

1 **Title:** A deep learning transformer model predicts high rates of undiagnosed rare disease in large electronic health
2 systems

3

4 **Authors:** Daniel M. Jordan, Ha My T. Vy, Ron Do* .

5

6 **Affiliation:** Center for Genomic Data Analytics, Charles Bronfman Institute for Personalized Medicine,
7 Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York,
8 NY, USA

9 * Corresponding author

10

11 **Correspondence:**

12 Ron Do, PhD

13 Annenberg Building, Floor 18 Room 80A

14 1468 Madison Ave

15 New York, NY-10029

16 Phone Number: 212-241-6206 | Fax Number: 212-849-2643

17 Email: ron.do@mssm.edu

18

19

20 **Abstract**

21 It is estimated that as many as 1 in 16 people worldwide suffer from rare diseases. Rare disease patients
22 face difficulty finding diagnosis and treatment for their conditions, including long diagnostic odysseys,
23 multiple incorrect diagnoses, and unavailable or prohibitively expensive treatments. As a result, it is likely
24 that large electronic health record (EHR) systems include high numbers of participants suffering from
25 undiagnosed rare disease. While this has been shown in detail for specific diseases, these studies are
26 expensive and time consuming and have only been feasible to perform for a handful of the thousands of
27 known rare diseases. The bulk of these undiagnosed cases are effectively hidden, with no straightforward
28 way to differentiate them from healthy controls. The ability to access them at scale would enormously
29 expand our capacity to study and develop drugs for rare diseases, adding to tools aimed at increasing
30 availability of study cohorts for rare disease. In this study, we train a deep learning transformer algorithm,
31 RarePT (Rare-Phenotype Prediction Transformer), to impute undiagnosed rare disease from EHR
32 diagnosis codes in 436,407 participants in the UK Biobank and validated on an independent cohort from
33 3,333,560 individuals from the Mount Sinai Health System. We applied our model to 155 rare diagnosis
34 codes with fewer than 250 cases each in the UK Biobank and predicted participants with elevated risk for
35 each diagnosis, with the number of participants predicted to be at risk ranging from 85 to 22,000 for
36 different diagnoses. These risk predictions are significantly associated with increased mortality for 65%
37 of diagnoses, with disease burden expressed as disability-adjusted life years (DALY) for 73% of
38 diagnoses, and with 72% of available disease-specific diagnostic tests. They are also highly enriched for
39 known rare diagnoses in patients not included in the training set, with an odds ratio (OR) of 48.0 in cross-
40 validation cohorts of the UK Biobank and an OR of 30.6 in the independent Mount Sinai Health System
41 cohort. Most importantly, RarePT successfully screens for undiagnosed patients in 32 rare diseases with
42 available diagnostic tests in the UK Biobank. Using the trained model to estimate the prevalence of
43 undiagnosed disease in the UK Biobank for these 32 rare phenotypes, we find that at least 50% of patients
44 remain undiagnosed for 20 of 32 diseases. These estimates provide empirical evidence of a high

45 prevalence of undiagnosed rare disease, as well as demonstrating the enormous potential benefit of using
46 RarePT to screen for undiagnosed rare disease patients in large electronic health systems.

47

48 **Introduction**

49 Rare diseases, also known as orphan diseases, are defined by the European Union as those affecting fewer
50 than one in 2,000 people, and in the United States as those affecting fewer than 200,000 people
51 nationwide^{1,2}. Rare diseases are collectively very common, and it is estimated that as many as 1 in 16
52 people (6.2%) suffer from one or more rare diseases³. This makes them a serious public health concern, as
53 rare disease patients are far less likely to receive accurate diagnoses or, once diagnosed, to have access to
54 effective treatments for their conditions⁴⁻⁶. This is due to the difficulty of studying rare diseases, a scarcity
55 of clinical expertise and diagnostic methods, as well as the unprofitability of developing drugs targeting
56 them. In fact, many of these diseases are so understudied and underdiagnosed that we do not know with
57 any certainty what their true prevalence is, and how many undiagnosed patients there may be³. One of the
58 primary reasons for all these problems is the difficulty of finding large enough populations of patients to
59 conduct well-powered studies on these diseases, either in the context of basic or translational research or
60 in the context of drug trials. This is a pressing problem, and tools such as MatchMaker Exchange, which
61 help researchers match similar cases to increase sample size for studies of rare diseases, are widely used⁷⁻
62 ⁹. Further development of these tools is also an active area of research, including expanding them to
63 include comparisons of phenotypic features mined from electronic health records (EHR) or imaging data⁹⁻
64 ¹². These tools are vital for rare disease research because undiagnosed rare disease patients masquerade as
65 healthy controls, making them invisible and inaccessible to researchers unless they can be revealed. There
66 is an urgent need for new approaches to reveal people suffering from hidden rare diseases in research and
67 drug trial cohorts, and in clinical practice.

68 In this study, we present such an approach, using a deep learning transformer model trained on
69 EHR data. Artificial intelligence (AI) language models based on the transformer architecture, such as
70 BERT (“Bidirectional Encoder Representations from Transformers” and GPT (“Generative Pretrained
71 Transformers”), have proved very successful at learning the relationships between concepts in natural
72 languages^{13,14}. Transformer models have also been successfully applied to problems in biology that are

73 not directly related to language processing in methods such as AlphaFold-2, AlphaMissense, DeepMAPS,
74 and Enformer¹⁵⁻¹⁸. One of the strength of the transformer architecture is that, with appropriate
75 tokenization and training schemes, transformers can correctly model concepts that are extremely rare, and
76 some transformer-based models have been shown to define a new word after seeing it only a small
77 number of times¹⁹⁻²². We have designed a modified transformer architecture to model phenotypic concepts
78 based on phenotypes derived from structured diagnosis codes in electronic health records (EHR), along
79 with a modified training procedure designed to maximize power to screen for missing rare diagnoses. The
80 resulting model, RarePT (Rare-Phenotype Prediction Transformer), was trained on EHR data from
81 436,407 individuals from the UK Biobank and validated on an independent cohort from 3,333,560
82 individuals from the Mount Sinai Health System in New York City, USA (**Table 1**). RarePT shows
83 remarkable power to recapitulate rare diagnoses, which is robust across different racial and ethnic groups,
84 different hospitals with different coding practices, and even different countries with different health care
85 standards and coding vocabularies. It also detects UK Biobank participants with undiagnosed rare disease,
86 enabling empirical measurement of the true prevalence of undiagnosed cases for rare diseases.

87 **Results**

88 *Model training and cross-validation*

89 We implemented a transformer model with a self-attention mechanism similar to AI language
90 models such as BERT and GPT, along with a “masked diagnosis modeling” training objective by analogy
91 to the “masked language modeling” objective used by some of these language models¹⁹. In this approach,
92 training examples consist of complete sequences with a single token removed, and the model is trained to
93 reconstruct the missing token. In the natural language processing case, this is a sequence of words with a
94 single word removed; in our case, it is a participant record with a single diagnosis removed (**Figure 1a**).
95 The model learns the meanings of tokens based on the context they appear in, resulting in embeddings
96 that cluster tokens that commonly appear together and tokens that appear in similar context. Models
97 trained with this objective are known to learn informative embeddings even for very rare tokens in many

98 cases^{19–22}. We made use of this feature to train a model to predict rare diagnoses, a critical need due to
99 underdiagnosis and understudying of rare diseases.

100 An additional advantage of the masked diagnosis modeling training objective for rare tokens is
101 that it allows us to weight the importance of tokens to the training objective independent of their
102 prevalence in the training corpus. This is because each training example specifies which token the model
103 must predict correctly to be scored as successful, and the model is not necessarily required to predict
104 every token in each example. The importance of each token to the training objective is determined by how
105 many examples have it as the masked token. In order to prevent very common diagnoses from dominating
106 the learned embeddings, we limited training examples to a fixed number of cases and controls for each
107 diagnosis. While we used 100 cases and controls for each diagnosis, in principle this is a tunable
108 parameter of the training process. Lower values allow rarer diagnoses to be included, while higher values
109 increase the amount of training data available.

110 In this study, we express diagnoses as phecodes²³. We determined phecodes from ICD-10 codes for
111 436,407 participants in the UK Biobank based on a standard mapping²⁴. We then filtered out all phecodes
112 with fewer than 100 cases and controls and constructed a training dataset consisting of 100 randomly
113 selected cases and 100 randomly selected controls for each phecode. The resulting training set consisted
114 of 259,400 training examples representing 1,297 query phecodes and 111,331 unique participants. For
115 each training example, input data included the following features:

- 116 1. The identity of the query phecode
- 117 2. All other phecodes for which the participant is considered a case
- 118 3. Age at recruitment
- 119 4. Sex reported from recruitment

120 These training examples were split into 5 subsamples for cross-validation, stratified so that each
121 participant appeared in only one split and so that each split contained a similar number of cases and

122 controls. Neural network architecture and other training hyperparameters were tuned on the training data
123 for each split using the Hyperband algorithm²⁵, and then the tuned model was trained on the same data
124 and tested on the held-out test data; see Methods for details of model tuning and training parameters. The
125 final tuned architecture is shown in in **Figure 1b**; details of training performance can be found in
126 **Supplementary Figure S1** and **Supplementary Table S1**. In general, the models performed well on the
127 test data and showed only minor loss of performance between training and test data.

128 *RarePT predicts rare diagnoses in the UK Biobank*

129 To test the ability of the trained model to predict rare diagnoses, we first selected all phecodes
130 appearing in fewer than 1 in 2,000 UK Biobank participants, corresponding to the definition of rare
131 diseases used by the European Union². There were 155 rare phecodes meeting this criterion, shown in
132 **Supplementary Table S2**. Not all phecodes that are rare in the UK Biobank represent phenotypes that
133 meet the definition of rare diseases in the general population. One reason for this is the known bias of the
134 UK Biobank population towards healthier and older participants^{26–28}, which reduces the apparent
135 prevalence of many diseases, especially severe diseases with early onset. For example, phecode 315.3
136 “mental retardation” appears in fewer than 200 participants in the UK Biobank even though the disorders
137 it represents are much more common in the general population, which is likely because severe childhood
138 disorders are underrepresented in this cohort of healthy adults. It is also likely that many of these rare
139 phecodes correspond to diagnosis codes that rarely appear in electronic health records (EHR) despite the
140 conditions they refer to being common. Likely examples of this include 367.4 “presbyopia” and 523.1
141 “gingivitis.” Nevertheless, even if not all of these phecodes represent phenotypes that are rare in the
142 general population, they do represent phenotypes that are rare in the data used to train our model, and the
143 model’s performance on these phecodes is informative about how our methodology handles rare
144 phenotypes. In total, 21,636 of the 436,407 participants tested have one or more of these rare diagnoses,
145 giving them a cumulative prevalence of 5.0%. This matches the estimated cumulative prevalence of 1.5-

146 6.2% for rare diseases in the general population³, suggesting that our selection of rare phecodes does
147 accurately capture the population distribution of rare diseases.

