁶⁹.

410

411 We note several limitations of our work. First, age at diagnosis information in EHR data may be
412 an imperfect proxy for true age at onset, particularly for less severe diseases that may be detected
413 as secondary diagnoses; although perfectly accurate age at onset information would be ideal, our
414 study shows that that imperfect age at diagnosis information is sufficient to draw meaningful
415 conclusions. Second, raw EHR data may be inaccurate and/or difficult to parse¹; again, although
416 perfectly accurate EHR data would be ideal, our study shows that imperfect EHR data is

417 sufficient to draw meaningful conclusions. Third, our ATM approach incurs substantial
418 computational cost (Supplementary Table 16); however, analyses of biobank-scale data sets are
419 computationally tractable, with our main analysis requiring only 4.7 hours of running time.
420 Finally, we have applied ATM to a UK population of predominantly European ancestry; it is of
421 interest to apply ATM to diverse populations^{30,31}. Despite these limitations, ATM is a powerful
422 approach for identifying age-dependent comorbidity profiles and disease subtypes.
423

424 Acknowledgements

425
426 This research has been conducted using the UK Biobank Resource; application number 12788.
427

428 Funded by Wellcome (BST00080- H503.01 to XJ, 100956/Z/13/Z to GM, [https://](https://wellcome.org)
429 wellcome.org); the Li Ka Shing Foundation (to GM, <https://www.lksf.org>); NIH grants R01
430 HG006399, R01 MH101244, and R37 MH107649 (to ALP); The Alan Turing Institute
431 (<https://www.turing.ac.uk>), Health Data Research UK (<https://www.hdruk.ac.uk>), the Medical
432 Research Council UK (<https://mrc.ukri.org>), the Engineering and Physical Sciences Research
433 Council (EPSRC <https://epsrc.ukri.org>) through the Bayes4Health programme Grant
434 EP/R018561/1, and AI for Science and Government UK Research and Innovation (UKRI,
435 <https://www.turing.ac.uk/research/asg>) (to CH); BHF Chair award CH/12/2/29428 (to XJ). This
436 work was supported by core funding from the: British Heart Foundation (RG/13/13/30194;
437 RG/18/13/33946), BHF Cambridge Centre of Research Excellence (RE/13/6/30180) and NIHR
438 Cambridge Biomedical Research Centre (BRC-1215-20014). The funders had no role in study
439 design, data collection and analysis, decision to publish, or preparation of the manuscript.
440

441 This work uses data provided by patients and collected by the NHS as part of their care and
442 Support. Computation used the Oxford Biomedical Research Computing (BMRC) facility, a
443 joint development between the Wellcome Centre for Human Genetics and the Big Data Institute
444 supported by Health Data Research UK and the NIHR Oxford Biomedical Research Centre. The
445 views expressed are those of the authors and not necessarily those of the NHS, the NIHR or the
446 Department of Health. We thank Kushal Dey, Luke Kelly and Yunlong Jiao for the discussion.
447

448 Data availability

449 UK Biobank data is publicly available at <https://www.ukbiobank.ac.uk/>.

451 Code availability

452 Open-source software implementing the ATM method is available at [https://github.com/Xilin-](https://github.com/Xilin-Jiang/ATM)
453 [Jiang/ATM](https://github.com/Xilin-Jiang/ATM). BOLT-LMM is available at <https://alkesgroup.broadinstitute.org/BOLT-LMM/>.
454 Heritability and genetic correlation analysis were performed using LDSC, which is available at

455 <https://github.com/bulik/ldsc>. PLINK v1.9, which was used for F_{ST} and association tests, is
456 available at <https://www.cog-genomics.org/plink/>.

457

458

459 Methods

460 Age-dependent topic model (ATM)

461 Our Age-dependent topic model (ATM) is a Bayesian hierarchical model to infer latent risk
462 profiles for common diseases. The model assumes that each individual possesses several age-
463 evolving disease profiles (topic loadings), which summarise the risk over age for multiple
464 diseases that tend to co-occur within an individual's lifetime, namely the age specific multi-
465 morbidity profiles. At each disease diagnosis, one of the disease profiles is first chosen based on
466 individual weights of profile composition (topic weights), the disease is then sampled from this
467 profile conditional on the age of the incidence.

468

469 We constructed a Bayesian hierarchical model to infer K latent risk profiles for D distinct
470 common diseases. Each latent risk profile (comorbidity topics) is age-evolving and contains risk
471 trajectories for all D diseases considered. Each individual might have a different number of
472 diseases, while the disease risk is determined by the weighted combination of latent risk topics.

473 The indices in this note are as follows:

- 474 • $s = 1, \dots, M$;
- 475 • $n = 1, \dots, N_s$;
- 476 • $i = 1, \dots, K$;
- 477 • $j = 1, \dots, D$;

478 where M is the number of subjects, N_s is the number of records within s^{th} subject, K is the
479 number of topics, and D is the total number of diseases we are interested in. The plate notation of
480 the generative model is summarised in Supplementary Figure 1:

- 481 • $\theta \in R^{M \times K}$ is the topic weight for all individuals (referred to as patient topic weights),
482 each row of which ($\in R^K$) is assumed to be sampled from a Dirichlet distribution with
483 parameter α . α is set as a hyper parameter: $\theta_s \sim Dir(\alpha)$.
- 484 • $z \in \{1, 2, \dots, K\}^{\sum_s N_s}$ (referred to as diagnosis-specific topic probability) is the topic
485 assignment for each diagnosis $w \in \{1, 2, \dots, D\}^{\sum_s N_s}$. Note the total number of
486 diagnoses across all patients are $\sum_s N_s$. The topic assignment for each diagnosis is
487 generated from a categorical distribution with parameters equal to s^{th} individual topic
488 weight: $z_{sn} \sim Multi(\theta_s)$.
- 489 • $\beta(t) \in F(t)^{K \times D}$ is the topic loading which is $K \times D$ functions of age t . $F(t)$ is the
490 class of functions of t . At each plausible t , the following is satisfied: $\sum_j \beta_{ij}(t) = 1$. In
491 practice we use softmax function to ensure above is true and add smoothness by constrain
492 $F(t)$ to be spline or polynomial functions: $\beta_{ij}(t) = \frac{\exp(p_{ij}^T \phi(t))}{\{\sum_{j=1}^D \exp(p_{ij}^T \phi(t))\}}$, where
493 $p_{ij} = \{p_{ijd}\}$; $d = 1, 2, \dots, P$; P is the degree of freedom that controls the smoothness;
494 $\phi(t)$ is polynomial and spline basis for age t .

- 495 • $w \in \{1, 2, \dots, D\}^{\sum_s N_s}$ are observed diagnoses. The n^{th} diagnosis of s^{th} individual w_{sn}
496 is sampled from the topic $\beta_{z_{sn}}(t)$ chosen by z_{sn} : $w_{sn} \sim Multi(\beta_{z_{sn}}(t_{sn}))$, here t_{sn} is the
497 age of the observed age at diagnosis of the observed diagnosis w_{sn} .

498

499 The values of interest in this model are global topic parameter β , individual (patient) level topic
500 weight θ , and diagnosis-specific topic probability z . Based on the generative process above, we
501 notice that each patient is independent conditional on α . Therefore, the inference of θ and
502 z (discussed below) could be performed by looping each individual in turn.

503

504 The key element in our model is age-evolving risk profiles, which is achieved by model the
505 comorbidity trajectories $\beta(t) \in F(t)^{K \times D}$ as functions of age. The functionals $F(t)$ considered
506 are linear, quadratic, cubic polynomials, and cubic splines with one, two and three knots.

507

508 Inference of ATM

509 The variables of interest are global topic parameter $\beta(t)$, individual (patient) level topic weight
510 θ , and diagnosis-specific topic probability z of each diagnosis. We could adopt an EM strategy,
511 where in the E-step we first estimate posterior distribution of θ and z , then in the M-step we
512 estimate β which maximises the evidence lower bound (ELBO).

513

514 The details of the inference is explained in Supplementary Note. In summary, in a Bayesian
515 setting, the model could be evaluated by the evidence function $p(w|\alpha, \beta)$. The best $\beta(t)$ is found
516 by maximise the evidence function, while for θ and z we aim to find or approximate their
517 posterior distribution $p(z, \theta | w, \alpha, \beta)$. Given that the posterior distribution is intractable, we use
518 variational distribution $q(z, \theta)$ to approximate them. Now we could write the evidence function
519 as:

$$520 \quad p(w | \alpha, \beta) = L(z, \theta, \beta, \alpha) + KL(q||p),$$

521 here $KL(q||p) = - \int_{z, \theta} q(z, \theta) \ln \frac{p(z, \theta | w, \alpha, \beta)}{q(z, \theta)}$ is the KL divergence. Since KL divergence is

522 always positive, $L(z, \theta, \beta, \alpha)$ is a lower bound of the evidence function:

$$523 \quad L(z, \theta, \beta, \alpha) = E_q \{ \ln p(w, z, \theta | \alpha, \beta) - \ln q(z, \theta) \}.$$

524

525 When finding the posterior of θ and z , we want $\ln q(z, \theta)$ to be as close to the posterior
526 $p(z, \theta | w, \alpha, \beta)$ as possible. Since $KL(q||p) = 0$ when $q(z, \theta) = p(z, \theta | w, \alpha, \beta)$, this could
527 be achieved by minimising $KL(q||p)$ or maximise $L(z, \theta, \beta, \alpha)$. The most commonly used form
528 of $q(z, \theta)$ assume the distribution is factorised, which might cause instability when signal-to-
529 noise ratio is low⁷⁰. Therefore, more accurate inference methods such as collapsed variational
530 inference is considered⁴⁰. Comparison of the evidence lower bound $L(z, \theta, \beta, \alpha)$ shows collapsed
531 variational inference is consistently more accurate than LDA (Supplementary Figure 8).

532 Therefore we choose the collapsed variational inference⁴⁰. The collapsed variational inference is
533 achieved by integrate out θ from the likelihood function $p(w, z, \theta | \alpha, \beta)$ and find the

534 approximated posterior distribution $q(z)$. For detailed derivation, the comparison between
535 collapsed variational inference and mean-field variational inference, and update algorithms, see
536 Supplementary Note.

537
538 When finding the $\beta(t)$ that maximises the evidence function, we again maximise $L(z, \theta, \beta, \alpha)$.
539 Maximising $L(z, \theta, \beta, \alpha)$ with respect to $\beta(t)$ does not have an analytical solution due to its
540 softmax structure. We use local variational methods and numeric optimisation to find the
541 distribution of $\beta(t)$. Details are provided in Supplementary Note.

