

Deep learning for polygenic prediction: The role of heritability, interaction type and sample size

Jason Grealey^{1,2}, Gad Abraham^{1,3,^}, Guillaume Méric^{1,4}, Rodrigo Cánovas¹, Martin Kelemen^{3,5,6}, Shu Mei Teo^{1,3}, Agus Salim^{1,2}, Michael Inouye^{1,3,5-9,*,#}, Yu Xu^{3,5,6,*,#}

1. Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute, 75 Commercial Rd, 7 Melbourne 3004, Victoria, Australia
2. Department of Mathematics and Statistics, La Trobe University, VIC 3086, Australia
3. Cambridge Baker Systems Genomics Initiative, Department of Public Health and Primary Care, University of Cambridge, Cambridge CB1 8RN, UK
4. Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Australia
5. British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, CB1 8RN, UK
6. Victor Phillip Dahdaleh Heart and Lung Research Institute, University of Cambridge, Cambridge, UK
7. Health Data Research UK Cambridge, Wellcome Genome Campus and University of Cambridge, Cambridge, CB10 1SA, UK
8. British Heart Foundation Cambridge Centre of Excellence, Department of Clinical Medicine, University of Cambridge, Cambridge, UK
9. The Alan Turing Institute, London, United Kingdom

[^] Current address: CSL Innovation

*Corresponding author: YX (yx322@medschl.cam.ac.uk), MI (Michael.inouye@baker.edu.au)

#Contributed equally

Keywords: Polygenic scores, genomic prediction, neural networks, non-linear interactions, machine learning, deep learning, encoding

33 **Abstract**

34

35 Polygenic scores (PGS), which aggregate the effects of genetic variants to estimate
36 predisposition for a disease or trait, have potential clinical utility in disease prevention and
37 precision medicine. Recently, there has been increasing interest in using deep learning (DL)
38 methods to develop PGS, due to their strength in modelling complex non-linear relationships
39 (such as GxG) that conventional PGS methods may not capture. However, the perceived value
40 of DL for polygenic scores is unclear. In this study, we assess the underlying factors impacting
41 DL performance and how they can be better utilised for PGS development. We simulate large-
42 scale realistic genotype-to-phenotype data, with varying genetic architectures of phenotypes
43 under quantitative control of three key components: (a) total heritability, (b) variant-variant
44 interaction type, and (c) proportion of non-additive heritability. We compare the performance
45 of one of most common DL methods (multi-layer perceptron, MLP) on varying training sample
46 sizes, with two well-established PGS methods: a purely additive model (pruning and
47 thresholding, P+T) and a machine learning method (Elastic net, EN). Our analyses show EN
48 has consistently better overall performance across traits of different architectures and training
49 data of different sizes. However, MLP saw the largest performance improvements as sample
50 size increases. MLP outperformed P+T for most traits and achieves comparable performance
51 as EN for numerous traits at the largest sample size assessed (N=100k), suggesting DL may
52 offer some advantages in future when they can be trained on biobanks of millions of samples.
53 We further found that one-hot encoding of variant input can improve performance of every
54 method, particularly for traits with non-additive variance. Overall, we show how different
55 underlying factors impact how well methods leverage non-additivity for polygenic prediction.

56

57 Introduction

58
59 Polygenic scores (PGS), which aggregate the effects of many genetic variants into a single
60 number, have become an important tool to predict the genetic predisposition of an individual
61 towards a phenotype and have been shown to have promising utility such as in disease
62 prevention and precision medicine¹⁻³. There is increasing interest in using deep learning (DL)
63 approaches to develop PGS of complex traits⁴⁻¹¹. Known as universal function
64 approximators^{12,13}, the value of deep learning models is in their ability to model complex non-
65 linear effects among genetic variants and their flexibility in combining with other non-genetic
66 factors for subsequent applications (e.g. disease related biomarkers and environmental factors
67 for disease risk models).

68
69 Human traits, including quantitative traits and diseases, are heritable to varying degrees and
70 many of them have been found to have a highly polygenic architecture (i.e., their variance is
71 accounted for by many thousands or even millions of genetic variants genome-wide)¹⁴. While
72 studies have shown that for most phenotypes^{15,16} the associated variants contribute largely in a
73 linearly additive manner, non-linear interaction effects (GxG) are present and sometimes make
74 a substantial contribution to the genetic variation of phenotypes, e.g. autoimmune diseases^{7,17}.