148 After training on a 111,311-participant subset of UK Biobank data constructed to force each
149 phecode to have prevalence of 50%, we measured RarePT's performance in the full UK Biobank dataset
150 of 436,407 participants. We arbitrarily chose a threshold of 0.95 in the model's probability score output,
151 so that participants with a score of 0.95 or higher in a given phecode were treated as positive predictions
152 for that phecode. With this definition, across all five cross-validated models, we generate specific positive
153 predictions for each of our 155 rare phecodes. The number of positive predictions varied by phecode,
154 ranging between 85 and 22,000 with a median of 2,135 positive predictions per phecode (**Supplementary**
155 **Table 3**). These positive predictions are broadly distributed across participants rather than being
156 concentrated in a small group of unhealthy participants, with no participant receiving more than 29
157 positive predictions and 41% of participants (177,484) receiving a positive prediction for at least one of
158 the 155 phecodes. **Figure 2** shows the performance of the 5 cross-validated models at predicting rare
159 phecodes in the full dataset, excluding each model's training data. We measured prediction performance
160 using diagnostic odds ratio (OR), defined as the ratio between the odds of a participant having a diagnosis
161 given a positive prediction from the model and the odds of a participant having a diagnosis given a
162 negative prediction from the model. The median OR for a positive prediction across all 155 rare phecodes
163 and across the five models trained in cross-validation was 48.0. Some specific phecodes reached a median
164 OR over 20,000, and the lowest median OR for any rare phecode was 5.13 (**Figure 2a, Supplementary**
165 **Table S3**). These values compare favorably to many commonly used diagnostic tests, where diagnostic
166 odds ratios in the range of 20-50 are considered very good^{29,30}. Similarly, the positive predictive value
167 (PPV) for cases is nearly 40% for some phecodes, which is well within the range of a useful screening test
168 (**Supplementary Figure S2, Supplementary Table S3**). Because PPV depends on the prevalence of the
169 condition within the test population, we expect this number to increase further when applying this method
170 in situations where the prior expectation of encountering a given diagnosis is increased, such as in

171 patients with undiagnosed rare conditions or patients who carry rare genetic variants. Importantly, the
172 predictions are able to distinguish not only between cases and controls but also between cases for one
173 phecode and cases for another, indicating that RarePT is making specific predictions for each phecode
174 rather than measuring general health (**Figure 2b**).

175 *Model trained on UK Biobank is predictive in an independent EHR cohort*

176 We applied the trained RarePT model to an independent dataset derived from the Mount Sinai
177 Data Warehouse (MSDW), consisting of anonymized EHR for a cohort of 3,333,560 patients seen in the
178 Mount Sinai Health System in New York City. We determined phecodes for these participants in the same
179 way as for the UK Biobank participants, but using a mapping designed for the US clinical modification to
180 ICD-10 (ICD-10CM) rather than the international standard ICD-10 system used by UK hospitals²⁴. 151 of the
181 155 phecodes determined to be rare in the UK Biobank cohort were present in the MSDW cohort. As a
182 health system based cohort, this cohort is expected to be significantly less healthy than the UK Biobank³¹,
183 and therefore we expect most phecodes to have higher prevalence than in the UK Biobank cohort.
184 Nevertheless, a majority of the rare phecodes we tested (86/151; 57%) still had prevalence less than 1 in
185 2,000 in the MSDW cohort (**Supplementary Table S2**). Likewise, we expect more positive predictions
186 for each phecode, both due to the dataset being over 7-fold larger and due to participants being less
187 healthy in general.

188 For these phecodes in the MSDW cohort, RarePT produced between 100 and 721,000 positive
189 predictions per phecode, with a median of 11,500, and produced at least one positive prediction in 47% of
190 participants (1,518,757). These predictions performed similarly to the predictions for the UK Biobank
191 cohort, with a median OR of 30.6 across all 151 phecodes (**Figure 2a, Supplementary Table S4**).
192 Performance for individual phecodes was also strongly correlated across the two datasets (Pearson $r =$
193 0.456 , $p = 5.03 \times 10^{-40}$, t-test; **Figure 2c**). The fact that performance is similar across the two datasets
194 indicates that RarePT's predictions are based on features that are robust to different methodologies for
195 sample ascertainment and data collection, rather than features that are only informative in the specialized

196 context of the UK Biobank. This replication is especially remarkable given the extensive differences
197 between the two cohorts: in addition to one being a population-based cohort of healthy volunteers and the
198 other being a health system cohort, these cohorts are also from different countries with different standard
199 medical practices, different billing structures and coding systems, and different distributions of race,
200 ethnicity, and genetic ancestry. This indicates the wide applicability of the RarePT method and suggests
201 that its performance does not depend on specific features of diagnosis coding in a particular health
202 system.

203 *Rare disease predictions are associated with mortality, disease burden, and known diagnostic biomarkers*

204 To further demonstrate that RarePT is capturing clinically relevant signals of disease rather than
205 bioinformatic artifacts related to diagnosis coding, we performed regression analyses to test the
206 association of positive predictions with mortality, disability, and, where available, known diagnostic
207 biomarkers. We retrieved the latest mortality data for UK Biobank participants as of October 2023, and
208 performed Cox proportional hazard regression to test whether a positive prediction is associated with
209 mortality, controlling for age, sex, and self-reported ethnicity. 101 phecodes (65% of phecodes tested) had
210 a significant association ($p < 0.05$) with increased mortality, of which 93 (60%) remained significant after
211 Bonferroni correction for 155 phecodes ($p < 0.00032$). The median phecode had a regression coefficient
212 of 0.70, corresponding to a hazard ratio of 2.01, or a twofold increase in mortality rate (**Supplementary**
213 **Table S5**).

214 Next, we estimated Disability Adjusted Life Years (DALY) and its two components, Years of Life
215 Lost (YLL) and Years Living with Disability (YLD), for 80 conditions for all UK Biobank participants³².
216 These measurements represent the number of years lost to both mortality and disability as a result of
217 illness and are used as a measure of disease burden, particularly in the Global Burden of Disease study³³.
218 We performed linear regressions with DALY, YLD, and YLL as the dependent variables to test whether a
219 positive prediction is associated with greater disease burden. In all of these regressions, we controlled for
220 age, sex, and self-reported ethnicity. 113 phecodes (73% of phecodes tested) had a significant association

221 (p < 0.05) with increased estimated DALY, and 106 (68%) remained significant after Bonferroni
222 correction for 155 phecodes (p < 0.00032). 134 phecodes (87%) had a significant association with
223 increased estimated YLD individually, 133 (86%) after Bonferroni correction; 110 phecodes (71%) had a
224 significant association with increased estimated YLL individually, 106 (68%) after Bonferroni correction.
225 For the median phecode, a positive prediction was associated with an increase in estimated DALY of 1.1
226 years (**Supplementary Table S5**).

227 To identify diagnostic biomarkers, we used the SNOMED-CT vocabulary of clinical terms^{34,35} to
228 identify phenotypes whose clinical definition includes laboratory tests that are available for large numbers
229 of participants in the UK Biobank. We identified 75 defined relationships between 32 rare phecodes and
230 23 laboratory tests (**Supplementary Table S6**). These tests were performed as part of the UK Biobank
231 recruitment process and were generally not returned to participants or their physicians, so the availability
232 of a test result does not indicate that it was ordered by a physician and the result of a test was not visible
233 to the physicians responsible for entering diagnoses into the participants' EHR. Since RarePT makes its
234 predictions using only diagnosis codes and has no access to physician-ordered laboratory tests except
235 through diagnosis codes, this means our model's predictions are independent of these test results. This is
236 in contrast to health system based cohorts, including our MSDW cohort, where diagnostic tests are
237 ordered and administered in the context of treating the patient, so that the presence and timing of a test are
238 informative about the judgment of the health care providers and the test result forms part of the diagnostic
239 criteria³⁶.

240 For each of these 75 relationships, we performed a logistic regression to test whether a confident
241 case prediction is associated with abnormal test results, again controlling for age, sex, and ethnicity. 54 of
242 these regressions, representing 72% of these relationships, had a result that was in the expected direction
243 and statistically significant (p < 0.05), and 45 (60%) remained significant after Bonferroni correction for
244 75 regressions (p < 0.00067). The median regression coefficient was 0.57, corresponding to an OR of
245 1.77. In other words, for the median diagnostic test, a participant with a positive prediction from RarePT

246 had 77% higher odds of having an abnormal test result. In 100 random permutations of pcode-
247 laboratory test relationships, no permutation showed as many Bonferroni-significant associations ($p <$
248 0.01) (**Figure 3a, Supplementary Table S7-S8**). We additionally performed linear regression for each of
249 these relationships, testing for a relationship between the model prediction and the quantitative test result.
250 43 of these regressions, representing 57% of these relationships, had a result that was in the expected
251 direction and statistically significant, and 38 (51%) remained significant after Bonferroni correction.
252 Again, 0 out of 100 random permutations showed as many Bonferroni-significant associations ($p < 0.01$)
253 (**Figure 3a, Supplementary Table S7-8**).

254 Taken together, these analyses demonstrate that positive predictions from RarePT do not merely
255 predict diagnosis codes for rare diseases, but also capture clinically and biologically relevant features
256 relevant to the diagnoses and to health outcomes more generally.

257 **Disease predictions suggest high rates of underdiagnosis for rare diseases**

258 It has been demonstrated for many diseases, both rare and common, that only a fraction of
259 affected individuals actually have a diagnosis annotated in their EHR³⁷⁻⁴³. As a result, it is likely that
260 many of the participants annotated as controls in our dataset are actually undiagnosed cases. In order to
261 evaluate RarePT's performance in these undiagnosed cases, we repeated the regression analyses of
262 mortality and estimated DALY restricting to participants labelled as controls, so that participants who had
263 the corresponding diagnosis in their EHR were excluded. The mortality analysis produced similar results
264 after excluding known diagnosed cases: 101 pcodes (67% of pcodes tested) had a significant
265 association ($p < 0.05$) with increased mortality, of which 93 (60%) remained significant after Bonferroni
266 correction for 155 pcodes ($p < 0.00032$). The median regression coefficient for the proportional hazard
267 regression on mortality was 0.86, corresponding to a hazard ratio of 2.4. The DALY analysis also
268 produced similar results: 114 pcodes (74% of pcodes tested) had a significant association ($p < 0.05$)
269 with increased estimated DALY, and 106 (68%) remained significant after Bonferroni correction for 155
270 pcodes ($p < 0.00032$). 131 pcodes (85%) had a significant association with increased estimated YLD

271 individually, 126 (81%) after Bonferroni correction; 110 phecodes (71%) had a significant association
272 with increased estimated YLL individually, 104 (67%) after Bonferroni correction. For the median
273 phecode, a positive prediction was associated with an increase in DALY of 1.5 years in controls. These
274 results demonstrate that RarePT predictions are associated with health outcomes even when a diagnosis is
275 not present in the EHR, suggesting that RarePT identifies clinically relevant features even in undiagnosed
276 individuals and may be identifying undiagnosed cases.

277 We next repeated the logistic regression analysis testing RarePT predictions against abnormal test
278 results. As expected, excluding known cases reduced the significance of many, but not all, of these
279 regressions. Nevertheless, 47 of these regressions, representing 63% of these relationships, had a result
280 that was in the expected direction and statistically significant ($p < 0.05$), and 36 (48%) remained
281 significant after Bonferroni correction for 75 regressions ($p < 0.00067$). The median regression coefficient
282 was 0.45, corresponding to an OR of 1.57. As with the regressions that included cases, in 100 random
283 permutations of phecode-laboratory test relationships, no permutation showed as many Bonferroni-
284 significant associations ($p < 0.01$) (**Figure 3b, Supplementary Tables S10-S11**). We also repeated the
285 linear regression analysis testing RarePT predictions against quantitative test results, excluding both
286 known cases and participants with abnormal test results. 29 of these regressions, representing 39% of
287 these relationships, had a result that was in the expected direction and statistically significant, with 21
288 (28%) remaining significant after Bonferroni correction. Again, 0 out of 100 random permutations
289 showed as many Bonferroni-significant associations (**Figure 3b, Supplementary Tables S10-S11**). This
290 analysis supports the conclusion that RarePT's predictions are predictive not only of existing rare
291 diagnoses, but also of undiagnosed cases.