542
543 We extract topic weights at patient-level and diagnosis-level from the posterior distribution
544 inferred from the data. Our model has the desired property that each patient and patient-diagnosis
545 are assigned to comorbidity topics. The model estimates the posterior distribution $q(z)$, which is
546 a categorical distribution (Supplementary Note.) Several metrics related to topic assignments
547 could be derived from the $q(z)$:

- 548 • Each patient-diagnosis (incident disease) has a diagnosis-specific topic probability, which
549 is computed as $E_q\{z_n\}$.
- 550 • Each patient has a posterior topic assignment θ_s , which is a dirichlet distribution $\theta_s \sim$
551 $Dir(\alpha + \sum_{n=1}^{N_S} E_q\{z_n\})$. The topic weights of each patient is the mode of this
552 Dirichlet distribution $\frac{\sum_{n=1}^{N_S} E_q\{z_n\}}{\sum_{i=1}^K \sum_{n=1}^{N_S} E_q\{z_{ni}\}}$ (we used $\alpha = 1$). The value is used as the patient
553 low-rank representation of disease history, for analysis including PRS association with
554 comorbidity within cases and G x Topic interaction analysis.
- 555 • The average topic assignments of disease j is the mean over all incidences
556 $\overline{E_q\{z_{sn} \in \{w_{sn}=j\}\}}$. This metric is used to measure which comorbidity topic a disease is
557 associated with (Figure 4B), and it is equivalent to a weighted average of topic loadings
558 (for the specific weighted average expression, see equation 5 of Supplementary Note). A
559 disease assigned to multiple topics is considered to have comorbidity subtypes.
- 560 • A hard assignment of a patient-diagnosis to a subtype is based on the max value of the
561 vector $E_q\{z_n\}$. The incident disease is assigned to topic $argmax_i (E_q\{z_{ni}\})$.

562 563 **Metrics for evaluating ATM**

564 ATM is evaluated for different purposes, which requires different metrics (Supplementary Table
565 1). Here we list the details of the four metrics considered: *Prediction odds ratio*, *Evidence Lower*
566 *Bound (ELBO)*, *AURPC*, and *Co-occurrence odds ratio*.

567
568 *Prediction odds ratio*: To compare models of different topic numbers and configuration of age
569 profiles, we compare the prediction odds ratio of each model. Briefly, prediction odds ratio is
570 defined on 20% held-out test data as the odds that the true diseases are within the top 1%
571 diseases predicted by ATM (trained on 80% of the training set and uses earlier diagnoses as

572 input), divided by the odds that the true diseases are within the top 1% of diseases ranked by
573 prevalence.

574

575 Specifically, we separate UK Biobank patients into a training set (80%) and a testing set (20%).
576 On the training set, we estimate the comorbidity topic loadings. On the testing set, we fix the
577 topic loadings and infer the patient topic weights to predict the next disease in chronological
578 order. The topic loadings are estimated using the n diseases and compute the risk rank of
579 diseases at the age of the $n+1$ disease. The odds ratio is computed by the odds of the $n+1$ disease
580 being in the top 1% of diseases versus being in the top 1% most prevalent diseases. We use the
581 top 1% most prevalent diseases instead of randomly chosen diseases as it represents a naive
582 prediction model that predicts disease based on prevalence. The patient topic weights
583 computation is in section Inference of ATM and the risk is computed as the linear combination
584 of topics using topic weights as coefficients. We also compute the prediction odds ratio using the
585 LDA model. We repeat the procedure for 10 times for each model configuration.

586

587 We compared the prediction odds ratio for topic number between 5 to 20, with linear, quadratic
588 polynomial, cubic polynomial, and splines with one, two and three knots. We also compare the
589 ATM model with the LDA model of topic number between 5 to 20.

590

591 *Evidence Lower Bound (ELBO)*: ELBO evaluated the accuracy of the variational inference
592 method on a specific data set. Mathematical expression of ELBO for ATM is presented in
593 equation 9 in Supplementary Note. To find the best model that fit to the entire dataset, we
594 evaluate the ELBO for models with topic numbers between 5 to 20, 25, 30, and 50 topics and age
595 profiles configured by linear, quadratic polynomial, cubic polynomial, and splines with one, two
596 and three knots. Each model is run for 10 times with random initialisations. We choose the
597 model that has the highest ELBO after converging.

598

599 *AURPC*: To evaluate whether a model could capture the comorbidity subtypes in simulation
600 analysis, we compute the precision, recall, and area under precision-recall curve (AUPRC) to
601 correctly classify disease diagnosis to be from the topic that it is generated from. The topic of
602 each diagnosis is determined by diagnosis-specific topic probability. Note we could only
603 evaluate AUPRC in simulations where the truth is known.

604

605 *Co-occurrence odds ratio*: To verify that the comorbidity profiles that the model captured are
606 capturing diseases that are more likely to present within the same individual, we estimate the
607 odds ratio of the disease duo, trio, quartet, and quintet that are captured by the topic versus that
608 of random combinations. We divide the population into an 80% training set and a 20% testing
609 set. We trained the ATM model with five random initialisations and kept the inference with the
610 highest ELBO. Each disease is assigned to a topic by the highest average topic assignments.
611 (section Inference of ATM) We focus on the top 100 diseases ranked by prevalence to avoid the

612 combination being too rare to appear in the population. In the testing set, we computed the odds
613 of individuals who have all diseases in the comorbidities versus the odds implied if all diseases
614 are independent (computed as the product of disease prevalence). The odds ratio is computed for
615 all combinations of duo, trio, quartet, and quintet that are assigned to the same topics. We
616 perform the same analysis using PCA for comparison.

617
618

619 **Simulations of ATM method.**

620 To test whether the algorithm could assign disease diagnosis to correct comorbidity profiles, we
621 simulated disease from two disease topics within a population of 10,000, using following
622 parameters:

- 623 • $M = 10,000$;
- 624 • $\overline{N_S} = 6.1$;
- 625 • $N_S \sim \text{exp}\{\overline{N_S}\}$;
- 626 • $D = 20$;
- 627 • $K = 2$;

628 Here M is the number of individuals in the population, $\overline{N_S}$ is the average number of diseases for
629 each individual, D is the total number of diseases, K is the number of comorbidity topics. The
630 distribution of disease number per-individual N_S is sampled from an exponential distribution,
631 which matches those from UK Biobank data (Supplementary Figure 26). According to equation
632 3.1 in Ghorbani et al.⁷⁰, whether the topic model could capture the true latent structure is
633 determined by the information signal-to-noise ratio and could be evaluated with limits $M \rightarrow$
634 ∞ ; $D \rightarrow \infty$; $\frac{D}{M} \rightarrow \delta$, where δ is a constant. Therefore we choose D and M at scales that make $\frac{D}{M}$
635 approximately similar to those of the UK Biobank dataset (Samples size = 282,957; distinct
636 disease number = 349).

637

638 The simulated topics loadings are constructed as follows:

- 639 • All but K diseases are simulated to be associated with comorbidity profiles. Each of them
640 has a risk period of 30 years and overlaps for 10 years with the next disease. For
641 example, if disease 1 has a risk period from 30 to 59 years of age, disease 2 will have a
642 risk period between 50 to 79 years of age. When the risk period reaches the maximal age,
643 the truncated part will be carried to the next disease to create diseases with shorter risk
644 period. All risk periods are assigned a value 1.
- 645 • K diseases that are not associated with comorbidity are simulated to span all topics. The
646 values of these diseases are sampled from $Unif(0, \frac{0.1}{K})$ for each topic. Here K is the
647 number of topics.
- 648 • The age profiles are then normalised at each age point to ensure $\sum_{j=1}^D \beta_j(t) = 1$ for all
649 t . With this constraint we could sample a disease at each age t using a multinomial
650 probability with the topic loading as the parameter. The age range of the simulated topics

651 is 30 to 81 years of age, which is the minimal and maximal age at diagnosis of incident
652 disease in the UK Biobank population. An example of a simulated topic is shown in
653 Supplementary Figure 27.

654
655 For each individual, we sampled the Dirichlet parameter α from a gamma distribution (shape =
656 50, rate = 50). Topic loadings are sampled from the Dirichlet distribution for each patient as the
657 generative process. For each patient, we first sample the number of diseases N_S . For each
658 incident disease, we sample the disease age from uniform distribution between age 30 to 81 and
659 a topic from the topic loading. We then choose the incident disease based on the age at diagnosis
660 from the chosen topic. The procedure follows the generative process described in Supplementary
661 Note.

662
663 Since in real data we only use the first age at diagnosis for diseases that are recorded repeatedly
664 within the same patient, we filter the simulated diseases accordingly. The filtered data are fed
665 into the inference functions to infer the latent topics and disease assignments. The inferred topics
666 resemble the true topics used to simulate diseases as shown in Supplementary Figure 27. For the
667 initialisation of each inference, we first sample β and θ from the Dirichlet distribution of non-
668 informative hyperparameters, then initialise other variables parameters following the generative
669 process. The variational inference converged where the relative increase of ELBO is below 10^{-6} .

670
671 To simulate disease having distinct comorbidity subtypes, we first simulate diseases using the
672 procedure described above. We consider two scenarios: (1) the subtype of diseases have the
673 same age at diagnosis distribution. (2) the subtypes of disease have distinct age at diagnosis
674 distribution.

675
676 We create diseases with distinct comorbidity profiles by combining diseases that are sampled
677 from distinct topics and labelling them as a single disease. We first chose one disease (**disease**
678 **A**) then sampled a proportion of a second disease (**disease B**) to label as **disease A**. The
679 proportion is varied to create a different sample size ratio of the two subtypes. In scenario one,
680 **disease B** is a disease that has the exact same age age distribution as **disease A** but from the other
681 topic. In scenario two, **disease B** is from the other topic and has a different age distribution (age
682 at diagnosis moves up for 20 years, 10 years, or 5 years, respectively) than **disease A**. After
683 changing the labels of **disease B** to be the same as **disease A**, we used the inference procedure
684 described as above to get the posterior distribution.

685
686 To evaluate whether a model could capture the comorbidity subtypes, we compute the precision,
687 recall, and area under precision-recall curve (AUPRC) to correctly classify incident **disease B** to
688 be from the topic that it is generated from. The topic of each diagnosis is determined by
689 diagnosis-specific topic probability. We use other diseases from the topic of **disease B** to
690 benchmark the topic label. Topic modelling on the simulated data is performed with both ATM

691 and LDA (both implemented using collapsed variational inference for fair comparison) to
692 compare the performances.