75
76 It has been shown that common machine learning methods, such as elastic net and gradient
77 boosting trees, can capture GxG in the genetic prediction of common traits and diseases^{5,18,19},
78 frequently improving PGS performance. While these methods do not explicitly model
79 interaction terms, GxG can still be captured to an extent through variant encoding or inherently
80 non-linear structures of the model. Deep learning methods readily model complex non-linear
81 relationships and have recently been proposed for PGS development of various human traits<sup>4-
82 7,20,21</sup>. DL methods have been found to improve PGS of several traits and diseases, such as breast
83 cancer⁶, Alzheimer's disease¹⁰ and type 1 diabetes⁷, but in many cases substantially improved
84 performance over simpler machine learning models has not been found^{4,5,7}. DL methods may
85 also be susceptible to confounding by joint tagging effects, whereby GxG is in fact attributed
86 to unaccounted additive genetic variants, and only provide moderate improvements in
87 prediction performance even under extreme genetic architectures.²²

88
89 Despite substantial efforts, it remains unclear under what conditions (if at all) DL may offer an
90 advantage over simple approaches to polygenic score construction. Here, we investigate how
91 and to what extent key factors of genetic architecture and sample size affect the performance
92 of PGS models and in particular, under what circumstances DL methods outperform linear
93 models. To answer these questions, we simulated genotype data of 100,000 individuals with
94 realistic linkage disequilibrium (LD) patterns, and phenotypes whose genetic architectures
95 were of varying: (a) total heritability (broad sense), (b) types of GxG interaction, and (c)
96 proportion of non-additive heritability. We compared the performance of a suite of common
97 methods for PGS development of these simulated phenotypes, which included a univariate
98 linear model (pruning and thresholding), a regularized linear regression (elastic net), and a deep
99 learning approach (multi-layer perceptron). We also investigated the impact of training data
100 sizes and variant encoding types on the performance of these methods. Our findings inform
101 study designs and methodology selection for future PGS development.

102
103
104

105 Methods

106 107 Simulating genotypes

108
109 HAPGEN2²³ was used to simulate genotypes with realistic linkage disequilibrium (LD)
110 patterns. As an empirical reference panel from which to draw haplotypes, we utilised
111 chromosome 22 (171,457 variants in total) from 99 individuals of the phase 3 of the 1000
112 Genomes project²⁴ (Finnish subset). This reference panel was used to generate 100,000
113 simulated individual haplotypes which after conversion to genotypes contained 100,455
114 variants (after keeping variants with minor allele frequency (MAF) between 1% and 40%). No
115 variant was found to violate the Hardy-Weinberg equilibrium ($p < 10^{-6}$) using PLINK²⁵.

116 117 Simulating phenotypes

118
119 Phenotypes (in this study, continuous traits) were simulated using the simulated genotypes
120 above, where genotypes were coded in a minor allele dosage format $\{0, 1, 2\}$. For each
121 phenotype, a total of 1,000 variants were randomly chosen ($\sim 1\%$ of the total variants in the
122 simulated dataset) and were given an effect size randomly drawn from a normal distribution
123 with a standard deviation σ_β (others have effect sizes of 0):

$$124 \quad \beta_j \sim N(0, \sigma_\beta) \quad (1).$$

125 After all the 1,000 variants were given a linear effect size, as drawn from equation (1) with σ_β
126 initialised at 0.01, these effect sizes were used to scale the non-additive heritability for the trait.
127 Of these 1,000 causal variants, if the trait was influenced by GxG, 500 of them (250 variant-
128 variant pairs are randomly sampled) were given non-additive effects, which were modelled
129 according to the two locus interaction types from Li and Reich²⁶. A non-additive effect was
130 simulated under a given combination of effect alleles for both variants according to four
131 interaction types: threshold (“T”), recessive/recessive (“RR”), exclusive/or (“XOR”), and
132 heterozygote/heterozygote (“HH”, previously named as “m16”²⁶) (**Figure 4**). If an individual
133 contains this specific combination of effect alleles, these variants will exhibit a variant
134 interaction effect on their phenotype.

135
136 The GxG interaction effect for a given pair of variants k of sample i is determined by the
137 following equation:

$$138 \quad I_{i,k} = Z_{i,k} \gamma_k \quad (2)$$

139 where: $\gamma_k \sim N(0, \sigma_\gamma)$

140
141 where $Z_{i,k}$ is an indicator function for the GxG interaction types (**Figure 4**) for interacting
142 variant pair k of sample i and is either 1 or 0 depending on the combination of genotypes and
143 the GxG interaction type; if $Z_{i,k} = 1$ (i.e. a given interaction type exhibits), the interaction
144 effect size is drawn from a normal distribution. σ_γ is initialised at 0.01 and scaled with respect
145 to the total (i.e. additive and non-additive) genetic variation within the phenotype to control
146 the level of non-additive variation contributing towards the phenotype (See below).