292 In order to estimate the number of these undiagnosed cases that exist in the UK Biobank dataset,
293 we first identified participants whose test results show that they are unlikely to be undiagnosed cases for a
294 particular phecode. We defined this category of "confirmed controls" as participants whose test results fell
295 within 1 standard deviation of the population mean for a particular test. This is possible because these

296 tests were administered in an unbiased way to a large cross-section of participants, and the presence of a
297 negative test result does not indicate that a physician ordered the test to rule out a suspected diagnosis. We
298 then measured RarePT's performance based on these confirmed controls and the observed diagnosed
299 cases. Assuming that RarePT performs similarly for unconfirmed controls and undiagnosed cases as for
300 confirmed controls and diagnosed cases, the prevalence of undiagnosed cases can be estimated by
301 comparing the expected number of false positives among unconfirmed controls to the actual number of
302 unconfirmed controls predicted as cases (**Figure 3c, Supplementary Note 1**).

303 We estimated the number of undiagnosed cases and the fraction of actual cases that are
304 undiagnosed for each rare phecode with an associated diagnostic test, using a bootstrap sampling
305 procedure to obtain 95% confidence intervals (**Figure 3d-e, Supplementary Table S12**). The estimated
306 proportion of undiagnosed cases varied widely by phecode, but nearly three-quarters of phecodes tested
307 ($23/32 = 72\%$) had an estimate greater than 20%. Even more remarkably, nearly two-thirds of phecodes
308 tested ($20/32 = 63\%$) had more undiagnosed cases than diagnosed cases, and over a third ($12/32 = 38\%$)
309 had a bootstrap confidence interval entirely above the number of diagnosed cases. The median estimated
310 rate of underdiagnosis across all phecodes tested was 83%, meaning that we estimate 83% of cases are
311 undiagnosed for the median rare phecode. This analysis suggests that there are a very large number of
312 undiagnosed cases of rare diseases in large population biobanks like the UK Biobank. Furthermore, it
313 suggests that RarePT is able to predict some of these hidden undiagnosed cases, allowing them to be
314 identified for the first time.

315 **Discussion**

316 Here we present RarePT, a transformer-based phenotype prediction method designed to predict
317 rare disease diagnoses based on diagnosis codes present in a patient's electronic health records (EHR). We
318 apply this method to predicting rare disease in the UK Biobank, and find that a very large fraction of rare
319 disease cases are undiagnosed. Our method adds to a growing collection of phenotype prediction methods
320 that use machine learning to clean and extend EHR data for downstream analysis⁴⁴⁻⁴⁸. Our method is

321 distinct from other approaches in that it focuses specifically on rare disease. It is typically difficult to train
322 machine learning approaches for rare disease because the low prevalence of these diseases limits the
323 availability of training data. We overcame this difficulty using a “masked diagnosis modeling” approach
324 inspired by the approaches used to train AI language models such as BERT¹⁹. This approach learns about
325 diagnoses by identifying which other diagnoses are most likely to appear in similar contexts, allowing it
326 to learn informative features even for rare diagnoses. In addition to this training strategy, we reweighted
327 our training data to give equal importance to rare and common diagnoses, boosting our power to predict
328 rare diseases.

329 The trained RarePT model is highly predictive of a wide range of rare disease diagnoses, showing
330 the promise of our deep learning approach as a screening test for specific rare diagnoses that could be
331 applied in a clinical setting in the future. Across all rare phecodes, RarePT’s predictions are associated
332 with a median diagnostic odds ratio (OR) of 48.0 in cross-validation – that is, participants predicted to
333 have a rare diagnosis by our model are 48.0 times more likely to have that diagnosis in their EHR
334 compared to participants without such a prediction. For some specific rare diseases, this performance is
335 even better, with the top 10% of diagnoses achieving a diagnostic OR over 350 in cross-validation and the
336 top 5% achieving a diagnostic OR over 1,500 in cross-validation. These values compare favorably to
337 many diagnostic tests currently in standard clinical use. Remarkably, this performance is replicated in a
338 completely independent cohort of patients at the Mount Sinai Health System in New York, which
339 represents not only a different health system but an entirely different country with different medical
340 practices and different standards for diagnostic coding and billing. The ability to predict rare disease
341 diagnoses in this independent cohort shows the power and transferability of this approach.

342 In addition to successfully predicting rare diagnoses in participants’ EHR, RarePT also provides
343 new evidence that a substantial number of participants may suffer from rare diseases without a diagnosis
344 appearing in the EHR. This reinforces the known fact that many diagnoses are missing from EHR, due to
345 biases in diagnosis, inconsistent use of billing codes, incomplete or fragmented patient records, and other

346 issues⁴⁹⁻⁵³. This effect has previously been quantified for a variety of diseases, both common and rare, and
347 often found to be substantial. For example, biobank studies have estimated that up to 75% of patients with
348 erythropoetic protoporphyria (EPP)³⁹, approximately 85% of patients with familial
349 hypercholesterolemia⁵⁴, and approximately 90% of patients with glycated hemoglobin (HbA_{1c}) levels
350 indicating diabetes⁴² remain undiagnosed. This is especially problematic for rare diseases, due to the
351 known difficulty of correctly diagnosing rare diseases and the long diagnostic odyssey experienced by
352 many rare disease patients^{3-6,55}. We hypothesized that RarePT would correctly predict many of these
353 undiagnosed cases, causing them to appear as false positives despite actually being correct predictions.
354 For rare diseases where relevant biomarkers were available, these biomarkers consistently showed a
355 significant excess of abnormal values and more extreme values within the normal range in individuals
356 predicted positive by RarePT, supporting the hypothesis that many of RarePT's predictions are actually
357 undiagnosed cases. Consistent with literature estimates for common diseases, we estimate the prevalence
358 of undiagnosed cases in rare diseases to be remarkably high: 72% of phecodes we tested appeared to have
359 an underdiagnosis rate above 20%, and 63% of phecodes we tested were consistent with a majority of
360 cases being undiagnosed. While these numbers may be higher than the true rate in the general population
361 due to the UK Biobank being biased towards healthier participants, who are less likely to seek out and
362 receive diagnoses than the general population^{26,27,31}, both the existence and magnitude of this
363 phenomenon are consistent with previous results on underdiagnosis of diseases in EHR. The RarePT
364 model allows us to measure this underdiagnosis systematically across a range of rare diseases, which has
365 not previously been possible, as well as to identify specific individuals who may be suffering from rare
366 disease and have a missing or incorrect diagnosis.

367 There are many potential practical applications for RarePT. One of these is as a phenotype
368 imputation step in a data preprocessing pipeline for downstream bioinformatic analysis, which has
369 previously been unavailable for rare diagnoses^{31,56,57}. Another application is in collecting rare disease
370 cohorts for research studies or drug trials. Due to the rarity of rare diseases, identifying multiple patients

371 with the same disease is a pressing problem in rare disease research, and the international research
372 community has developed several tools to address it⁷⁻⁹. RarePT can augment or support these tools by
373 allowing researchers to rapidly search EHR for patients who are likely to have a particular disease or
374 patients who are phenotypically similar to another specific patient. It also has the potential to be
375 developed into a clinical screening test for rare diseases, particularly in patients with a specific risk factor
376 such as family history of disease or a genetic risk allele.

377 There are several limitations and areas of further development for this approach. First, phecodes
378 are designed for phenome-wide association (PheWAS) studies primarily targeting common phenotypes
379 and are not specifically designed to target rare diseases. While some specific rare phenotypes may be
380 inaccessible to RarePT for this reason, recent studies have shown that vocabularies for common disease,
381 including phecodes, do capture information about rare phenotypes⁵⁸⁻⁶⁰, and we identified phecodes that
382 were rare in both the UK Biobank and MSDW cohorts. Future versions of this approach could increase
383 the resolution for rare phenotypes by using phenotype ontologies designed to represent rare diseases, such
384 as the Human Phenotype Ontology (HPO)⁶¹ or the OrphaNet disease ontology⁶². However, there are
385 tradeoffs involved in this choice. Using a more fine-grained vocabulary for rare disease would
386 dramatically increase the complexity of the model and its computational requirements. The analyses we
387 present here require generation of millions of phenotype predictions, which took several hours of GPU
388 time under the current RarePT architecture and would quickly become infeasible if the phenotype
389 encoding used a complex hierarchical structure with a vocabulary many times larger. Choosing a
390 phenotype encoding scheme that can more precisely capture rare phenotypes could also harm the model's
391 ability to capture information about rare and common disease in the same vocabulary, which may reduce
392 the power of the model for rare disease and its transferability to other cohorts beyond its training set.

393 Second, the ICD-10 diagnosis codes we use to derive phecodes are known to be noisy and
394 unreliable^{52,63,64}. We have relied on established methods for automated phenome-wide phenotyping, many
395 of which use diagnosis codes in spite of their limitations because more reliable sources of data are either

396 difficult to access in an automated way or are not available phenome-wide^{31,53}. This is particularly true for
397 rare diseases, due to the difficulty of finding specialized experts who are capable of reviewing individual
398 patient charts in detail to arrive at a confident diagnosis, particularly at scale⁶⁵⁻⁶⁷. In spite of this, the use
399 of these automated phenotyping approaches could be a concern in our analysis of undiagnosed cases,
400 since both our identification of participants with a disease diagnosis and our identification of patients with
401 abnormal test results are based on these potentially unreliable automated procedures. It is possible that
402 some of the supposedly undiagnosed cases we identified were actually diagnosed cases where the
403 diagnosis escaped detection by our automated phenotyping process. It is also possible that the availability
404 of certain test results and not others biased our analysis towards specific categories of phenotypes that are
405 not representative of rare diseases in general. For example, blood disorders and metabolic disorders
406 appear to be overrepresented among phenotypes with available diagnostic tests, while neoplasms and
407 neurological disorders are entirely absent (**Supplementary Table S6**). However, while diagnosis rates
408 may differ by category, there is no reason to suppose that categories with greater availability of diagnostic
409 tests are diagnosed at a lower rate than those with less availability of diagnostic tests. Indeed, if anything,
410 the availability of simple diagnostic tests should increase the rate of diagnosis, making our estimates
411 conservative. It is likely that we could make more reliable determinations of both diagnosis status and
412 true phenotype by making use of other features available in the EHR, such as lab results, vitals,
413 medications, or unstructured physician's notes. We chose not to include these features in this analysis to
414 avoid circularity in training and analysis, as the diagnoses contained in the EHR are informed by the lab
415 results, vitals, and physician's notes from the same EHR. Without careful insulation of these different
416 modalities of data, any trained model or statistical analysis is likely to simply recapitulate the physician's
417 diagnostic criteria without gaining any predictive power for undiagnosed patients, a problem which
418 RarePT avoids by excluding these redundant data sources. Previous studies have also identified
419 undiagnosed cases using a longitudinal study design with direct physician involvement^{40,41}. RarePT could
420 facilitate this kind of analysis for specific diagnoses in future studies, following up on this broad
421 automated analysis with in-depth analysis of individual diseases incorporating specific clinical expertise.

422 Finally, there are many opportunities to improve on our model architecture. Deep learning and AI
423 is a rapidly evolving field, and the transformer-based architecture we used for this analysis may not be the
424 optimal way to learn the semantic structure of EHR diagnoses. Recent studies have proposed new ways of
425 derived phenotype embeddings, including extracting them from curated knowledge graphs or from
426 general-purpose large language models pretrained on non-EHR data⁶⁸⁻⁷⁰. There are also a variety of
427 approaches that have been used to process time series data from EHR in machine learning applications,
428 including neural network models designed for time series data such as recurrent neural networks and
429 using pretrained large language models to process EHR^{71,72}. While incorporating newer and more
430 sophisticated approaches may improve the model, they may also promote overfitting and reduce
431 transferability of the model across health systems and datasets, as well as slowing down training and
432 prediction. RarePT is both transferable and tractable, essential properties for a method designed to process
433 large health system datasets.