693

694 We evaluate the subtype classification with varying values for three simulation parameters:

- 695 ● ratio of sample sizes between the two subtypes. We change the ratio of the two subtypes
696 by a grid between 0 to 0.9 with a step size 0.1. The default value of sample size ratio is
697 set as 0.1 in other simulations to test for other parameters that have impacts on the
698 precision and recall.
- 699 ● Simulated population size. We simulated population sizes equal to 200, 500, 1000, 2000,
700 5000, and 10,000. The default population size is 10,000 in other simulations.
- 701 ● Number of distinct diseases. We simulated datasets with 20, 30, 40, and 50 distinct
702 diseases, with 2, 3, 4 and 5 underlying disease topics respectively. The default number of
703 distinct diseases is 20 in other simulations.
- 704 ● Difference of age distribution. We considered three scenarios of subtype age distribution,
705 with 0, 10, and 20 years of difference in the average age at diagnosis.

706

707

708 **UK Biobank comorbidity data.**

709 We analysed comorbidity data from 282,957 UK Biobank samples with diagnoses for at least
710 two of the 348 focal diseases that we studied (see below). We use the hospital episode statistics
711 (HES) data within the UK Biobank dataset, which records diseases using the ICD-10/ICD-10CM
712 coding system. Codes started with letters from A to N are kept as they correspond to disease
713 code (opposed to procedure codes). The disease records were mapped from ICD-10/ICD-10CM
714 codes to PheCodes using a three-step procedure: Firstly, we map the first four letters of each
715 ICD-10 records to the phecodes, using the map file downloaded from phewascatalog.org;
716 Secondly, we map the remaining records using ICD-10CM map file downloaded from
717 phewascatalog.org; Lastly, we map remaining records to a collapsed ICD-10CM mapping
718 system which only use the first four character of ICD-10CM codes. We also noticed an ICD-
719 10/ICD-10CM code could map to multiple PheCodes. When a single ICD-10/ICD-10CM code s
720 mapped to more than one PheCodes, we only kept the Phecode that are mapped to the most ICD-
721 10 codes (i.e. PheCode is constructed by combining ICD-10 that represent similar diseases. The
722 Phecode that represent a larger number of ICD-10 codes are more likely to be a well defined
723 disease, which we chose to keep.), which ensure that one ICD-10(CM) code only maps to one
724 PheCode. Using the procedure above, we mapped 99.7% ICD-10/ICD-10CM code to PheCodes,
725 with 4,637,127 records in total.

726

727 The mapped Phecodes are filtered to keep only the first age at diagnosis for the same diseases
728 within a patient. The age at diagnosis for each record is computed as the difference between
729 month of birth to the episode starting date. We then computed the occurrence of each disease in
730 the UK Biobank and kept 348 that have more than 1,000 occurrences (Supplementary Table 4).

731 Starting with all 488,377 UK Biobank patients (including both European and non-European
732 ancestries), we filtered the patients to keep only those who have at least two distinct diseases
733 from the 348 focal diseases, as we are most interested in the comorbidity information. We treated
734 the death as an additional disease (8,666 records) to evaluate if certain comorbidities are more
735 likely to lead to fatal events. After these procedures, there are in total 1,726,144 distinct records
736 across 282,957 patients.

737
738 To name the topics inferred from the UK Biobank, we take the sum of average topic assignments
739 (section Inference of ATM) over diseases that are within each phecode system and extract the
740 top 3 systems. Most of comorbidity topics are named using the first three topics (e.g. CER:
741 cardiovascular, endocrine/metabolic, respiratory), except for topics that are predominantly
742 associated with one system (LGI: lower gastrointestinal; UGI: upper gastrointestinal; CVD:
743 cardiovascular).

744
745 We present focal diseases for each topic in two ways. Firstly, we filter each topic using the
746 profile mean value between age 30 to 81 to keep the top seven diseases. We chose seven for
747 visualisation, as we found more diseases would be harder to read on a plot. Secondly, we also
748 show seven diseases that have the highest average assignment to each topic. This will give a
749 picture of diseases that are not the most prevalent in the population but are predominantly
750 associated with the target topic.

751
752 To compare the comorbidity heterogeneity between age groups, we group the incidences for each
753 disease to two age groups: young group (<60 years of age) and old group (≥ 60 years of age). We
754 compute the average topic assignment of each group as described in section Inference of ATM.
755 Additionally, we inferred topics for male (984,554 records in 156,366 individuals) and female
756 (741,590 records in 126,591 individuals) populations respectively using a model with 10 topics
757 and spline function with one knot. We extract the average topic assignment for each disease, and
758 use Pearson's correlation to match the topics for both sexes to the topics inferred on the entire
759 population.

760
761 Each diagnosis could be assigned to a specific topic using max diagnosis-specific topic
762 probability. We focus our disease heterogeneity analysis on 52 diseases that have at least 500
763 incidences assigned to a secondary topic.

764
765 **UK Biobank genotype data.**

766 For all analyses except BOLT-LMM we use 488,377 UK Biobank participants. For BOLT-LMM
767 analyses, we constrain our analysis to 409,694 British Isle ancestry individuals to remove the
768 possibility that topics are capturing population structure. For F_{ST} analysis with PLINK we used
769 805,426 genotyped SNPs; for BOLT-LMM PRS analysis we used 727,882 genotyped SNP with
770 $MAF > 0.1\%$; for genetic correlation analysis using LDSC, we used 157,756 Genotyped SNPs

771 mapped to HapMap3 SNPs; for BOLT-LMM and subsequent LDSC analysis that use imputed
772 SNPs, we used 1,201,838 imputed SNPs mapped to HapMap3 SNPs.

773

774 **Polygenic risk scores (PRS) analysis.**

775 To exclude the possibility of population stratification, we compute PRS using mixed-effect
776 association on the British Isle ancestry group ($N = 409,694$) for the 10 heritable diseases that
777 have the highest heritability z-scores. We used a mixed model to estimate effect size
778 implemented by BOLT-LMM and constructed genome-wide PRS⁵⁰. For the computation of
779 PRS, we randomly sampled half of the British isle ancestry population ($N = 204,847$) for
780 computation efficiency (essential hypertension, arthropathy, asthma, and hypercholesterolemia)
781 or sampled 9 controls for each case to ensure case proportion at or above 10% as recommended
782 by BOLT-LMM (type 2 diabetes, varicose veins of lower extremity, hypothyroidism, other
783 peripheral nerve disorders, major depressive disorder, and GRED). We used PLINK to select
784 genotyped SNPs with $MAF > 0.1\%$ as recommended in BOLT-LMM. For each disease, we used
785 5-fold cross validation to estimate effect sizes using BOLT-LMM and computed the PRS on the
786 held-out testing set. The predictive PRS are then used to compute the excess PRS over different
787 topic loadings, by a linear regression where PRS is the response variable and topic weights is the
788 predictor.

789

790 We compute the relative risk for each percentile of PRS using the following formula:

$$791 \quad RR_{pt,s} = \frac{n_{pt,s} \times 100}{n_s},$$

792 where $RR_{pt,s}$ is the relative risk of s subtype for the pt^{th} PRS percentile (computed for the entire
793 population); $n_{pt,s}$ is the number of cases in s subtype that has PRS within the pt^{th} percentile; n_s
794 is the number of cases in the s subtype.

795

796 **Genetic correlation analysis.**

797 For each disease and disease subtype, we use a case-control matching strategy to construct data
798 to estimate coefficients for genetic correlation analysis. For each case in the disease group, we
799 pick four nearest neighbors (without replacement) from the control group, matching sex, BMI,
800 year of birth and 40 genetic principal components. The covariates are available within the UK
801 Biobank data set, over which we computed the principal components. We then compute the
802 Euclidean distance of the principal components to find the nearest neighbours in the population.
803 All cases are matched with four controls except for 401.1 essential hypertension which has a
804 sample size larger than 20% of the population. We match only one control for each hypertension
805 case.

806

807 We perform logistic regression with sex and top 10 principal components as covariates to
808 estimate the main variant effect of the 805,426 variants that are genotyped. We used PLINK 1.9
809 for association analysis⁷¹. With the summary statistics from the association analysis, we use

810 LDSC to map the summary statistics to HapMap3 SNPs and match the effect and non-effect
811 alleles^{2,72}. Since UK Biobank is mostly of British Isle ancestry, we use the pre-computed LD
812 score from the LDSC website. We estimated the heritability for each disease or disease subtype
813 which has more than 1000 incidences (378 diseases subtypes and diseases). We use 1000
814 incidence threshold as LDSC are more accurate with larger sample size. We focus on 71 disease
815 and 18 disease subtypes that have heritability z-score above 4 for genetic correlation analysis.

816

817 The genetic correlation is computed for each pair of disease/subtypes using the same summary
818 statistics and LD score regression. We report the estimate of genetic correlation and z-scores.
819 Additionally, for pairs that involve subtypes (disease-subtype or subtype-subtype), we compute
820 the excess genetic correlation, defined as the difference between the genetic correlation
821 involving subtypes and the genetic correlation involving all disease diagnoses. For example, the
822 genetic correlation between T2D-CER and hypertension-CVD is compared to the genetic
823 correlation between all T2D and all hypertension. The z-score and p-value of the genetic
824 correlation differences are reported. We note that genetic correlations between subtypes of the
825 same disease are compared to 1. We only reported p-values of excess genetic correlation when
826 both genetic correlation estimation has standard error <0.1 and at least one of the genetic
827 correlation has $|z\text{-score}|>4$.

828

829 To avoid potential collider effects where subtypes are defined by topic components that are
830 independent of the diseases, we further match cases in each subtype with controls that match the
831 topic loadings. We computed PCs from 23 variables (10 topic loadings, 10 PCs, year of birth,
832 sex, and BMI) and use the nearest neighbour procedure (by Euclidean Distance) to find controls
833 for each case. Here controls are chosen from individuals without the targeting disease, i.e. an
834 individual with one subtype of the target disease could not be a control for the other subtypes.
835 We performed the same analysis using this case-control matching procedure and compared the
836 genetic correlation with the case-control procedure described above. We perform the analysis for
837 four diseases that have evidence for genetic subtypes: asthma, type 2 diabetes,
838 hypercholesterolemia, and hypertension. For one subtype (hypertension-CVD), the heritability
839 (0.0313, s.e. = 0.0289) is below threshold after matching the topic, which was excluded in
840 genetic correlation analysis.