147
148 As well as GxG interaction effects and linear effects, there was also a proportion of noise in
149 the phenotype that genetics do not explain, i.e. a non-heritable contribution. This was modelled
150 as follows:

151
$$\epsilon_i \sim N(0, \sigma_{Noise}) \quad (3)$$

152

153 where σ_{Noise} is scaled to fix the heritability of the trait.

154

155 The above equations combine for the phenotype like so:

156
$$P_i = G_i + I_i + \epsilon_i \quad (4)$$

157
$$\text{where, } G_i = \sum_j^m x_{i,j} \beta_j, I_i = \sum_k^S Z_{i,k} \gamma_k,$$

158
$$\beta_j \sim N(0, \sigma_\beta), \quad \gamma_k \sim N(0, \sigma_\gamma), \quad \epsilon_i \sim (0, \sigma_{Noise}),$$

159 G_i , I_i and ϵ_i are the combined linear effects, combined non-additive effects for variants
160 exhibiting interactions and the noise for sample i respectively; $x_{i,j}$ is the number of effect
161 alleles present in variant j of sample i ; β_j is the linear effect size drawn from equation (1) for
162 the effect allele in variant j ; $Z_{i,k}$ determines if the k^{th} pair of interacting variants exhibits a
163 certain GxG interaction type in sample i where 4 interaction types are applied using the two
164 locus penetrance tables in this study (**Figure 4**). The phenotype value for sample i is
165 determined by summing all contributions from the linear effects of m simulated variants, the
166 GxG interaction effects of S pairs of simulated interaction variants, and the noise.

167

168 **Heritability simulation**

169

170 After the first step of initialization of phenotype values (i.e. its noise, non-additive and linear
171 components) as described above, we then performed linear regression to scale the contribution
172 of each component to control the non-additive contribution to the total heritability and its total
173 heritability or broad sense heritability for the purpose of simulating phenotypes of different
174 settings²⁷.

175

176 As described above, the sum of all genetic effects on a given phenotype for individual i is as
177 follows:

178

179
$$S_i = G_i + I_i.$$

180

181 To determine the proportion of non-additive heritability in the total genetic effects, we
182 performed the following linear regression across all individuals:

183
$$S_i \sim I_i \quad (5)$$

184 where the goodness-of-fit (R^2) of the regression determines the total non-additive contribution
185 to the heritability. For example, if the R^2 was 20%, then only 80% of the total heritability
186 would be narrow sense (linear additive) and the other 20% being non-additive, i.e. from GxG.
187 These non-additive effects are increased or decreased by scaling all the pairwise interacting
188 effects to obtain the required level of non-additive variance in the trait. The linear regression
189 was performed using Scikit-learn python package²⁸.

190

191 Similarly, we performed the following regression to determine the broad sense heritability of
192 a given trait:

193
$$P_i \sim \epsilon_i \quad (6)$$

194 with which, noise in the trait is increased or decreased to obtain the required broad sense
195 heritability.