434 In this paper we have shown that our deep learning phenotype prediction approach, RarePT, is
435 capable of modeling and predicting rare disease diagnoses on a phenome-wide basis, with performance
436 that compares favorably to diagnostic screening tests used in clinical settings. Remarkably, RarePT
437 achieves this performance not only in held-out segments of the UK Biobank cohort it was trained on, but
438 also on an entirely separate cohort of patients in the Mount Sinai Health System in New York City. This
439 demonstrates that the predictive features RarePT uses are not specific to the UK Biobank, but are robust
440 to differences in recruitment strategy, differences in race and ethnicity, and even differences in medical
441 practices and billing procedures between countries. In addition to capturing specific diagnoses, RarePT
442 predictions also predict clinical outcomes, including mortality, quality of life, and specific biomarkers
443 associated with rare disease. Finally, we used predicted phenotypes from the model to estimate the
444 prevalence of undiagnosed rare disease in the UK Biobank, showing that it is likely extremely high. This
445 kind of systematic phenome-wide analysis has not previously been possible for rare diseases, highlighting
446 the utility of RarePT to conduct large-scale studies on rare disease. The high rate of undiagnosed rare

447 disease in large population datasets like the UK Biobank also highlights the need for new methods like
448 RarePT to address the problem of undiagnosed rare disease and suggests a wide range of valuable clinical
449 and research applications.

450

451 **Author Contributions:** Dr. Jordan and Dr. Do had full access to all of the data in the study and take
452 responsibility for the integrity of the data and accuracy of the data analysis.

453 *Concept and design:* Jordan, Do.

454 *Acquisition, analysis, or interpretation of the data:* Jordan, Vy, Do.

455 *Drafting of the manuscript:* Jordan, Do.

456 *Critical revision of the manuscript for important intellectual concept:* Jordan, Do.

457 *Statistical analysis:* Jordan.

458 *Administrative, technical, or material support:* Do

459 *Supervision:* Do.

460

461 **Conflict of Interest Disclosures:** Dr. Do reported receiving grants from AstraZeneca, grants and non-
462 financial support from Goldfinch Bio, being a scientific co-founder, consultant and equity holder for
463 Pensieve Health (pending), and being a consultant for Variant Bio, all not related to this study.

464

465 **Materials and Correspondence:** Requests for materials and correspondence should be addressed to Dr.
466 Do.

467

468 **Data availability statement:** Summary data required to generate figures will be deposited in a public
469 repository prior to publication, and are available on request from the authors otherwise. Individual-level
470 data from the UK Biobank and the Mount Sinai Data Warehouse are governed by third-party data use
471 agreements and cannot be made available with this study. Researchers who qualify for access to
472 deidentified data under the policies of the UK Biobank and/or the Mount Sinai Health System can access
473 these data upon application to the respective institutions.

474

475 **Code availability statement:** Code required to train all models, run all analyses, and generate all figures
476 will be published in a public repository under an open-source license prior to publication, and is available
477 upon request from the authors otherwise.

478

479 **Funding/Support:** Dr. Do is supported by the National Institute of General Medical Sciences of the NIH
480 (R35-GM124836). This research has been conducted using the UK Biobank Resource under Application
481 Number 16218. This work was supported in part through the Mount Sinai Data Warehouse (MSDW)
482 resources and staff expertise provided by Scientific Computing and Data at the Icahn School of Medicine
483 at Mount Sinai.

484

485 **Disclaimer:** The content is solely the responsibility of the authors and does not necessarily represent the
486 official views of the National Institutes of Health.

487

488 **Methods**

489 *Data collection and preprocessing*

490 The primary training data were derived from the UK Biobank⁷³. For each participant, we
491 retrieved age at recruitment (field 21022), sex (field 31), and a list of all ICD-10 diagnosis codes recorded
492 across all inpatient hospital records (field 41270). We also retrieved self-reported ethnicity (field 21000),
493 body mass index (BMI, field 21001), blood pressure (fields 4079-4080), LDL cholesterol (field 30780),
494 total cholesterol (field 30690), blood glucose (field 30740), and glycated hemoglobin (HbA_{1C}, field
495 30750) as baseline cohort characteristics (shown in **Table 1**), and 23 diagnostic tests and biomarkers
496 assayed as part of the recruitment process (**Supplementary Table S6**), though these were not included in
497 the data used to train the model. For this study, we represented each ICD-10 code as a binary indicator
498 that could either be present or absent, ignoring the dates associated with each code. We mapped ICD-10
499 codes to version 1.2 of the Phecode classification, using the published mapping²⁴. This resulted in a
500 dataset of 436,407 participants, 239,711 female and 196,696 male, including 1,558 of the 1,570 phecodes
501 with defined ICD-10 code mappings. Demographic and clinical characteristics of this cohort are shown in
502 **Table 1**.

503 We constructed a balanced dataset of training examples consisting of 100 cases and 100 controls
504 of each phecode, excluding phecodes with fewer than 100 cases or fewer than 100 controls. This excluded
505 273 phecodes, leaving 1,297 unique query phecodes. These 273 phecodes remained in the training data as
506 diagnoses but were never used as the query in any training examples. Cases were defined as participants
507 whose phecode diagnoses contained the query phecode. Controls were defined as participants whose
508 phecode diagnoses did not include the query phecode or any phecodes listed as exclusions for the query
509 phecode. For sex-specific phecodes, controls were also required to be the correct sex for the query
510 phecode. For cases, the query phecode and all exclusion phecodes were removed from the diagnosis list
511 as part of preprocessing (**Figure 1a**). Phecode diagnoses were encoded using a many-hot encoding, with

512 phecodes ordered linearly by their numerical code; query phecodes were encoded using a one-hot
513 encoding with the same ordering of phecodes.

514 Training examples were randomly split into five equal subsets for five-fold cross validation.
515 Since the same individual can appear multiple times with different query phecodes, we required each
516 cross-validation subset to contain a unique set of individuals, so that each individual can only appear in
517 one subset. This prevents the model from improving performance by recognizing specific individuals
518 from the training set and recapitulating the known diagnoses of those individuals. Cross-validation
519 subsets were also required to contain similar numbers of cases and controls.

520 Our independent validation data were derived from the Mount Sinai Data Warehouse (MSDW), a
521 database of clinical and operational data derived from the electronic health records (EHR) systems of the
522 Mount Sinai Health System in New York City. These data are anonymized, standardized, and
523 preprocessed for use in clinical and translational research. For each patient in this database, we retrieved
524 age as of 2009 (the median date associated with the “age of recruitment” field in the UK Biobank),
525 physician-reported sex, and a list of all ICD-10-CM diagnosis codes recorded in the EHR. We mapped
526 these to phecodes using the published mapping for ICD-10-CM diagnosis codes, which is slightly
527 different from the mapping for the ICD-10 codes used in the UK Biobank²⁴. Patient records were
528 processed into examples suitable for the model in the same way as described above. The final size of this
529 cohort was 3,333,560. Demographic and clinical characteristics of this cohort are shown in **Table 1**.

530 Study protocols were approved by the Institutional Review Board at the Icahn School of
531 Medicine at Mount Sinai (New York City, NY, USA; GCO#07–0529; STUDY-11–01139) and all
532 participants provided informed consent. Use of data from the UK Biobank was approved with the UK
533 Biobank Resource under application number 16218.

534 *Model architecture, tuning, and training*

535 The model was implemented in Python using the Keras package⁷⁴. **Figure 1a** shows a schematic
536 of the model architecture. The input diagnosed phecodes feed into a stack of modified transformer
537 decoder modules, based on the TransformerDecoder layer implemented in the KerasNLP package⁷⁵. The
538 standard TransformerDecoder layer was modified to remove the causal attention mask that prevents the
539 self-attention layer from paying attention to positions that are later in the sequence than the token
540 currently being considered. Since the phecode encoding is ordered by phenotype category and we are
541 ignoring temporal sequencing, this causal mask would be inappropriate. Each decoder layer also contains
542 a cross-attention layer which takes input from the encoded query phecode, allowing the model to learn
543 attention relationships between the diagnosis phecodes and the query phecode. To adjust for demographic
544 variables, the demographic variables are passed through a single densely connected layer to transform
545 them into the same dimension as the phecodes, and then added to the output of the transformer layers and
546 normalized. Finally, the prediction is given by a dot product between the input query phecodes and the
547 demographics-adjusted output of the transformer layer, and transformed into a probability score using a
548 softmax function. The Python code implementing this architecture will be made publicly available along
549 with this publication.

550 Hyperparameter tuning was performed using the Hyperband algorithm, as implemented in the
551 KerasTune package^{25,76}. The list of hyperparameters and their final tuned values are found in
552 **Supplementary Table S13**. We randomly sampled 80% of the training data to use for training and used
553 the remaining 20% as the validation set for the hyperband algorithm, choosing the hyperparameters that
554 minimized the training loss function on the validation set. For five-fold cross-validation runs, this 20%
555 held out validation sample was contained within the training set and did not overlap the cross-validation
556 test set. The hyperband algorithm was run for up to 18 epochs per model, stopping if validation loss failed
557 to improve in 5 consecutive epochs. The final selected model was then trained for up to 54 epochs, and
558 the best epoch was selected based on validation loss. Finally, after tuning was complete, the held-out

559 validation set was added back into the training set and selected model was retrained for the selected
560 number of epochs on the complete training set. Again, the Python code implementing this training
561 procedure will be made publicly available along with this publication. We repeated this hyperparameter
562 tuning for each of the five training subsets produced by cross-validation as well as on the full training set,
563 and all six runs selected identical values for all hyperparameters. Models were trained using NVIDIA
564 A100 GPUs on the Mount Sinai local high-performance computing cluster, Minerva. Each cross-
565 validation run took approximately 5 hours of GPU time to tune and train, and the full model took
566 approximately 6 hours, for a total of approximately 31 hours.

567 *Prediction of rare phenotypes*

568 We calculated the prevalence of each phecode in the UK Biobank by dividing the number of
569 cases by the total number of participants. We identified 155 rare phecodes with prevalence less than 1 in
570 2,000, or 0.05%, corresponding to the European Union definition of a rare disease. We calculated model
571 predictions from each of the six trained models (five cross-validation models and one full-dataset model)
572 using each of these 155 phecodes as a query for 436,407 participants in the UK Biobank, excluding from
573 each model the participants contained in its own training set. We additionally produced predictions from
574 the full trained model for 3,333,560 patients in the MSDW cohort for each of the 151 rare phecodes that
575 were also present in that cohort. Generating model predictions for the UK Biobank cohort took
576 approximately 3 hours of GPU time for each of the six models, for a total of approximately 18 hours;
577 generating model predictions for the MSDW cohort took approximately 45 hours of GPU time. The
578 MSDW cohort took substantially longer because the dataset was too large for the model to fit into
579 memory and had to be broken up into batches.