841

842 **F_{ST} analysis.**

843 To evaluate the genetic heterogeneity between disease subtypes, we estimated the F_{ST} for 52
844 diseases that have at least 500 incidences assigned to a secondary topic. To test the statistical
845 significance of F_{ST} , we adopted a permutation strategy and sampled the same number of controls
846 of similar topic weights distribution for each subtype. The topic weights are matched by
847 sampling (without replacement) the same number of controls for each dominant topic weight
848 quartile of the cases (i.e. matching the topic that defines the subtype), which ensures the controls
849 have the same topic weight stratification as the disease subtypes. We then compute the F_{ST} across

850 the control groups matched for subtypes. We excluded three diseases, “hypertension”,
851 “hypercholesterolemia”, and “arthropathy”, from F_{ST} analysis as we do not have enough controls
852 that match topic weight distribution. The F_{ST} s are computed using PLINK 1.9's weighted mean
853 across all genotyped SNPs, which report F statistics across all subtypes.

854

855 We obtained 1,000 permutation samples and reported the permutation p-value. Under the
856 assumption that causal and non-causal variants have similar allele frequency differences across
857 the subtypes, F_{ST} could be a measure of causal genetic effect heterogeneity across subtypes.

858

859 **SNP x topic interaction test.**

860 For the diseases that have heritability z-score above 4 in the UK Biobank, we further investigated
861 whether there are interactions between genetic risk factors with the topic loadings. We used a fit
862 a logistic regression model using following model:

$$863 \quad \text{logit}(p) = \beta_0 + \beta_1 * T + \beta_2 * T^2 + \beta_3 * G + \beta_4 * G * T,$$

864 where T is individual topic weights for a specified topic, G is the genotype, and p is the
865 probability of getting the disease. We computed the test statistics under the null that $\beta_4 = 0$. We
866 used QQ plots to check that the test statistics are well calibrated for each disease-topic pair.

867

868 Since the simulation shows the interaction test is underpowered when the variant effects are
869 small, we focus on the set of SNP that reaches genome-wide significance level to increase power
870 to detect interaction effects. We performed LD-clumping using $r^2 > 0.6$ to remove variants that
871 are in strong LD with the lead variants. We computed the test statistics using the model above
872 (for testing $\beta_4 = 0$) and computed study-wise FDR across disease-topic pairs.

873

874 To verify the significant interactions, we divided cases into quartiles based on topic loading for
875 each disease-topic pair, and randomly sampled two controls that match the topic loading for each
876 case. We estimated the main effect sizes for all GWAS-SNP within each quartile of topic
877 loadings to capture effects that are modulated by topic weights. We focus on the SNPs that have
878 significant interaction test statistics computed in the previous step and compare it with
879 background SNPs that have genome-wide significant main effects but no interaction effect
880 ($P > 0.05$).

881

882 **Simulations of SNP x topic interaction**

883 We simulate comorbidity with genetics to test interaction between genetic and comorbidity
884 topics. We simulated 100 independent variants with MAF randomly sampled from $Unif(0, 0.5)$.
885 We assumed an additive model and simulated genotypes for the population using Hardy-
886 Weinberg equilibrium. We simulated three types of genetic effects on topic and diseases on topic
887 of the simulation framework described in Simulations of ATM method section:

- 888 • Genetics-topic effect: each variant is simulated to have an linear effect of 0.04 on the
889 topic loading. We choose this value as after normalising the topic, a regression of causal

- 890 variant to topic would have an effect size approximately 0.01 which is similar to our
891 observation in the UK Biobank. The number of variants that are causal to the topic varies
892 between 2 to 20. We simulated the effect on one topic by adding additive SNP effects and
893 normalise the topic loadings of each patient. The topic-disease causality is a natural
894 consequence following the generative process of sampling data.
- 895 ● Genetic-disease-topic effect: we simulated a heritable disease that is causal to the topic.
896 The disease is simulated with 20 causal variants each of effect size 0.15. We vary the
897 disease-to-topic causal effect from 0.05 to 0.5, with a default value of 0.1 in other
898 analyses (similar to the correlation we found in UK Biobank analysis). We simulated the
899 effect on one topic by adding additive causal disease effects and normalise the topic
900 loadings of each patient.
 - 901 ● The genetic effect could interact with the topic when contributing to disease risk. We
902 simulated four additional diseases to represent different structures (Supplementary Figure
903 22).
 - 904 ○ Genetic effects interact with topic loading on altering disease risk. The interaction
905 term is added to the mean of disease liability, which is sampled from a Gaussian
906 distribution. The disease is then sampled by a threshold on the liability, where the
907 incidence rate is by default 0.5. The interaction effect is varied from 0.4 to 4, with
908 default value equal to 2.
 - 909 ○ Pleiotropy effects are simulated with a variant that have both genetic-disease and
910 genetic-topic-disease effects. Both genetic and topic effects are added to the mean
911 of disease liability. A disease is sampled by a threshold with default incidence rate
912 equal to 0.5. The topic-disease effect is varied from 0.4 to 4, with default value
913 equal to 2.
 - 914 ○ Pleiotropy effect with nonlinear topic-disease effect. A quadratic term of topic-
915 disease effect added to the second model.
 - 916 ○ Pleiotropy effect with nonlinear genetic-disease effect. A quadratic term of
917 genetic-disease effect added to the second model.

918
919 For disease-topic or topic-disease causal effects, we simulated 50 repetition at each causal effect
920 size. For interaction analysis, we repeated 10 times at each parameter value, as there are more
921 SNPs for uncertainty estimation. The simulated disease sets are fed into the inference procedure
922 to infer the patient topic weights.

923

924

925

926

927 **Tables**

928

Acronym	Disease systems	Representative diseases	Number of associated diseases
NRI	Neoplasms, respiratory, infectious diseases	Secondary malignancy of lymph nodes; Pneumococcal pneumonia; Bacterial infection NOS	53
CER	Circulatory system, endocrine/metabolic, respiratory	Type 2 diabetes; Obesity; Chronic airway obstruction	41
SRD	Sense organs, respiratory, dermatologic	Cataract; Septal Deviations/Turbinates Hypertrophy; Benign neoplasm of skin	38
CVD	Cardiovascular disease	Hypercholesterolemia; Coronary atherosclerosis; Myocardial infarction	27
UGI	Upper gastrointestinal disease	Diaphragmatic hernia; Benign neoplasm of other parts of digestive system; Gastritis and duodenitis;	22
LGI	Lower gastrointestinal disease	Irritable Bowel Syndrome; Benign neoplasm of colon; Anal and rectal polyp;	13
FGND	Female genitourinary, neoplasms, digestive	Uterine leiomyoma; Malignant neoplasm of female breast; Hypothyroidism NOS	34
MGND	Male genitourinary, neoplasms, digestive	Urinary tract infection; Cancer of prostate; Other disorders of bladder	33
MDS	Musculoskeletal, digestive, symptoms	Back pain; Cholelithiasis; Other disorders of soft tissues	29
ARP	Arthropathy-related disease	Arthropathy NOS; Rheumatoid arthritis; Enthesopathy	26

929

930

931

932 **Table 1. Summary of 10 inferred disease topics in the UK Biobank.** For each topic, we list its

933 3-letter acronym, disease systems, representative diseases, and number of associated diseases

934 (defined as diseases with average diagnosis-specific topic probability >50% for that topic).

935 Topics are ordered by the Phecode system (see Figure 3). 316 of 348 diseases analysed are

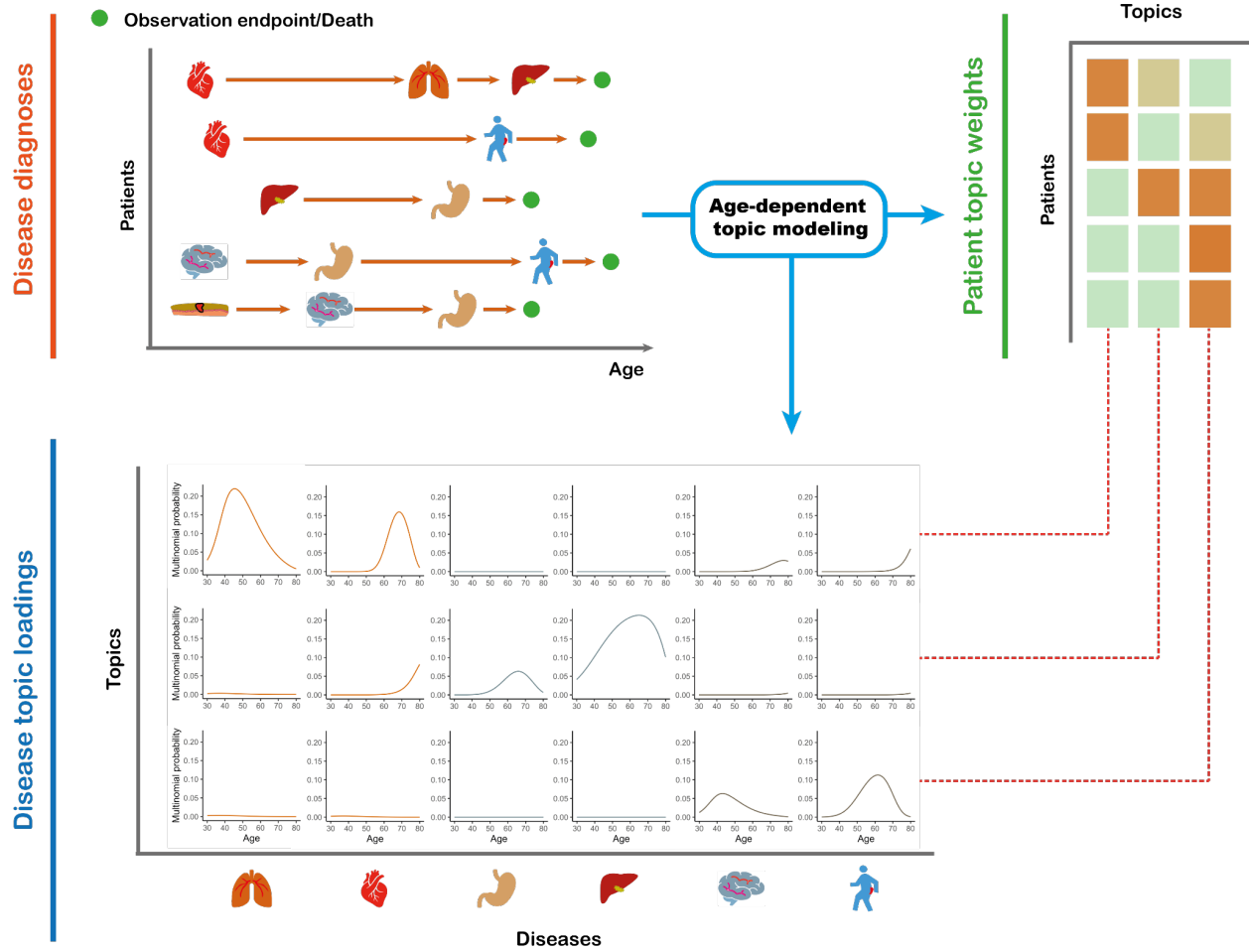
936 associated with a topic; the remaining 32 diseases do not have a topic with average diagnosis-

937 specific topic probability >50%.