196 **Summary of the simulated dataset**

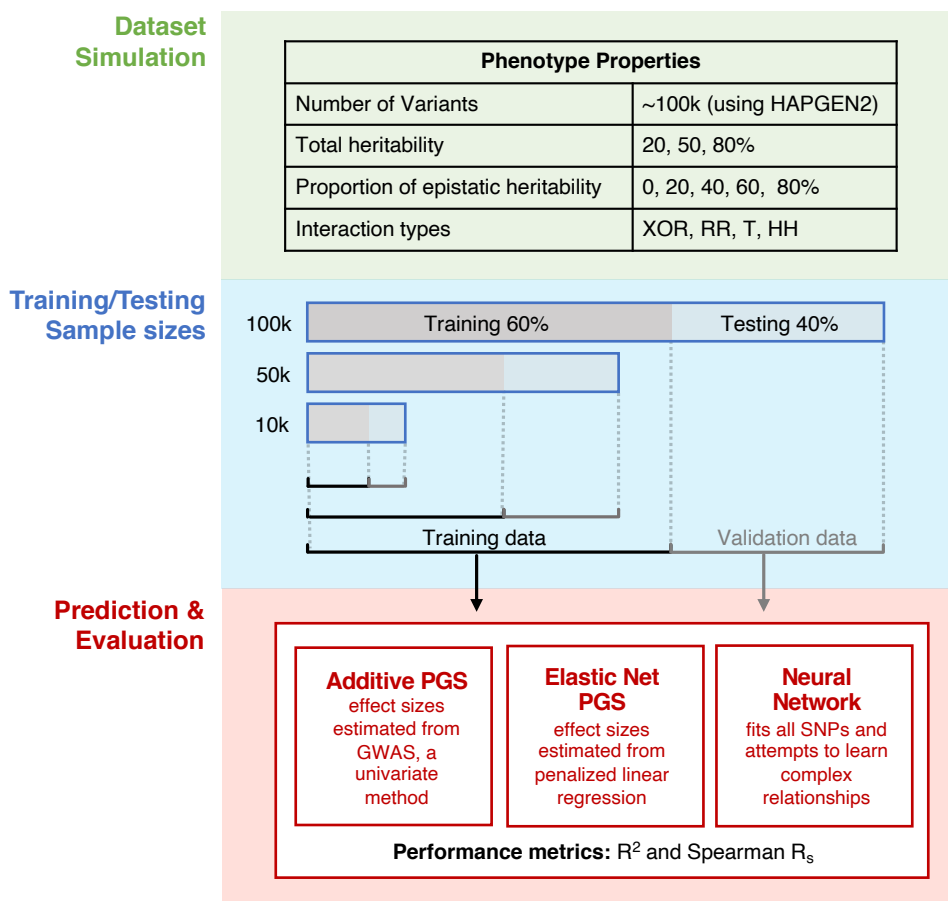
197

198 In total, we simulated 60 phenotypes of different settings, which were under control of three
199 parameters: (a) total heritability (20, 50, or 80%), (b) GxG interaction type (HH, RR, XOR, or
200 T)²⁶, and (c) proportion of non-additive heritability (0, 20, 40, 60, or 80%).

201

202 For each phenotype, 500 variants were randomly selected and given a linear contribution for
203 the phenotype. The remaining 500 variants were given linear and paired GxG interaction
204 effects for a given epistatic model, where variants are randomly selected to generate 250 non-
205 overlapping pairs. In the case of no non-additive effects, these 500 variants were only given
206 linear effects. These effects were summed for all 1000 variants, noise was added, and
207 heritability and the proportion of non-additive heritability were fixed. Three predictive models:
208 (i) additive PGS, (ii) elastic net PGS and (iii) feed forward neural networks, were used to
209 develop polygenic scores for these phenotypes, which are detailed below.

210



211

212 **Figure 1. Schematic study design for data simulation and genetic prediction.** Simulations of 100k genotypes were
213 generated using HAPGEN2 and subsampled to smaller datasets of 50k and 10k samples. Using these simulated genotypes,
214 traits with different settings of heritability, GxG interaction type, and proportion of non-additive heritability were generated.
215 The samples were split into training and testing sets in each dataset of 100k, 50k and 10k samples, after which they were used
216 to train and test the prediction methods (neural networks, elastic net PGS and additive PGS).

217 Three sample sets of different sizes (100k, 50k, and 10k) were randomly selected from the
218 simulated dataset for each phenotype, each of which was then split 60/40% into training and
219 testing sets (**Figure 1**). Then every prediction model used the same generated sample sets to
220 train PGS models and test their performance.

221 Additive PGS method P+T

222 This additive PGS method assumes that the genetic variants have linear additive effects on
223 PGS of the trait, and develops PGS of a trait using the weighted sum of genotypes of the
224 selected variants for that trait:

$$225 \quad S_i = \sum_{j=1}^m \beta_j x_{i,j} \quad (6)$$

226 where S_i is a polygenic score for individual i ; $x_{i,j}$ is the genotype dosage of variant j of the
227 individual i ; the β_j is the effect size of the variant j that is usually obtained through the
228 univariate statistical association tests on training data; the m variants were often selected
229 through a LD pruning/clumping and p-value thresholding step²⁹ (so this method is often named
230 as $P+T$). The software PLINK²⁵ was used to estimate the univariate effect sizes for each
231 simulated variant on the training data of a given simulated phenotype. Using these univariate
232 estimations, PRSice-2³⁰ was then employed to develop PGS of the phenotype on the training
233 data. PRSice-2 performs LD clumping to reduce the correlation amongst variants and then tests
234 thousands of optimised p-value filtered PGSs to obtain the most predictive PGS.