580 We quantified the performance of our models by diagnostic odds ratios and positive predictive
581 values. We arbitrarily chose a threshold probability score of 0.95 to represent a relatively high-confidence
582 case prediction, and treated predictions with probability score > 0.95 as predicted cases and ≤ 0.95 as

583 predicted controls. We quantified performance using odds ratio (OR) and positive predictive value (PPV),
584 as these are measures relevant to diagnostic screening tests²⁹. We calculated OR as

$$585 \quad OR = \frac{(TP + 0.5)/(FP + 0.5)}{(FN + 0.5)/(TN + 0.5)}$$

586 (1)

587 where TP is the number of true positive predictions (cases correctly predicted as cases), FP is the number
588 of false positive predictions (controls incorrectly predicted as cases), TN is the number of true negative
589 predictions (controls correctly predicted as controls), and FN is the number of false negative cases. In
590 other words, this is the ratio between the odds of a positive prediction being a case and the odds of a
591 negative prediction being a case. We added a correction of 0.5 to each count to correct for zeros⁷⁷. We
592 calculated PPV as

$$593 \quad PPV = \frac{TP}{TP + FN}$$

594 (2)

595 In other words, this is the probability that a positive prediction is a case. In all instances, we excluded
596 controls with an exclusion phecode, controls whose sex did not match the phecode, and all individuals
597 who were included in the training set of the cross-validated models.

598 *Mortality and DALY analyses*

599 We estimated disability-adjusted life years (DALY) and its components years lost to disability
600 (YLD) and years of life lost (YLL) for UKBB individuals using per-disease estimates from the 2019
601 Global Burden of Disease (GBD) study. We used the 80 non-overlapping non-communicable diseases that
602 account for the majority of a population's DALY as described by Jukarainen et al.^{32,33} GBD definitions of
603 specific diseases and conditions were used to label individuals affected by these diseases in the UK
604 Biobank. Estimates of disease burden in the UK from GBD were then applied to individuals with each

605 disease to produce estimated values of DALY, YLD, and YLL.³² These estimated values were tested
606 against RarePT predictions by linear regression. We retrieved a single prediction score by using the model
607 trained on the full dataset for individuals who were not included in the training set, and the appropriate
608 cross-validation model for individuals who were included in the training set (that is, the cross-validation
609 model whose training set did not include that individual). We turned this score into a binary prediction
610 using an arbitrary threshold of 0.95. We then performed linear regression testing the ability of this score
611 (independent variable) to predict DALY, YLD, or YLL (dependent variable), controlling for age, sex, and
612 self-reported ethnicity. We repeated this analysis both including all UK Biobank participants and
613 excluding known diagnosed cases and exclusions for each phecode.

614 We additionally retrieved date of death and date of recruitment from the UK Biobank (fields
615 40000 and 53) and performed an analysis of mortality using Cox proportional hazard regression. We
616 treated time from recruitment to death as a right-censored dependent variable, again using the binary
617 RarePT prediction as an independent variable along with age at recruitment, sex, and self-reported
618 ethnicity. As with DALY, we repeated this analysis both including all UK Biobank participants and
619 excluding known diagnosed cases and exclusions for each phecode. All regressions were performed in
620 Python using the statsmodels package⁷⁸.

621 *Biomarker and diagnostic test analysis*

622 We collected biomarkers and diagnostic tests associated with phecodes using the SNOMED-CT
623 database of clinical terms^{34,35}. We identified all ICD-10 codes that mapped to any of our 155 rare
624 phecodes and also mapped to a SNOMED-CT term with an “interprets” relationship (concept
625 363714003). The “interprets” relationship indicates that the concept represented by the diagnosis code has
626 an underlying evaluation that is “intrinsic to the meaning of” that concept⁷⁹. Examples of this kind of
627 relationship include the relationship between obesity and measured body weight, hypercholesterolemia
628 and total serum cholesterol, or thrombocytopenia and platelet count. In most cases, SNOMED-CT also
629 identifies the direction of the relationship using the “has interpretation” relationship (concept 363713009).

630 For example, hypercholesterolemia is interpreted as total serum cholesterol above reference range, while
631 thrombocytopenia is interpreted as platelet count below reference range. For each concept that was the
632 target of an “interprets” relationship, we manually searched for a corresponding measurement available in
633 the UK Biobank and a corresponding reference range. The result was a list of 75 relationships between
634 rare phecodes and UK Biobank data fields, encompassing 32 rare phecodes and 23 data fields, each with
635 an expected direction of relationship (above, below, or outside) and sex-specific reference ranges
636 **(Supplementary Table S6).**

637 As with the previously described regression analyses, we retrieved a single prediction score for
638 each of these 32 rare phecodes by using the model trained on the full dataset for individuals who were not
639 included in the training set, and the appropriate cross-validation model for individuals who were included
640 in the training set (that is, the cross-validation model whose training set did not include that individual).
641 We turned this score into a binary prediction using the same 0.95 threshold. We then performed two
642 regression analyses testing the ability of this binary prediction (independent variable) to predict the
643 corresponding data field (dependent variable), controlling for age, sex, and self-reported ethnicity. In the
644 first analysis, we used the expected direction of relationship and the reference range to construct a binary
645 variable indicating whether each individual had an abnormal result in the direction expected. We
646 performed logistic regression using this binary variable as the dependent variable. We repeated this
647 analysis both including all participants and excluding individuals labelled as cases or exclusions for each
648 phecode. This regression tests whether the model can predict individuals with abnormal test results
649 consistent with a diagnosis even in individuals labelled as controls. In the second analysis, we normalized
650 the values of each biomarker within the reference range so that the sample population for each sex had
651 mean 0 and variance 1 after excluding all individuals with values outside the reference range. Finally, we
652 aligned the values so that the expected direction of association was always positive, by multiplying them
653 by -1 for “below reference range” relationships and taking their absolute value for “outside reference
654 range” relationships. We performed standard linear regression using this normalized and aligned value as

655 the dependent variable. We repeated this analysis both including all participants and excluding both
656 individuals labelled as cases or exclusions and individuals with abnormal test results. This regression tests
657 whether the model can predict individuals with elevated or reduced results even if they are still within the
658 normal range. Regressions were performed in Python using the statsmodels package⁷⁸.

659 Finally, we identified a set of “confirmed controls” for each phecode, defined as individuals labelled as
660 controls who also had all associated results within the reference range and within one standard deviation
661 of the population mean for their sex. We consider these individuals very unlikely to be undiagnosed cases
662 incorrectly labelled as controls. We used the performance of our model on these confirmed controls to
663 estimate the false positive rate of our model for each of the 32 phecodes with available relationships to
664 test results. We then used this false positive rate to estimate the number of undiagnosed cases using the
665 following relationship:

$$666 \quad P_u = \frac{UP - FPR \times U}{TPR - FPR}$$

667 Where P_u represents the number of undiagnosed cases; U represents the total number of individuals with
668 unknown case-control status, excluding controls confirmed by laboratory tests but including undiagnosed
669 cases; UP represents the total number of unknowns predicted as cases by the model, again excluding
670 controls confirmed by laboratory tests but including undiagnosed cases; FPR represents the false positive
671 rate of the model as estimated from confirmed controls; and TPR represents the true positive rate of the
672 model as estimated from diagnosed cases. See **Supplementary Note 1** for derivation and discussion of
673 this relationship.

674 The Python code implementing all these analyses will be made publicly available along with this
675 publication.

676 *Software Package and Workflow*

677 For portability and reproducibility, model training and analysis code is formatted as a Snakemake
678 workflow.⁸⁰ This allows easy retraining of the RarePT model and reproduction of the analyses reported
679 here on any appropriately-formatted individual-level dataset. After creating an appropriately named and
680 formatted input data file and setting up Snakemake for their execution environment, users can train a new
681 model with a single command:

```
682 snakemake  
683 results/models/my_dataset.100_case_control_sample.full_from_5_fold_cv.  
684 seed_18
```

685 The number of cases and controls sampled, the number of cross-validation folds used for testing, and the
686 random seed can be changed by changing the appropriate values in the targeted filename. Likewise,

```
687 snakemake  
688 results/data/my_dataset.100_case_control_sample.5_fold_cv.seed_18.all_  
689 rare_predictions_with_cv.parquet
```

690 trains a model and uses cross-validation to generate model predictions for all rare phecodes, and

```
691 snakemake  
692 results/data/uk_biobank.100_case_control_sample.full_from_5_fold_cv.se  
693 ed_18.vs.my_dataset.all_rare_predictions.parquet
```

694 uses the UK Biobank trained model to generate model predictions for all rare phecodes in a user dataset.

695 Snakemake can be configured for many different high-performance computing and cloud computing
696 environments and, when properly configured, automatically manages resource requirements and package
697 dependencies.

698 The Snakemake workflow will be published with acceptance of this manuscript in a peer-
699 reviewed journal. Prior to formal publication, it is available on request from the authors.

700

701

702

703 **Table 1. Baseline demographic and clinical characteristics of cohorts.**

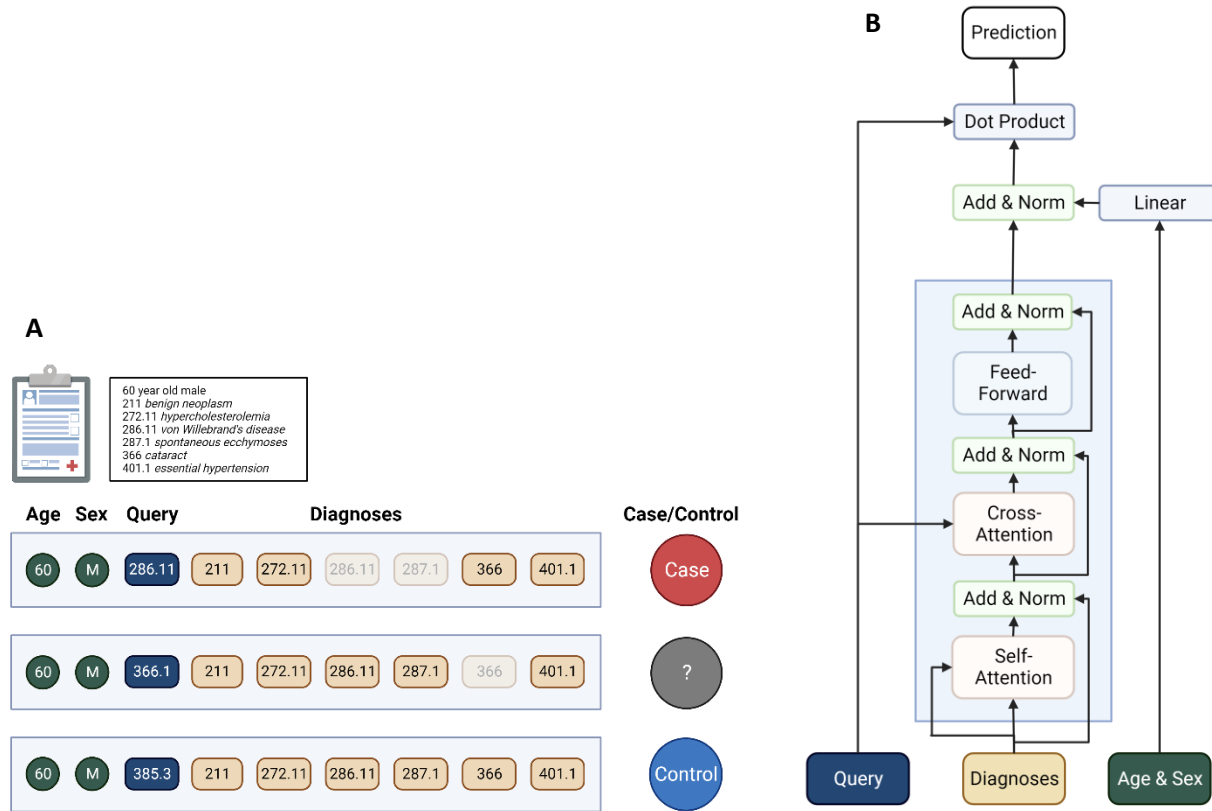
Characteristic		No. (%)	
		UK Biobank (N=436,407)	Mount Sinai Data Warehouse (N=3,333,560)
Women		239,611 (55%)	1,861,995 (56%)
Men		196,696 (45%)	1,471,565 (44%)
Age, mean (SD), yrs		57 (8)	33 (22)
Self- Reported Race and/or Ethnicity	<i>White</i>	411,074 (94%)	1,118,704 (34%)
	<i>Black</i>	6,900 (1.6%)	385,815 (12%)
	<i>Asian</i>	9,695 (2.2%)	185,233 (5.6%)
	<i>Hispanic or Latino</i>	0 (0%)	367,924 (11%)
	<i>Multiple, Other, or Unknown</i>	8,738 (2.0%)	383,984 (12%)
BMI, mean (SD), kg/m²		27 (5)	26 (7)
SBP, mean (SD), mmHg		138 (19)	123 (17)
DBP, mean (SD), mmHg		82 (11)	74 (10)
LDL-C, mean (SD), mg/dL		138 (34)	101 (35)
Total cholesterol, mean (SD), mg/dL		220 (45)	181 (45)
HbA_{1c}, median (IQR), %		5.5 (0.4)	5.6 (1.6)
Glucose, mean (SD), mg/dL		92 (22)	132 (55)
Hypertension (phecode 401)		151,254 (34%)	497,655 (15%)
Type 2 Diabetes (phecode 250.2)		42,335 (9.7%)	211,779 (6.4%)
Breast Cancer (phecode 174)		436,407 (5.2%)	43,212 (1.3%)