938

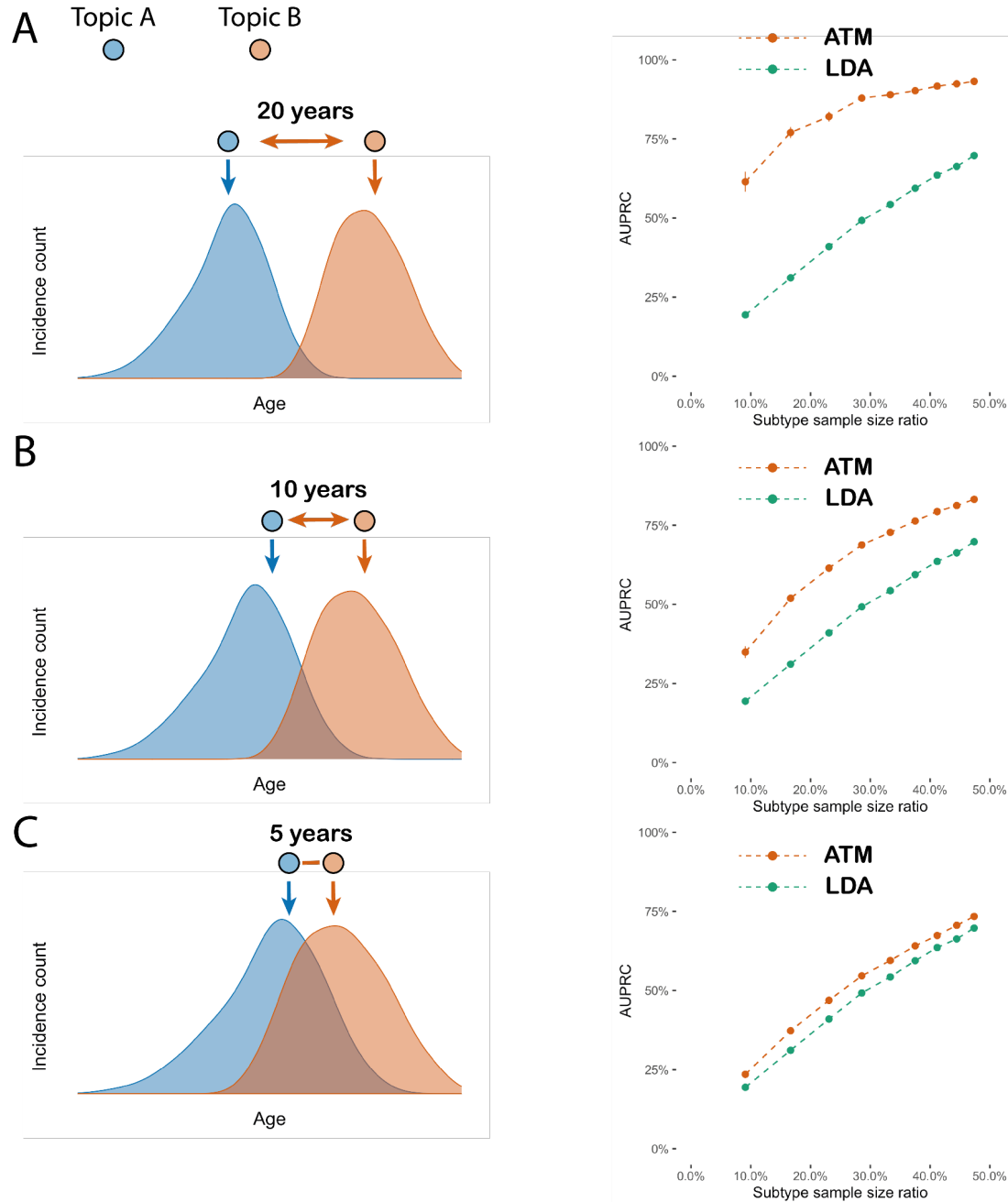
939

940 **Figures**



941 **Figure 1: ATM provides an efficient way to represent longitudinal comorbidity data.** Top
942 left: input consists of disease diagnoses as a function of age. Top right: ATM assigns a topic
943 weight to each patient. Bottom: ATM infers age-dependent topic loadings.
944

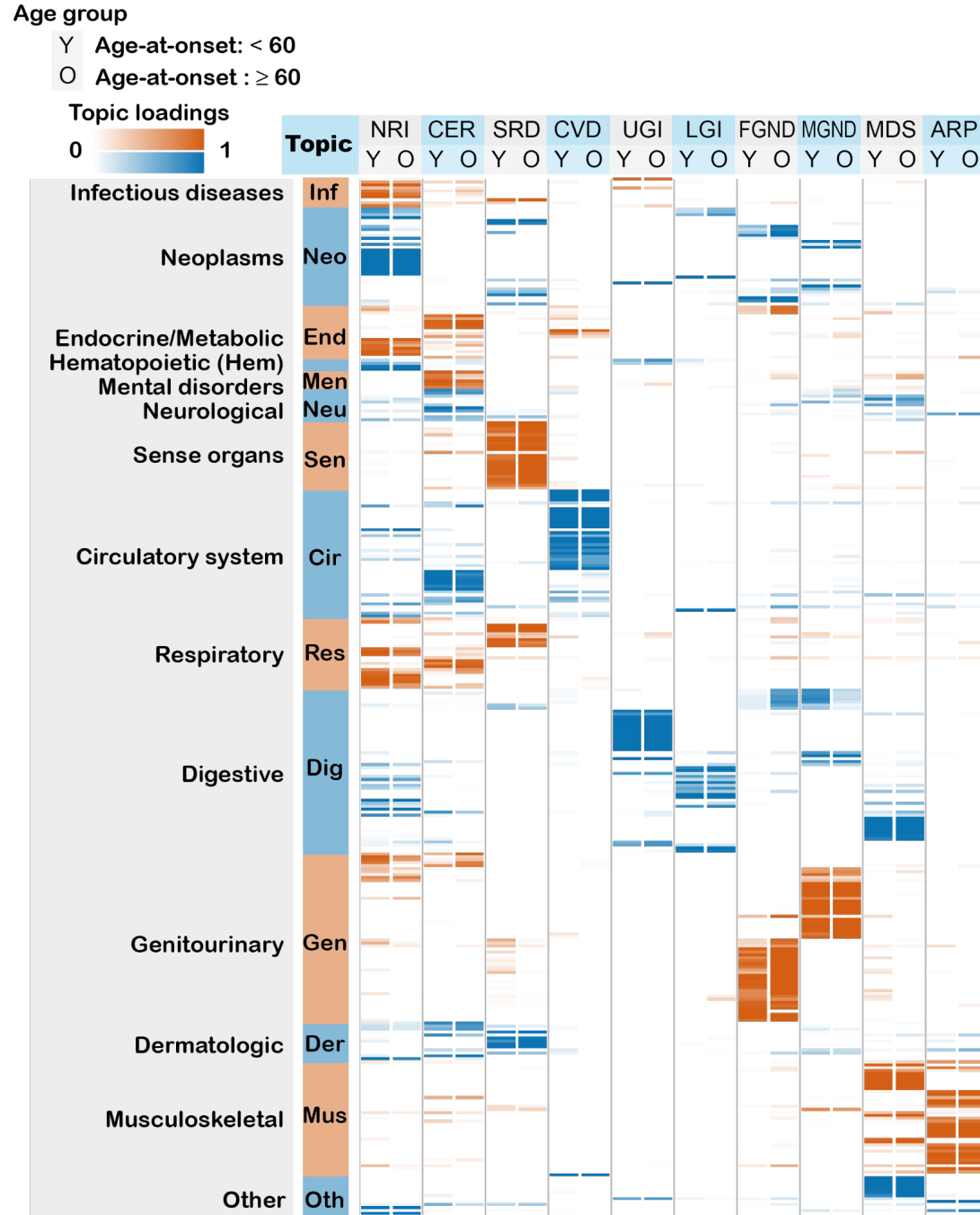
945
946
947
948



949

950 **Figure 2: ATM outperforms LDA in simulations with age-dependent effects.** In simulations
951 at different levels of age-dependent effects (left panels), we report the area under the precision
952 and recall curve (AUPRC) for ATM vs. LDA as a function of subtype sample size proportion
953 (the proportion of diagnoses belonging to the smaller subtype) (right panels). Each dot represents
954 the mean of 100 simulations of 10,000 individuals. Error bars denote 95% confidence intervals.
955 (A) 20-year difference in age at diagnosis for the two subtypes. (B) 10-year difference in age at
956 diagnosis for the two subtypes. (C) 5-year difference in age at diagnosis for the two subtypes.
957 Numerical results are reported in Supplementary Table 2.