704 BMI = body mass index, SBP = systolic blood pressure, DBP = diastolic blood pressure, LDL-C = LDL
705 cholesterol, HbA_{1c} = Hemoglobin A_{1c}. For the UK Biobank cohort, demographic characteristics, vitals,
706 and blood tests were measured at recruitment. In the Mount Sinai Data Warehouse cohort, there was no
707 distinct recruitment visit, so individual ages were measured as age in 2009, which corresponds to the
708 median recruitment date for the UK Biobank cohort, and individual vitals and blood tests were measured
709 as the mean across all available measurements from each individual's electronic health records (EHR). In
710 both cohorts, diagnoses were measured using phecodes derived from the presence or absence of specific
711 ICD-10 or ICD-10-CM diagnosis codes across each individual's EHR; see methods for details.

712

713

714 **Figure 1. Schematics of masked phenotype modeling training procedure and RarePT model**
715 **architecture.**

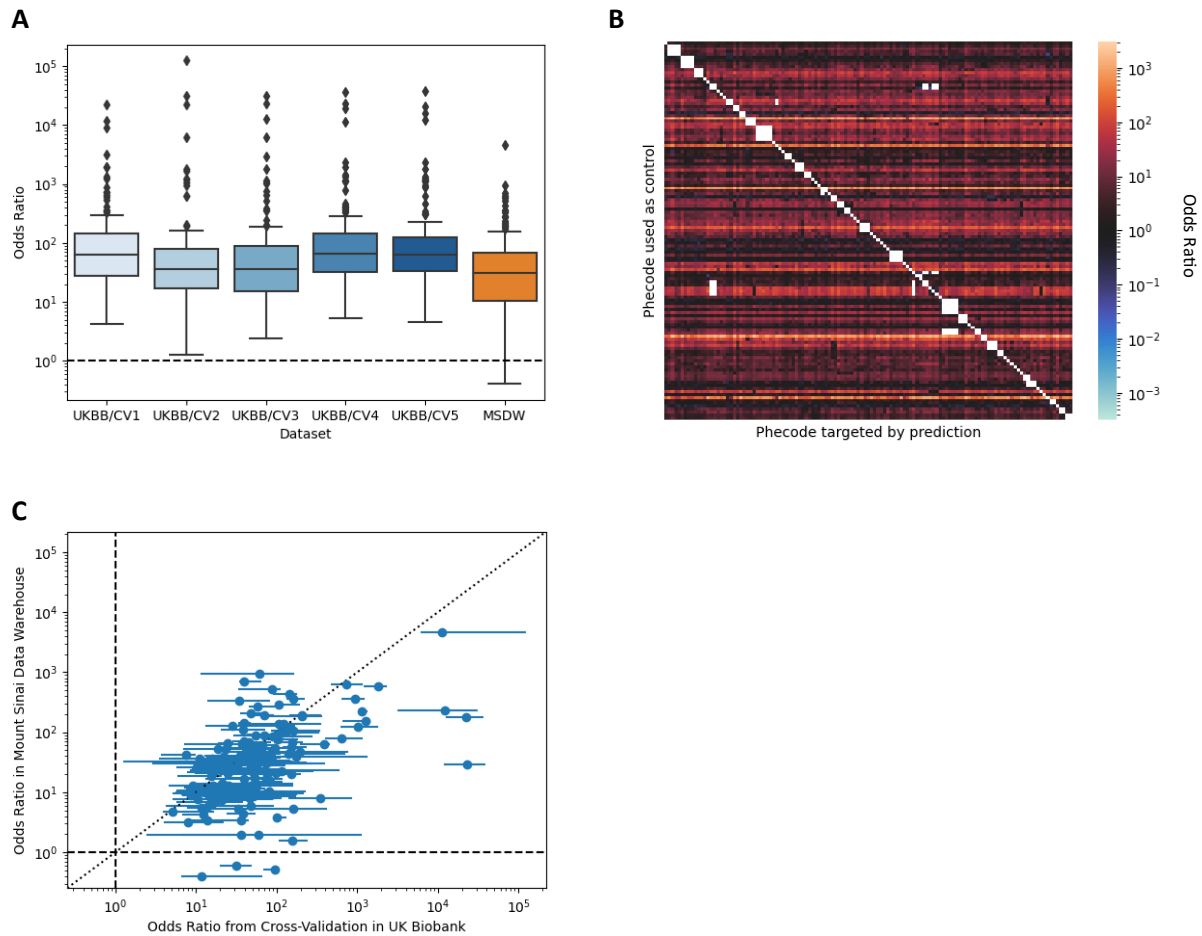


716 (a) Masked phenotype modeling training procedure. Each individual is represented by demographic
717 characteristics (age and sex) and the set of all diagnoses present in their EHR across all encounters,
718 represented as phecodes. A single training example consists of all demographic characteristics, a “query”
719 phecode indicating the phecode to be trained on, and the target case/control designation. The example is
720 considered a “Case” if the individual’s EHR contains the query phecode, “Unknown” if the individual’s
721 EHR does not contain the query phecode but does contain a phecode defined as an exclusion for the query
722 phecode, and “Control” if the individual’s EHR contains neither the query phecode nor any phecode
723 defined as an exclusion for the query phecode. In all instances, the query phecode and all phecodes
724 defined as exclusions for the query phecode are hidden before training or prediction. (b) Rare-Phenotype
725 Prediction Transformer (RarePT) architecture. We apply a simplified transformer architecture to the

726 diagnoses. Diagnoses are fed into a transformer decoder in a many-hot encoding, in a single step with no
727 sequence. The transformer decoder consists of a self-attention layer, a cross-attention attending to the
728 query phecocode, and a feed-forward layer, all connected by residual skip connections, as in the standard
729 transformer architecture. Demographic characteristics are interpreted by a dense linear layer and then
730 applied as an adjustment to the transformer output. The final prediction is the dot product of the query
731 phecocode input with the demographic-adjusted transformer output.

732

733 **Figure 2. Performance of RarePT to predict diagnosed cases of 155 rare phecodes.**



734

735

736 (a) Box and whisker plots showing distribution of diagnostic odds ratio (ratio of odds in predicted cases

737 to odds in predicted controls) across rare phecodes. The five blue boxes labelled “UKBB/CV1-5” show

738 the performance of the five models trained using five-fold cross-validation within the UK Biobank,

739 excluding each model’s training set; median OR across all five cross-validation sets was 48 (full range

740 1.27-39,000). The orange box labeled “MSDW” shows the performance of the model trained using the

741 full UK Biobank dataset on the Mount Sinai Data Warehouse dataset, an independent cohort from the

742 Mount Sinai Health System in New York; median OR for this cohort was 31 (full range 0.41-4,600).

743 Performance on the independent dataset is only slightly reduced, despite extensive differences between

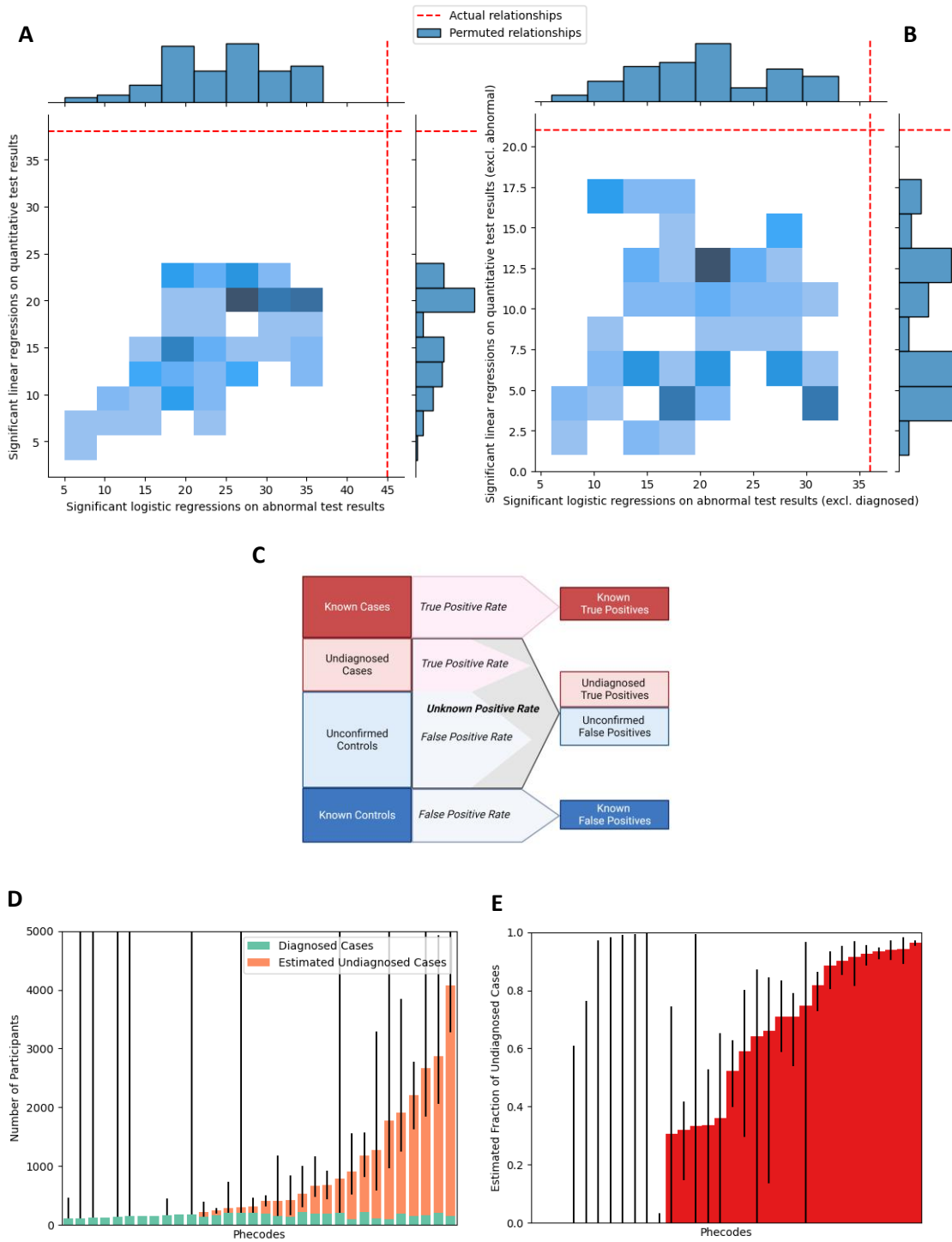
744 these datasets and health systems, demonstrating the robustness of the approach. Boxes show 1st quartile,
745 median, and 3rd quartile; whiskers extend to 80% of interquartile range; dashed line shows random
746 chance. (b) Heatmap showing odds ratios (median across five-fold cross-validation) for distinguishing
747 between diagnosed cases of one phecode and diagnosed cases of a second phecode. Columns represent
748 different case phecodes, rows represent different control phecodes, both ordered so that phecodes in the
749 same biological category are grouped together. For each cell, the RarePT models trained in cross-
750 validation were presented with a dataset consisting of all diagnosed cases of both phecodes in UK
751 Biobank, excluding the model's own training set. The model was asked to predict case status for the case
752 phecode, treating cases of the control phecode as controls. Red indicates that the model's case predictions
753 are enriched for cases of the correct phecode, while blue indicates that the model's case predictions are
754 depleted for cases of the correct phecode. Blank (white) cells indicate phecodes that are labelled as
755 exclusions in the definition of the target phecode. While some comparisons are more successful than
756 others, mostly depending on the identity of the phecode used as controls, RarePT successfully predicts
757 cases across a wide range of rare phecodes, even when the background is cases for other rare phecodes.
758 Blue This demonstrates that RarePT's predictions are specific to each specific phecode, and are not
759 primarily predicting categories of diagnoses or general health. (c) Scatterplot comparing RarePT's
760 performance for rare phecodes on UK Biobank and Mount Sinai Data Warehouse cohorts. Each data
761 point represents a single phecode; horizontal error bars represent the results of five-fold cross validation
762 within the UK Biobank cohort. Performance for specific phecodes is strongly correlated between cohorts
763 (Pearson $r = 0.456$, $p = 5.03 \times 10^{-40}$, t-test), demonstrating that the specific features used to differentiate
764 different phecodes are robust to differences between cohorts, including differences in diagnosis coding,
765 recruitment, and ethnic and racial composition. Dashed lines show random chance within each cohort;
766 dotted line shows equal performance between cohorts.