958



959

960

961

962

963

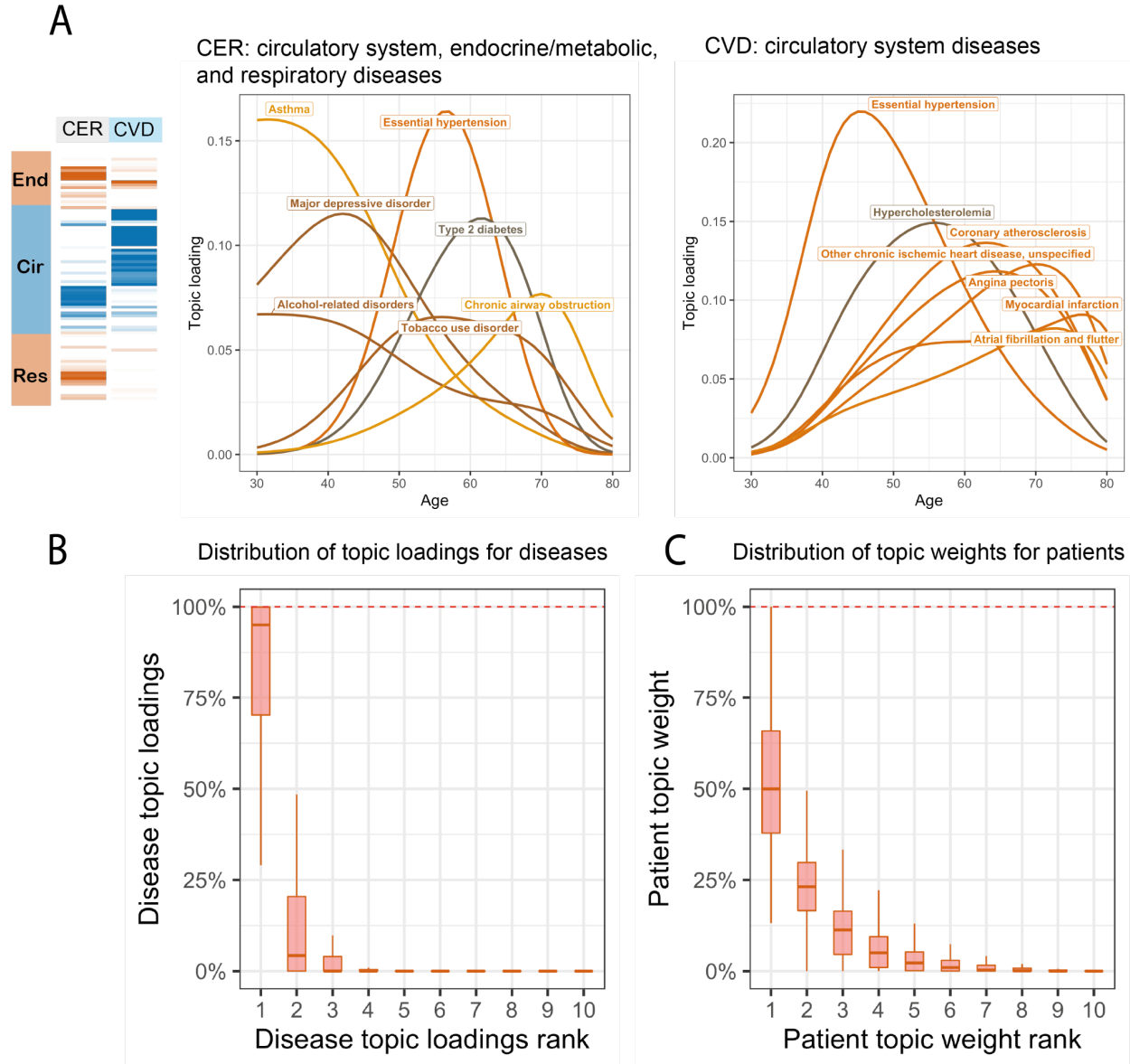
964

965

966

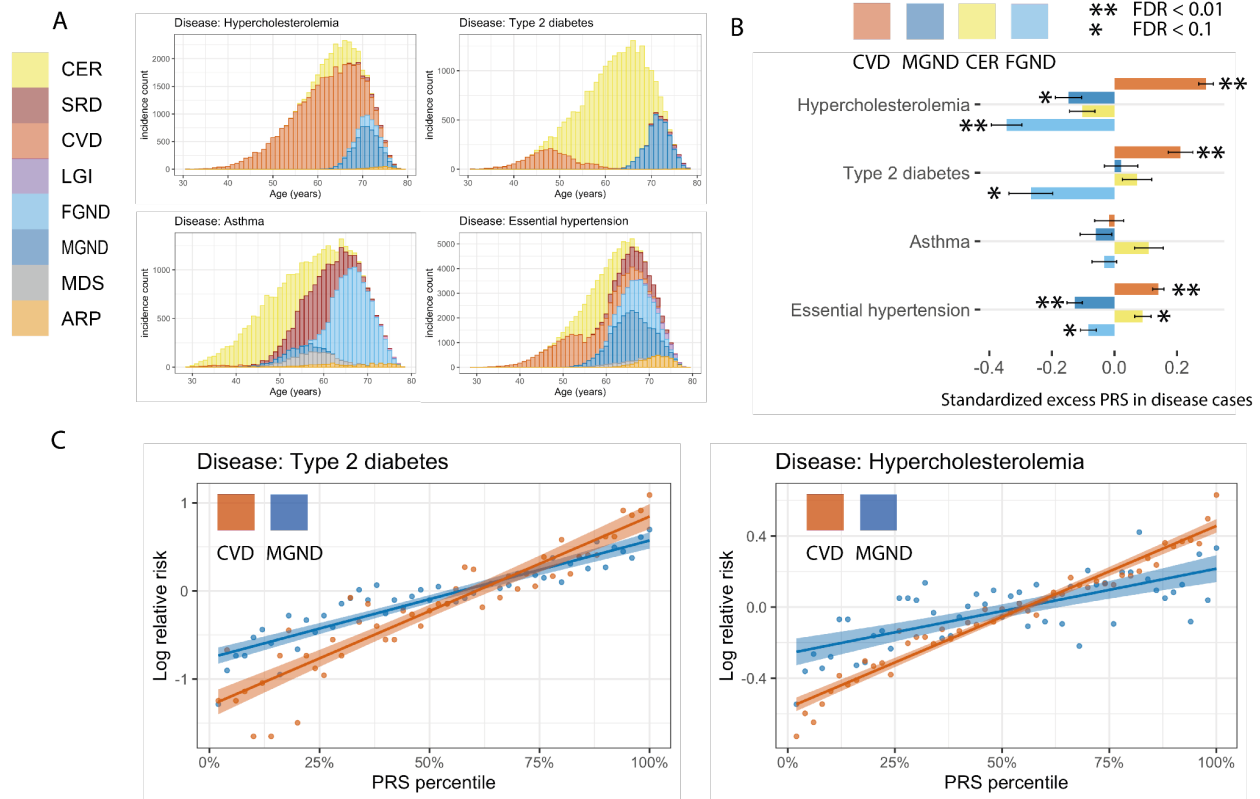
Figure 3. Age-dependent topic loadings of 10 inferred disease topics across 348 diseases in the UK Biobank. We report topic loadings averaged across younger ages (age at diagnosis < 60) and older ages (age at diagnosis > 60). Row labels denote disease categories ordered by Phecode systems, with alternating blue and red color for visualisation purposes; “Other” is a merge of five Phecode systems: “congenital anomalies”, “symptoms”, “injuries & poisoning”, “other tests”, and “death” (which is treated as an additional disease, see Methods). Topics are ordered by the corresponding Phecode system. Further details on the 10 topics are provided in Table 1. Further

967 details on the diseases discussed in the text (type 2 diabetes and breast cancer) are provided in
968 Supplementary Figure 6. Numerical results are reported in Supplementary Table 4.
969
970



971
972 **Figure 4. Topic loadings capture age-dependent comorbidities.** (A) Age-dependent topic
973 loadings for two representative topics, CER and CVD; for each topic, we include the top seven
974 diseases with highest topic loadings. Results for all 10 topics are reported in Supplementary
975 Figure 13. (B) Box plot of disease topic loading as a function of rank; disease topic loadings are
976 computed as a weighted average across all values of age at diagnosis. (C) Box plot of patient
977 topic weight as a function of rank. Numerical results are reported in Supplementary Table 5.
978

979



980

981 **Figure 5. Polygenic risk scores vary across disease subtypes defined by distinct topics. (A)**

982 Stacked barplots of age-dependent subtypes (defined by topics) for 4 representative diseases

983 (type 2 diabetes, asthma, hypercholesterolemia, and essential hypertension); for each disease, we

984 include all subtypes with at least one diagnosis. Results for all 52 diseases are reported in

985 Supplementary Figure 16. (B) Standardised excess PRS values in disease cases (s.d. increase in

986 PRS per unit increase in patient topic weight) for 4 representative diseases and 4 corresponding

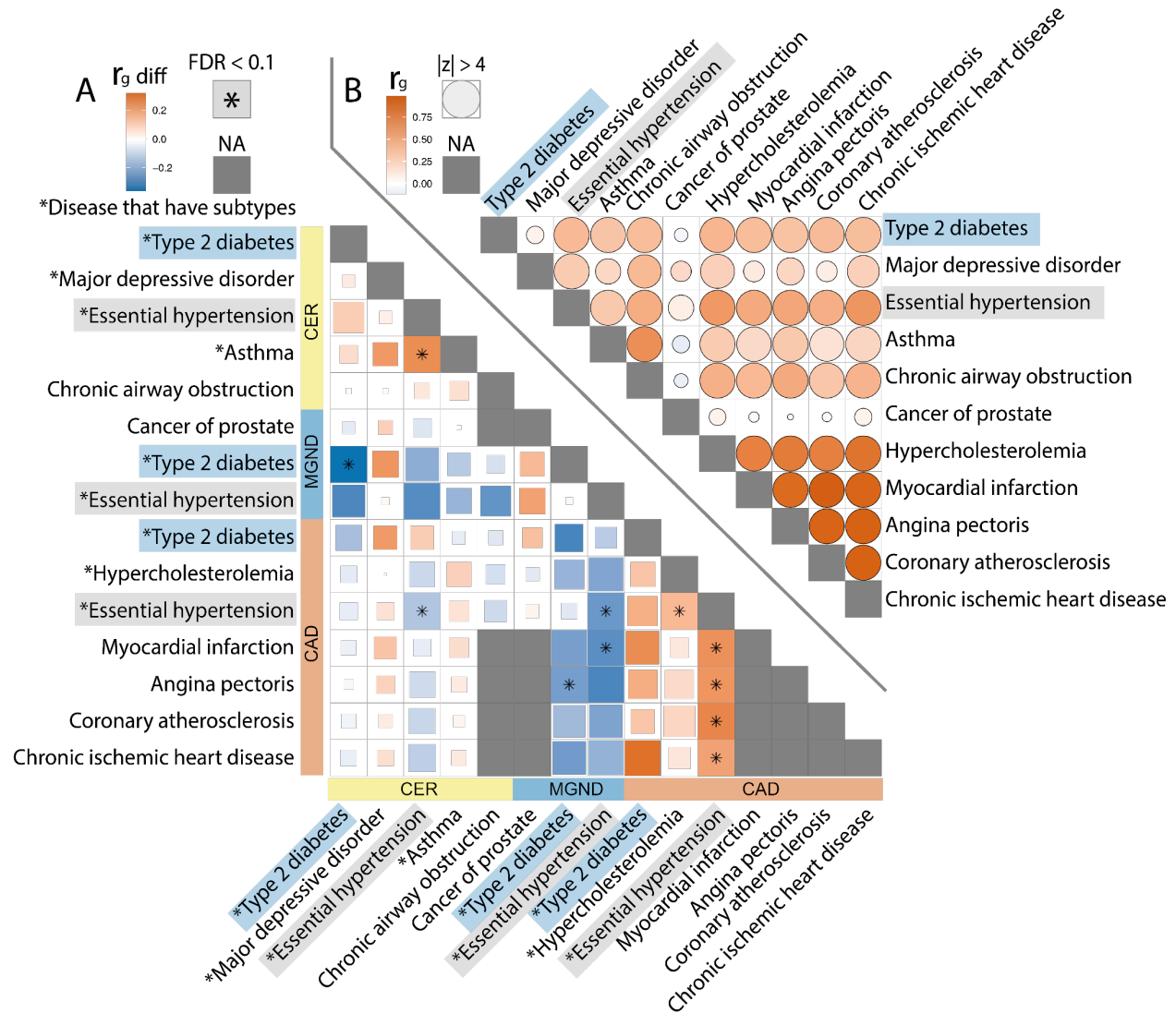
987 topics. (C) Relative risk for cases of type 2 diabetes and hypercholesterolemia of CVD and

988 MGND subtypes (vs. controls) across PRS percentiles. Each point spans 2 PRS percentiles.

989 Lines denote regression on log scale. Error bars denote 95% confidence intervals. Numerical

990 results are reported in Supplementary Table 6.