767

768

769

770 **Figure 3. Estimated performance of RarePT to predict undiagnosed cases of 32 rare phecodes with**
 771 **available relevant laboratory tests.**



773 *(a-b) Regression analysis showing RarePT predicts diagnostic test results in UK Biobank participants.*
774 Plots show number of significant regressions after Bonferroni correction for 75 tests in known relevant
775 diagnostic tests (red dashed lines) and 100 random permutations of test-disease relationships (blue boxes).
776 X axis shows logistic regression of abnormal test results vs. RarePT prediction, while Y axis shows linear
777 regression of quantitative test results vs. RarePT prediction, both controlling for age at recruitment, sex,
778 and self-reported ethnicity. Panel a shows results among all UK Biobank participants; Panel b excludes
779 known cases, and the Y axis of panel b additionally excludes participants with abnormal test results. For
780 both analyses, the number of significant regressions were lower than the actual observed results for all
781 100 permutations in both analyses. This demonstrates that RarePT predicts relevant clinical features of
782 disease, not only the presence of a diagnosis. Individuals with a positive prediction from RarePT but no
783 diagnosis likely have phenotypic profiles that are similar to the predicted rare diagnosis, and may include
784 undiagnosed cases. *(c) Illustration of estimation of undiagnosed cases.* All individuals can be classified as
785 known cases with EHR diagnosis, known controls with confirmed normal test results (defined as
786 participants within 1 standard deviation of the population mean for all tests associated with a given
787 phecode), or unknown. Unknown individuals are a mixture of undiagnosed cases and unconfirmed
788 controls in an unknown proportion. We assume that the behavior of the model on unknown individuals
789 ("unknown positive rate") is a mixture of its behavior on cases ("true positive rate") and controls ("false
790 positive rate"). With this assumption, we estimate the proportion of cases and controls in the unknown
791 group based on the relationship between these three positive prediction rates. See Supplementary Note 1
792 for discussion and derivation of this relationship. *(d-e) Number of undiagnosed cases and fraction of*
793 *cases undiagnosed for 32 rare phecodes, estimated from RarePT's performance on known cases and*
794 *controls.* Error bars represent bootstrap 95% confidence intervals. While the presence of undiagnosed
795 cases detected by RarePT varies across different rare phecodes, a majority of phecodes tested have an
796 estimate above zero, and half are estimated to have a higher number of undiagnosed cases than diagnosed
797 cases. This highlights the importance of developing methods to screen for undiagnosed cases of rare
798 disease.

800 **List of Supplementary Items**

801 **Supplementary Note S1.** Estimation of undiagnosed cases.

802 **Supplementary Figure S1.** Receiver operating characteristic (ROC) curves for training and testing cross-
803 validation datasets in UK Biobank.

804 **Supplementary Figure S2.** Positive predictive value of RarePT predictions for 155 rare phecodes in UK
805 Biobank and MSDW cohorts.

806 **Supplementary Table S1.** Cross-validation performance metrics for training and testing.

807 **Supplementary Table S2.** List of 155 rare phecodes in UK Biobank and Mount Sinai Data Warehouse.

808 **Supplementary Table S3.** Model performance by phecode for 155 rare phecodes in UK Biobank cohort.

809 **Supplementary Table S4.** Model performance by phecode for 151 rare phecodes in MSDW cohort.

810 **Supplementary Table S5.** Regression results for mortality, DALY, and related measures.

811 **Supplementary Table S6.** List of diagnostic tests relevant to rare phecodes.

812 **Supplementary Table S7.** Regression results for 75 diagnostic tests known to be relevant to rare
813 phecodes.

814 **Supplementary Table S8.** Regression results for 100 permutations of phecode-test relationships.

815 **Supplementary Table S9.** Regression results for mortality, DALY, and related measures in controls only.

816 **Supplementary Table S10.** Regression results for 75 diagnostic tests known to be relevant to rare
817 phecodes in controls only.

818 **Supplementary Table S11.** Regression results for 100 permutations of phecode-test relationships in
819 controls only.

820 **Supplementary Table S12.** Estimated numbers of undiagnosed cases for 32 rare phecodes.

821 **Supplementary Table S13.** Description of model hyperparameters and final tuned values.

822

823 **References**

- 824 1. Herder, M. What Is the Purpose of the Orphan Drug Act? *PLOS Med.* **14**, e1002191 (2017).
- 825 2. Richter, T. *et al.* Rare Disease Terminology and Definitions—A Systematic Global Review: Report
826 of the ISPOR Rare Disease Special Interest Group. *Value Health* **18**, 906–914 (2015).
- 827 3. Ferreira, C. R. The burden of rare diseases. *Am. J. Med. Genet. A.* **179**, 885–892 (2019).
- 828 4. Jaffe, A., Zurynski, Y., Beville, L. & Elliott, E. Call for a national plan for rare diseases. *J. Paediatr.*
829 *Child Health* **46**, 2–4 (2010).
- 830 5. Nutt, S. & Limb, L. Survey of patients' and families' experiences of rare diseases reinforces calls
831 for a rare disease strategy. *Soc. Care Neurodisability* **2**, 195–199 (2011).
- 832 6. Molster, C. *et al.* Survey of healthcare experiences of Australian adults living with rare diseases.
833 *Orphanet J. Rare Dis.* **11**, 30 (2016).
- 834 7. Azzariti, D. R. & Hamosh, A. Genomic Data Sharing for Novel Mendelian Disease Gene Discovery:
835 The Matchmaker Exchange. *Annu. Rev. Genomics Hum. Genet.* **21**, 305–326 (2020).
- 836 8. Philippakis, A. A. *et al.* The Matchmaker Exchange: A Platform for Rare Disease Gene Discovery.
837 *Hum. Mutat.* **36**, 915–921 (2015).
- 838 9. Boycott, K. M., Azzariti, D. R., Hamosh, A. & Rehm, H. L. Seven years since the launch of the
839 Matchmaker Exchange: The evolution of genomic matchmaking. *Hum. Mutat.* **43**, 659–667 (2022).
- 840 10. Laurie, S. *et al.* The RD-Connect Genome-Phenome Analysis Platform: Accelerating diagnosis,
841 research, and gene discovery for rare diseases. *Hum. Mutat.* **43**, 717–733 (2022).
- 842 11. Hsieh, T.-C. *et al.* GestaltMatcher: Overcoming the limits of rare disease matching using facial
843 phenotypic descriptors. 2020.12.28.20248193 Preprint at <https://doi.org/10.1101/2020.12.28.20248193>
844 (2021).
- 845 12. Frederiksen, S. D. *et al.* Rare disorders have many faces: in silico characterization of rare disorder
846 spectrum. *Orphanet J. Rare Dis.* **17**, 76 (2022).
- 847 13. Koroteev, M. V. BERT: A Review of Applications in Natural Language Processing and
848 Understanding. Preprint at <https://doi.org/10.48550/arXiv.2103.11943> (2021).
- 849 14. Zhang, M. & Li, J. A commentary of GPT-3 in MIT Technology Review 2021. *Fundam. Res.* **1**, 831–
850 833 (2021).
- 851 15. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–
852 589 (2021).
- 853 16. Cheng, J. *et al.* Accurate proteome-wide missense variant effect prediction with AlphaMissense.
854 *Science* **0**, eadg7492 (2023).
- 855 17. Ma, A. *et al.* Single-cell biological network inference using a heterogeneous graph transformer.
856 *Nat. Commun.* **14**, 964 (2023).

- 857 18. Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range
858 interactions. *Nat. Methods* **18**, 1196–1203 (2021).
- 859 19. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional
860 Transformers for Language Understanding. in *Proceedings of the 2019 Conference of the North American*
861 *Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long*
862 *and Short Papers)* 4171–4186 (Association for Computational Linguistics, 2019). doi:10.18653/v1/N19-
863 1423.
- 864 20. Yu, W. *et al.* Dict-BERT: Enhancing Language Model Pre-training with Dictionary. Preprint at
865 <https://doi.org/10.48550/arXiv.2110.06490> (2022).
- 866 21. Liu, Q., McCarthy, D. & Korhonen, A. Second-order contexts from lexical substitutes for few-shot
867 learning of word representations. in *Proceedings of the Eighth Joint Conference on Lexical and*
868 *Computational Semantics (*SEM 2019)* 61–67 (Association for Computational Linguistics, 2019).
869 doi:10.18653/v1/S19-1007.
- 870 22. Schick, T. & Schütze, H. BERTRAM: Improved Word Embeddings Have Big Impact on
871 Contextualized Model Performance. Preprint at <http://arxiv.org/abs/1910.07181> (2020).
- 872 23. Bastarache, L. Using Phecodes for Research with the Electronic Health Record: From PheWAS to
873 PheRS. *Annu. Rev. Biomed. Data Sci.* **4**, 1–19 (2021).
- 874 24. Wu, P. *et al.* Mapping ICD-10 and ICD-10-CM Codes to Phecodes: Workflow Development and
875 Initial Evaluation. *JMIR Med. Inform.* **7**, e14325 (2019).
- 876 25. Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A. & Talwalkar, A. Hyperband: A Novel Bandit-
877 Based Approach to Hyperparameter Optimization. *J. Mach. Learn. Res.* **18**, 6765–6816 (2017).
- 878 26. Fry, A. *et al.* Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank
879 Participants With Those of the General Population. *Am. J. Epidemiol.* **186**, 1026–1034 (2017).
- 880 27. Schoeler, T. *et al.* Participation bias in the UK Biobank distorts genetic associations and
881 downstream analyses. *Nat. Hum. Behav.* **7**, 1216–1227 (2023).
- 882 28. Hilton, B. *et al.* Laboratory diagnosed microbial infection in English UK Biobank participants in
883 comparison to the general population. *Sci. Rep.* **13**, 496 (2023).
- 884 29. Deeks, J. J. Systematic reviews of evaluations of diagnostic and screening tests. *BMJ* **323**, 157–
885 162 (2001).
- 886 30. Glas, A. S., Lijmer, J. G., Prins, M. H., Bonsel, G. J. & Bossuyt, P. M. M. The diagnostic odds ratio: a
887 single indicator of test performance. *J. Clin. Epidemiol.* **56**, 1129–1135 (2003).
- 888 31. Beesley, L. J. *et al.* The emerging landscape of health research based on biobanks linked to
889 electronic health records: Existing resources, statistical challenges, and potential opportunities. *Stat.*
890 *Med.* **39**, 773–800 (2020).
- 891 32. Jukarainen, S. *et al.* Genetic risk factors have a substantial impact on healthy life years. *Nat. Med.*
892 **28**, 1893–1901 (2022).