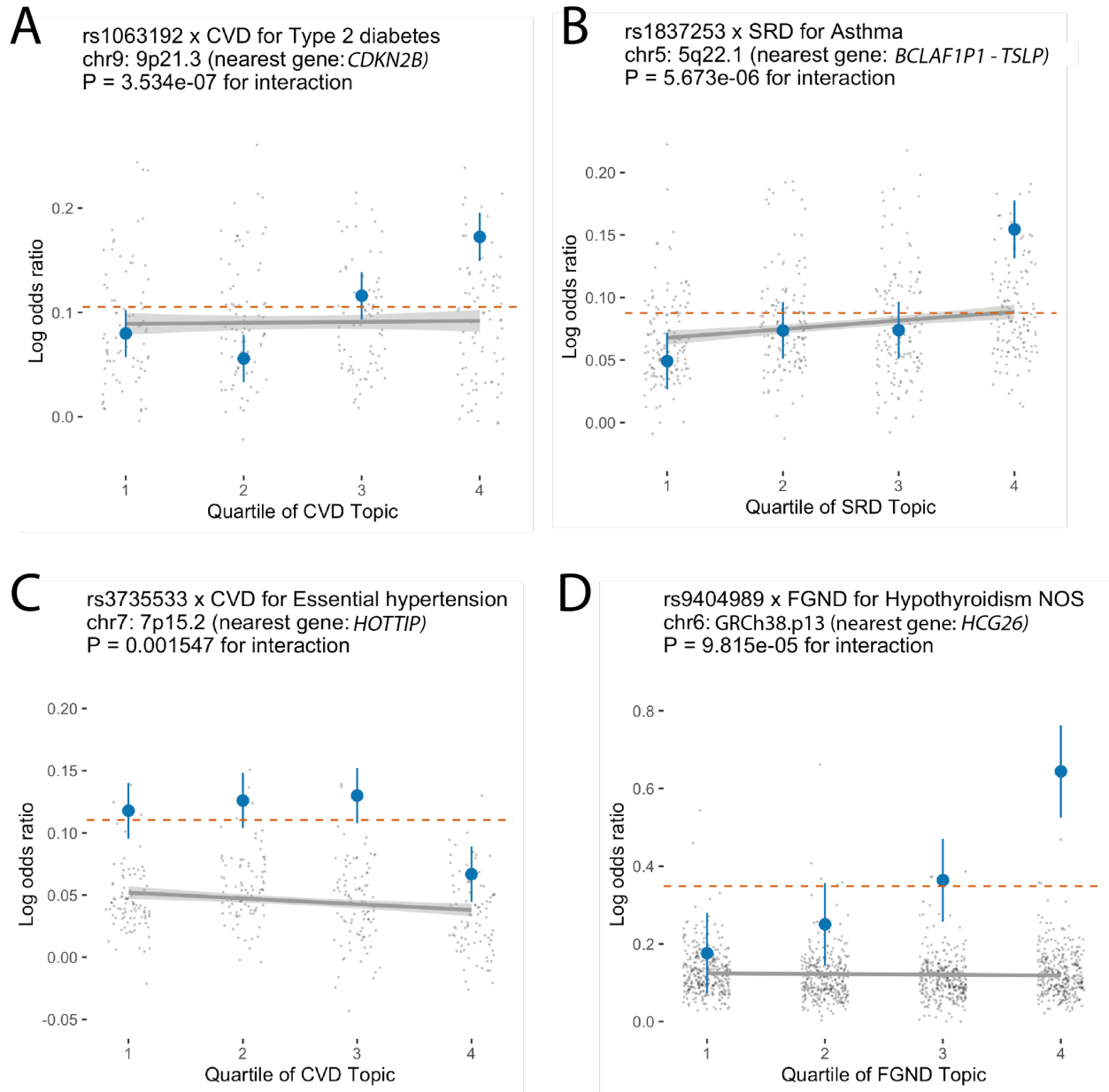
991



992
993
994
995
996
997
998
999
1000
1001
1002

Figure 6. Genetic correlations vary across disease subtypes defined by distinct topics. (a)

Excess genetic correlations for pairs of 15 disease subtypes or diseases (9 disease subtypes (denoted with asterisks) + 6 diseases without subtypes), relative to genetic correlations between the underlying diseases. Full square with asterisk denotes FDR < 0.1; less than full squares have area proportional to z-scores for difference. Grey squares denote NA (pair of diseases without subtypes or pair of same disease subtype or disease). (b) Genetic correlations between the underlying diseases. Full circle denotes $|z\text{-score}| > 4$ for nonzero genetic correlation; less than full circles have area proportional to $|z\text{-score}|$. Numerical results are reported in Supplementary Table 9.



1003
1004
1005
1006
1007
1008
1009
1010
1011
1012

Figure 7. Examples of SNP x topic interaction effects on disease phenotypes. For each example, we report main SNP effects (log odds ratios) specific to each quartile of topic weights across individuals, for both the focal SNP (blue dots) and background SNPs for that disease and topic (genome-wide significant main effect ($P < 5 \times 10^{-8}$) but non-significant SNP x topic interaction effect ($P > 0.05$); grey dots). Dashed red lines denote aggregate main SNP effects for each focal SNP. Error bars denote 95% confidence intervals. Grey lines denote linear regression of grey dots, with grey shading denoting corresponding 95% confidence intervals. Numerical results are reported in Supplementary Table 14.

1013 References

- 1014 1. Abul-Husn, N. S. & Kenny, E. E. Personalized Medicine and the Power of Electronic Health
1015 Records. *Cell* **177**, 58–69 (2019).
- 1016 2. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits.
1017 *Nat. Genet.* **47**, 1236–1241 (2015).
- 1018 3. Wang, K., Gaitsch, H., Poon, H., Cox, N. J. & Rzhetsky, A. Classification of common
1019 human diseases derived from shared genetic and environmental determinants. *Nat. Genet.*
1020 **49**, 1319–1325 (2017).
- 1021 4. Zhao, W. *et al.* Identification of new susceptibility loci for type 2 diabetes and shared
1022 etiological pathways with coronary heart disease. *Nat. Genet.* **49**, 1450–1457 (2017).
- 1023 5. Zhu, Z. *et al.* A genome-wide cross-trait analysis from UK Biobank highlights the shared
1024 genetic architecture of asthma and allergic diseases. *Nat. Genet.* **50**, 857–864 (2018).
- 1025 6. Turley, P. *et al.* Multi-trait analysis of genome-wide association summary statistics using
1026 MTAG. *Nat. Genet.* **50**, 229–237 (2018).
- 1027 7. O'Connor, L. J. & Price, A. L. Distinguishing genetic correlation from causation across 52
1028 diseases and complex traits. *Nat. Genet.* **50**, 1728–1734 (2018).
- 1029 8. Cortes, A., Albers, P. K., Dendrou, C. A., Fugger, L. & McVean, G. Identifying cross-
1030 disease components of genetic risk across hospital data in the UK Biobank. *Nat. Genet.* **52**,
1031 126–134 (2019).
- 1032 9. Morrison, J., Knoblauch, N., Marcus, J. H., Stephens, M. & He, X. Mendelian randomization
1033 accounting for correlated and uncorrelated pleiotropic effects using genome-wide summary
1034 statistics. *Nat. Genet.* **52**, 740–747 (2020).
- 1035 10. Peyrot, W. J. & Price, A. L. Identifying loci with different allele frequencies among cases of
1036 eight psychiatric disorders using CC-GWAS. *Nat. Genet.* **53**, 445–454 (2021).
- 1037 11. Mattheisen, M. *et al.* Identification of shared and differentiating genetic architecture for

- 1038 autism spectrum disorder, attention-deficit hyperactivity disorder and case subgroups. *Nat.*
1039 *Genet.* **54**, 1470–1478 (2022).
- 1040 12. Cortes, A. *et al.* Bayesian analysis of genetic association across tree-structured routine
1041 healthcare data in the UK Biobank. *Nat. Genet.* **49**, 1311–1318 (2017).
- 1042 13. Zhang, H. *et al.* Genome-wide association study identifies 32 novel breast cancer
1043 susceptibility loci from overall and subtype-specific analyses. *Nat. Genet.* **52**, 572–581
1044 (2020).
- 1045 14. Mansour Aly, D. *et al.* Genome-wide association analyses highlight etiological differences
1046 underlying newly defined subtypes of diabetes. *Nat. Genet.* **53**, 1534–1542 (2021).
- 1047 15. Hautakangas, H. *et al.* Genome-wide analysis of 102,084 migraine cases identifies 123 risk
1048 loci and subtype-specific risk alleles. *Nat. Genet.* **54**, 152–160 (2022).
- 1049 16. Srebro, N. & Shraibman, A. Rank, Trace-Norm and Max-Norm. in *Learning Theory* 545–
1050 560 (Springer Berlin Heidelberg, 2005).
- 1051 17. Candès, E. & Recht, B. Exact matrix completion via convex optimization. *Commun. ACM*
1052 **55**, 111–119 (2012).
- 1053 18. Yan, J. & Pollefeys, M. A General Framework for Motion Segmentation: Independent,
1054 Articulated, Rigid, Non-rigid, Degenerate and Non-degenerate. in *Computer Vision – ECCV*
1055 *2006* 94–106 (Springer Berlin Heidelberg, 2006).
- 1056 19. Ma, Y., Derksen, H. & Hong, W. Segmentation of multivariate mixed data via Lossy data
1057 coding and compression. *IEEE Trans. Pattern Anal. Mach. Intell.* **29**, 1546–1562 (2007).
- 1058 20. Rao, S., Tron, R., Vidal, R. & Ma, Y. Motion segmentation in the presence of outlying,
1059 incomplete, or corrupted trajectories. *IEEE Trans. Pattern Anal. Mach. Intell.* **32**, 1832–
1060 1845 (2010).
- 1061 21. Liu, G. & Yan, S. Latent Low-Rank Representation for subspace segmentation and feature
1062 extraction. in *2011 International Conference on Computer Vision* 1615–1622
1063 (ieeexplore.ieee.org, 2011).