- 893 33. Vos, T. *et al.* Global burden of 369 diseases and injuries in 204 countries and territories, 1990–
894 2019: a systematic analysis for the Global Burden of Disease Study 2019. *The Lancet* **396**, 1204–1222
895 (2020).
- 896 34. Millar, J. The Need for a Global Language - SNOMED CT Introduction. *Stud. Health Technol.*
897 *Inform.* **225**, 683–685 (2016).
- 898 35. Wasserman, H. & Wang, J. An Applied Evaluation of SNOMED CT as a Clinical Vocabulary for the
899 Computerized Diagnosis and Problem List. *AMIA. Annu. Symp. Proc.* **2003**, 699–703 (2003).
- 900 36. Agniel, D., Kohane, I. S. & Weber, G. M. Biases in electronic health record data due to processes
901 within the healthcare system: retrospective observational study. *BMJ* **361**, k1479 (2018).
- 902 37. Forrest, I. S. *et al.* Genetic and phenotypic profiling of supranormal ejection fraction reveals
903 decreased survival and underdiagnosed heart failure. *Eur. J. Heart Fail.* **24**, 2118–2127 (2022).
- 904 38. Appadurai, V. *et al.* Apparent underdiagnosis of Cerebrotendinous Xanthomatosis revealed by
905 analysis of ~60,000 human exomes. *Mol. Genet. Metab.* **116**, 298–304 (2015).
- 906 39. Dickey, A. K. *et al.* Evidence in the UK Biobank for the underdiagnosis of erythropoietic
907 protoporphyria. *Genet. Med.* **23**, 140–148 (2021).
- 908 40. Shoemark, A. *et al.* Genome sequencing reveals underdiagnosis of primary ciliary dyskinesia in
909 bronchiectasis. *Eur. Respir. J.* **60**, (2022).
- 910 41. Damrauer, S. M. *et al.* Association of the V122I Hereditary Transthyretin Amyloidosis Genetic
911 Variant With Heart Failure Among Individuals of African or Hispanic/Latino Ancestry. *JAMA* **322**, 2191–
912 2202 (2019).
- 913 42. Anderson, J. J. *et al.* Ethnic differences in prevalence of actionable HbA1c levels in UK Biobank:
914 implications for screening. *BMJ Open Diabetes Res. Care* **9**, e002176 (2021).
- 915 43. Sturm, A. C. *et al.* Clinical Genetic Testing for Familial Hypercholesterolemia. *J. Am. Coll. Cardiol.*
916 **72**, 662–680 (2018).
- 917 44. An, U., Cai, N., Dahl, A. & Sankararaman, S. AutoComplete: Deep Learning-Based Phenotype
918 Imputation for Large-Scale Biomedical Data. in *Research in Computational Molecular Biology* (ed. Pe’er,
919 I.) 385–386 (Springer International Publishing, 2022). doi:10.1007/978-3-031-04749-7_38.
- 920 45. Stekhoven, D. J. & Bühlmann, P. MissForest—non-parametric missing value imputation for
921 mixed-type data. *Bioinformatics* **28**, 112–118 (2012).
- 922 46. Li, R., Chen, Y. & Moore, J. H. Integration of genetic and clinical information to improve
923 imputation of data missing from electronic health records. *J. Am. Med. Inform. Assoc.* **26**, 1056–1063
924 (2019).
- 925 47. Shaw, D. M. *et al.* Phenome risk classification enables phenotypic imputation and gene discovery
926 in developmental stuttering. *Am. J. Hum. Genet.* **108**, 2271–2283 (2021).

- 927 48. Bernardini, M., Doynychko, A., Romeo, L., Frontoni, E. & Amini, M.-R. A novel missing data
928 imputation approach based on clinical conditional Generative Adversarial Networks applied to EHR
929 datasets. *Comput. Biol. Med.* **163**, 107188 (2023).
- 930 49. Beaulieu-Jones, B. K. *et al.* Characterizing and Managing Missing Structured Data in Electronic
931 Health Records: Data Analysis. *JMIR Med. Inform.* **6**, e8960 (2018).
- 932 50. Wells, B. J., Chagin, K. M., Nowacki, A. S. & Kattan, M. W. Strategies for Handling Missing Data in
933 Electronic Health Record Derived Data. *eGEMs* **1**, 1035 (2013).
- 934 51. Madden, J. M., Lakoma, M. D., Rusinak, D., Lu, C. Y. & Soumerai, S. B. Missing clinical and
935 behavioral health data in a large electronic health record (EHR) system. *J. Am. Med. Inform. Assoc.* **23**,
936 1143–1149 (2016).
- 937 52. Köpcke, F. *et al.* Evaluation of data completeness in the electronic health record for the purpose
938 of patient recruitment into clinical trials: a retrospective analysis of element presence. *BMC Med. Inform.*
939 *Decis. Mak.* **13**, 1–8 (2013).
- 940 53. Hripcsak, G. & Albers, D. J. Next-generation phenotyping of electronic health records. *J. Am.*
941 *Med. Inform. Assoc.* **20**, 117–121 (2013).
- 942 54. Abul-Husn, N. S. *et al.* Genetic identification of familial hypercholesterolemia within a single U.S.
943 health care system. *Science* **354**, (2016).
- 944 55. Marwaha, S., Knowles, J. W. & Ashley, E. A. A guide for the diagnosis of rare and undiagnosed
945 disease: beyond the exome. *Genome Med.* **14**, 23 (2022).
- 946 56. Basile, A. O. & Ritchie, M. D. Informatics and machine learning to define the phenotype. *Expert*
947 *Rev. Mol. Diagn.* **18**, 219–226 (2018).
- 948 57. Onnela, J.-P. Opportunities and challenges in the collection and analysis of digital phenotyping
949 data. *Neuropsychopharmacology* **46**, 45–54 (2021).
- 950 58. McArthur, E., Bastarache, L. & Capra, J. A. Linking rare and common disease vocabularies by
951 mapping between the human phenotype ontology and phecodes. *JAMIA Open* **6**, ooad007 (2023).
- 952 59. Kingdom, R. *et al.* Rare genetic variants in genes and loci linked to dominant monogenic
953 developmental disorders cause milder related phenotypes in the general population. *Am. J. Hum. Genet.*
954 **109**, 1308–1316 (2022).
- 955 60. Bakker, O. B. *et al.* *Linking common and rare disease genetics through gene regulatory networks.*
956 2021.10.21.21265342 <https://www.medrxiv.org/content/10.1101/2021.10.21.21265342v2> (2021)
957 doi:10.1101/2021.10.21.21265342.
- 958 61. Köhler, S. *et al.* The Human Phenotype Ontology in 2021. *Nucleic Acids Res.* **49**, D1207–D1217
959 (2021).
- 960 62. Vasant, D. *et al.* *ORDO: An Ontology Connecting Rare Disease, Epidemiology and Genetic Data.*
961 (2014).

- 962 63. Connolly, B. *et al.* Natural Language Processing – Overview and History. in *Pediatric Biomedical*
963 *Informatics: Computer Applications in Pediatric Research* (ed. Hutton, J. J.) 203–230 (Springer, 2016).
964 doi:10.1007/978-981-10-1104-7_11.
- 965 64. Smoller, J. W. The use of electronic health records for psychiatric phenotyping and genomics.
966 *Am. J. Med. Genet. B Neuropsychiatr. Genet.* **177**, 601–612 (2018).
- 967 65. Griggs, R. C. *et al.* Clinical research for rare disease: Opportunities, challenges, and solutions.
968 *Mol. Genet. Metab.* **96**, 20–26 (2009).
- 969 66. Stoller, J. K. The Challenge of Rare Diseases. *Chest* **153**, 1309–1314 (2018).
- 970 67. Banerjee, J. *et al.* Machine learning in rare disease. *Nat. Methods* 1–12 (2023)
971 doi:10.1038/s41592-023-01886-z.
- 972 68. De Freitas, J. K. *et al.* Phe2vec: Automated disease phenotyping based on unsupervised
973 embeddings from electronic health records. *Patterns* **2**, 100337 (2021).
- 974 69. Kane, M. J. *et al.* A Compressed Language Model Embedding Dataset of ICD 10 CM Descriptions.
975 2023.04.24.23289046 Preprint at <https://doi.org/10.1101/2023.04.24.23289046> (2023).
- 976 70. Sanjak, J., Zhu, Q. & Mathé, E. A. Clustering rare diseases within an ontology-enriched
977 knowledge graph. 2023.02.15.528673 Preprint at <https://doi.org/10.1101/2023.02.15.528673> (2023).
- 978 71. Xie, F. *et al.* Deep learning for temporal data representation in electronic health records: A
979 systematic review of challenges and methodologies. *J. Biomed. Inform.* **126**, 103980 (2022).
- 980 72. Kraljevic, Z. *et al.* Foresight - Generative Pretrained Transformer (GPT) for Modelling of Patient
981 Timelines using EHRs. (2022). doi:10.48550/arXiv.2212.08072.
- 982 73. Sudlow, C. *et al.* UK Biobank: An Open Access Resource for Identifying the Causes of a Wide
983 Range of Complex Diseases of Middle and Old Age. *PLOS Med.* **12**, e1001779 (2015).
- 984 74. Chollet, F. & et al. Keras. (2015).
- 985 75. Watson, M., Qian, C., Bischof, J., Chollet, F. & et al. KerasNLP. (2022).
- 986 76. O’Malley, T., Bursztein, E., Long, J., Chollet, F. & et al. KerasTuner. (2019).
- 987 77. Gart, J. J. Alternative Analyses of Contingency Tables. *J. R. Stat. Soc. Ser. B Methodol.* **28**, 164–
988 179 (1966).
- 989 78. Seabold, S. & Pertkold, J. statsmodels: Econometric and statistical modeling with python. in
990 *Proceedings of the 9th Python in Science Conference* (2010).
- 991 79. International Health Terminology Standards Development Organisation,. SNOMED-CT Editorial
992 Guide. <http://snomed.org/eg>.
- 993 80. Mölder, F. *et al.* Sustainable data analysis with Snakemake. Preprint at
994 <https://doi.org/10.12688/f1000research.29032.2> (2021).