- 1064 22. Liu, Z. *et al.* Efficient Low-rank Multimodal Fusion with Modality-Specific Factors. *arXiv*
1065 *[cs.AI]* (2018).
- 1066 23. Chen, Y. & Chi, Y. Harnessing Structures in Big Data via Guaranteed Low-Rank Matrix
1067 Estimation: Recent Theory and Fast Algorithms via Convex and Nonconvex Optimization.
1068 *IEEE Signal Process. Mag.* **35**, 14–31 (2018).
- 1069 24. Chen, Y. & Candès, E. J. The projected power method: An efficient algorithm for joint
1070 alignment from pairwise differences. *Commun. Pure Appl. Math.* **71**, 1648–1714 (2018).
- 1071 25. Jia, G. *et al.* Estimating heritability and genetic correlations from large health datasets in the
1072 absence of genetic data. *Nat. Commun.* **10**, 1–11 (2019).
- 1073 26. Tanigawa, Y. *et al.* Components of genetic associations across 2,138 phenotypes in the UK
1074 Biobank highlight adipocyte biology. *Nat. Commun.* **10**, 4064 (2019).
- 1075 27. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human
1076 phenotypes. *Nat. Genet.* **53**, 1415–1424 (2021).
- 1077 28. Jia, G. *et al.* Discerning asthma endotypes through comorbidity mapping. *Nat. Commun.*
1078 **13**, 1–19 (2022).
- 1079 29. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
1080 *Nature* **562**, 203–209 (2018).
- 1081 30. of Us Research Program, A. The ‘All of Us’ Research Program. *N. Engl. J. Med.* **381**, 668–
1082 676 (2019).
- 1083 31. Ishigaki, K. *et al.* Large-scale genome-wide association study in a Japanese population
1084 identifies novel susceptibility loci across different diseases. *Nat. Genet.* **52**, 669–679
1085 (2020).
- 1086 32. Siggaard, T. *et al.* Disease trajectory browser for exploring temporal, population-wide
1087 disease progression patterns in 7.2 million Danish patients. *Nat. Commun.* **11**, 1–10 (2020).
- 1088 33. Posey, J. E. *et al.* Resolution of Disease Phenotypes Resulting from Multilocus Genomic
1089 Variation. *N. Engl. J. Med.* **376**, 21–31 (2017).

- 1090 34. Cook, E. K. *et al.* Comorbid and inflammatory characteristics of genetic subtypes of clonal
1091 hematopoiesis. *Blood Adv* **3**, 2482–2486 (2019).
- 1092 35. Udler, M. S. *et al.* Type 2 diabetes genetic loci informed by multi-trait associations point to
1093 disease mechanisms and subtypes: A soft clustering analysis. *PLoS Med.* **15**, e1002654
1094 (2018).
- 1095 36. Wani, B., Aziz, S. A., Ganaie, M. A. & Mir, M. H. Metabolic Syndrome and Breast Cancer
1096 Risk. *Indian J. Med. Paediatr. Oncol.* **38**, 434–439 (2017).
- 1097 37. Blei, Ng & Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.* (2003).
- 1098 38. Pritchard, J. K., Stephens, M. & Donnelly, P. Inference of population structure using
1099 multilocus genotype data. *Genetics* **155**, 945–959 (2000).
- 1100 39. Bishop, C. M. *Pattern Recognition and Machine Learning*. (Springer New York, 2006).
- 1101 40. Teh, Y., Newman, D. & Welling, M. A collapsed variational Bayesian inference algorithm for
1102 latent Dirichlet allocation. *Adv. Neural Inf. Process. Syst.* **19**, (2006).
- 1103 41. Grau, J., Grosse, I. & Keilwagen, J. PRROC: computing and visualizing precision-recall and
1104 receiver operating characteristic curves in R. *Bioinformatics* **31**, 2595–2597 (2015).
- 1105 42. Wu, P. *et al.* Mapping ICD-10 and ICD-10-CM Codes to Phecodes: Workflow Development
1106 and Initial Evaluation. *JMIR Med Inform* **7**, e14325 (2019).
- 1107 43. Suvila, K. *et al.* Early Onset Hypertension Is Associated With Hypertensive End-Organ
1108 Damage Already by MidLife. *Hypertension* HYPERTENSIONAHA11913069 (2019).
- 1109 44. Wong, B. *et al.* Cardiovascular Disease Risk Associated With Familial
1110 Hypercholesterolemia: A Systematic Review of the Literature. *Clin. Ther.* **38**, 1696–1709
1111 (2016).
- 1112 45. Shah, M. S. & Brownlee, M. Molecular and Cellular Mechanisms of Cardiovascular
1113 Disorders in Diabetes. *Circ. Res.* **118**, 1808–1829 (2016).
- 1114 46. Shah, A. D. *et al.* Type 2 diabetes and incidence of cardiovascular diseases: a cohort study
1115 in 1.9 million people. *The Lancet Diabetes & Endocrinology* **3**, 105–113 (2015).

- 1116 47. Dabelea, D. & Hamman, R. F. Elevated Cardiometabolic Risk Profile Among Young Adults
1117 With Diabetes: Need for Action. *Diabetes care* vol. 42 1845–1846 (2019).
- 1118 48. Wong, J. L. & Evans, S. E. Bacterial Pneumonia in Patients with Cancer: Novel Risk
1119 Factors and Management. *Clin. Chest Med.* **38**, 263–277 (2017).
- 1120 49. Falstie-Jensen, A. M. *et al.* Incidence of hypothyroidism after treatment for breast cancer—
1121 a Danish matched cohort study. *Breast Cancer Res.* **22**, 1–10 (2020).
- 1122 50. Loh, P.-R., Kichaev, G., Gazal, S., Schoech, A. P. & Price, A. L. Mixed-model association
1123 for biobank-scale datasets. *Nat. Genet.* **50**, 906–908 (2018).
- 1124 51. Loh, P.-R. *et al.* Efficient Bayesian mixed-model analysis increases association power in
1125 large cohorts. *Nat. Genet.* **47**, 284–290 (2015).
- 1126 52. Jiang, X., Holmes, C. & McVean, G. The impact of age on genetic risk for common
1127 diseases. *PLoS Genet.* **17**, e1009723 (2021).
- 1128 53. Weir, B. S. & Cockerham, C. C. ESTIMATING F-STATISTICS FOR THE ANALYSIS OF
1129 POPULATION STRUCTURE. *Evolution* **38**, 1358–1370 (1984).
- 1130 54. Bhatia, G., Patterson, N., Sankararaman, S. & Price, A. L. Estimating and interpreting FST:
1131 the impact of rare variants. *Genome Res.* **23**, 1514–1521 (2013).
- 1132 55. Scott, L. J. *et al.* A genome-wide association study of type 2 diabetes in Finns detects
1133 multiple susceptibility variants. *Science* **316**, 1341–1345 (2007).
- 1134 56. Maller, J. B. *et al.* Bayesian refinement of association signals for 14 loci in 3 common
1135 diseases. *Nat. Genet.* **44**, 1294–1301 (2012).
- 1136 57. Cheng, M., An, S. & Li, J. CDKN2B-AS may indirectly regulate coronary artery disease-
1137 associated genes via targeting miR-92a. *Gene* **629**, 101–107 (2017).
- 1138 58. Scott, R. A. *et al.* An Expanded Genome-Wide Association Study of Type 2 Diabetes in
1139 Europeans. *Diabetes* **66**, 2888–2902 (2017).
- 1140 59. Aragam, K. G. *et al.* Discovery and systematic characterization of risk variants and genes
1141 for coronary artery disease in over a million participants. *bioRxiv* (2021)

- 1142 doi:10.1101/2021.05.24.21257377.
- 1143 60. Indra, A. K. Epidermal TSLP: a trigger factor for pathogenesis of atopic dermatitis. *Expert*
1144 *Rev. Proteomics* **10**, 309–311 (2013).
- 1145 61. Kanai, M. *et al.* Genetic analysis of quantitative traits in the Japanese population links cell
1146 types to complex human diseases. *Nat. Genet.* **50**, 390–400 (2018).
- 1147 62. Oluwafemi, O. O. *et al.* Genome-Wide Association Studies of Conotruncal Heart Defects
1148 with Normally Related Great Vessels in the United States. *Genes* **12**, (2021).
- 1149 63. Blei, D. M. & Lafferty, J. D. A correlated topic model of Science. *Ann. Appl. Stat.* **1**, 17–35
1150 (2007).
- 1151 64. Zaitlen, N. *et al.* Informed conditioning on clinical covariates increases power in case-
1152 control association studies. *PLoS Genet.* **8**, e1003032 (2012).
- 1153 65. Sun, B. B. *et al.* Genetic regulation of the human plasma proteome in 54,306 UK Biobank
1154 participants. *bioRxiv* 2022.06.17.496443 (2022) doi:10.1101/2022.06.17.496443.
- 1155 66. Mostafavi, H. *et al.* Variable prediction accuracy of polygenic scores within an ancestry
1156 group. *Elife* **9**, e48376 (2020).
- 1157 67. Dumitrescu, L. *et al.* Evidence for age as a modifier of genetic associations for lipid levels.
1158 *Ann. Hum. Genet.* **75**, 589–597 (2011).
- 1159 68. Chatterjee, N., Shi, J. & García-Closas, M. Developing and evaluating polygenic risk
1160 prediction models for stratified disease prevention. *Nat. Rev. Genet.* **17**, 392–406 (2016).
- 1161 69. Lin, J. *et al.* Integration of biomarker polygenic risk score improves prediction of coronary
1162 heart disease in UK Biobank and FinnGen. *bioRxiv* (2022)
1163 doi:10.1101/2022.08.22.22279057.
- 1164 70. Ghorbani, B., Javadi, H. & Montanari, A. An Instability in Variational Inference for Topic
1165 Models. in *Proceedings of the 36th International Conference on Machine Learning* (eds.
1166 Chaudhuri, K. & Salakhutdinov, R.) vol. 97 2221–2231 (PMLR, 09--15 Jun 2019).
- 1167 71. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer

1168 datasets. *Gigascience* **4**, 7 (2015).

1169 72. International HapMap 3 Consortium *et al.* Integrating common and rare genetic variation in
1170 diverse human populations. *Nature* **467**, 52–58 (2010).

1171