235

236 Elastic net PGS method

237

238 Elastic net (EN) also assumes variants have a linear additive effect, estimated via penalised
239 regression, where all the variants are jointly fit together. SparSNP³¹ is a tool designed to fit
240 penalised linear models for genetic prediction, and was used to perform elastic-net regression
241 in this study. SparSNP minimizes the following loss function for estimating the effect sizes of
242 genetic variants:

243

$$244 \quad L(\beta_0, \beta) = \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \mathbf{x}_i^T \boldsymbol{\beta})^2 + \lambda_1 \sum_{j=1}^m |\beta_j| + \frac{\lambda_2}{2} \sum_{j=1}^m |\beta_j|^2 \quad (7)$$

245

246 where, y_i is the simulated phenotype value; x_i are the genotype dosages of the m variants for
247 sample i (i.e. the 100,455 simulated variants after quality control); β is the vector of effect
248 sizes for the m variants; β_0 is the intercept term; λ_1 and λ_2 are the penalties for a LASSO and
249 Ridge regularisation respectively. A 5-fold cross validation was performed for 10 times in
250 SparSNP to select the optimal λ_1, λ_2 pair on a given training set, where λ_2 was set as 0.2 and
251 λ_1 was identified from a default set of thirty options in SparSNP. Finally, effect sizes of
252 variants were estimated by minimizing equation (7) on the training set, which are then used to
253 construct the PGS using equation (6).

254

255 Neural networks

256

257 Multilayered perceptrons (MLPs; also called feed-forward neural networks) are one of most
258 common neural network architectures and can improve genetic prediction of quantitative traits,
259 e.g. blood cell traits⁵. MLPs do not make any assumptions about the distributions behind the
260 data they fit, and can be trained to approximate any smooth function in theory¹³. They usually
261 consist of nodes (functions) connected to many other layers through directed acyclic graphs¹³;
262 the output of a layer is used as the input to subsequent layers and element-wise transformed by
263 non-linear activation functions, which allows for it to model complex correlations. Given m

264 nodes in the L th layer, the output of a node i in the next or $(L + 1)$ th layer is calculated like
265 so:

$$266 \quad output_{i,L+1} = activation \left(\sum_{a=1}^n input_{a,L} weight_a + bias \right) \quad (8)$$

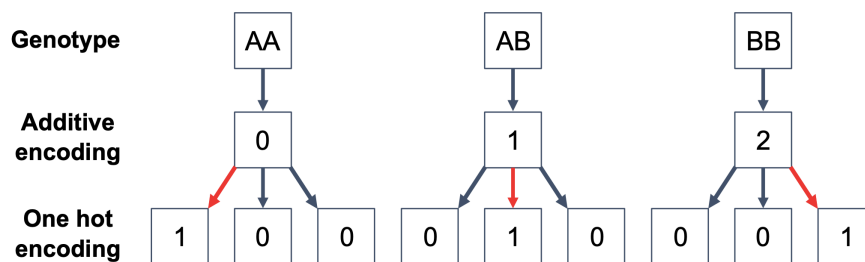
267 where $output_{i,L+1}$ is the output of a node i in the $(L + 1)$ th layer; $input_{a,L}$ is the input from
268 node a (from the total n nodes) in the previous layer. Each $input_{a,L}$ is multiplied by a weight
269 (i.e. $weight_a$) and added to a $bias$, which are then summed up and passed into an activation
270 function as the output of node i . As mentioned above these activation functions are used to
271 incorporate non-linearity to the modelling process. This process occurs from the first hidden
272 layer, to any hidden layers until the output is reached (**Supplementary Figures 13-14**). MLP
273 models were implemented with Keras (keras.io) and TensorFlow³² in our study.

274

275 Genotype Encoding

276

277 We considered two different types of genotype encoding, additive encoding (i.e. effect allele
278 dosage) and one hot encoding, as the input of the prediction models in this study (**Figure 2**).
279 The first encoding involved the use of Plink v1.9²⁵ to encode the variants into allelic dosages,
280 where the variants were input as counts of the effect alleles (i.e. 0, 1, 2). In “one hot” encoding,
281 variants were encoded into the absence or presence of their genotype classes.



282

283 **Figure 2. Schematics of genotype encoding.** This schematic shows the two different variants encodings used in this study
284 and shows how genotypes are represented in additive and one hot encodings.

285 Hyperparameter optimisation

286

287 An essential component in neural network model training, in particular MLP in this study, is
288 hyperparameter optimisation. Hyperparameters are variables that dictate the network’s
289 structure, its complexities and training process, which are set before the model training. Each
290 set of hyperparameters can perform differently on a given task, thus a search must be conducted
291 to determine the optimal set for each task. The hyperparameter search was conducted using
292 Talos³³ package which aids in performing the random search for the best set of hyperparameters
293 for a given task (i.e. predicting each phenotype) on the training data. Given a list of
294 hyperparameters to optimise from, Talos randomly searches this list to create numerous unique
295 combinations of hyperparameters (see details in **Supplementary Table 1**), which were used
296 to determine the best performing set of hyperparameters for a given phenotype.

297 Assessment of prediction accuracy

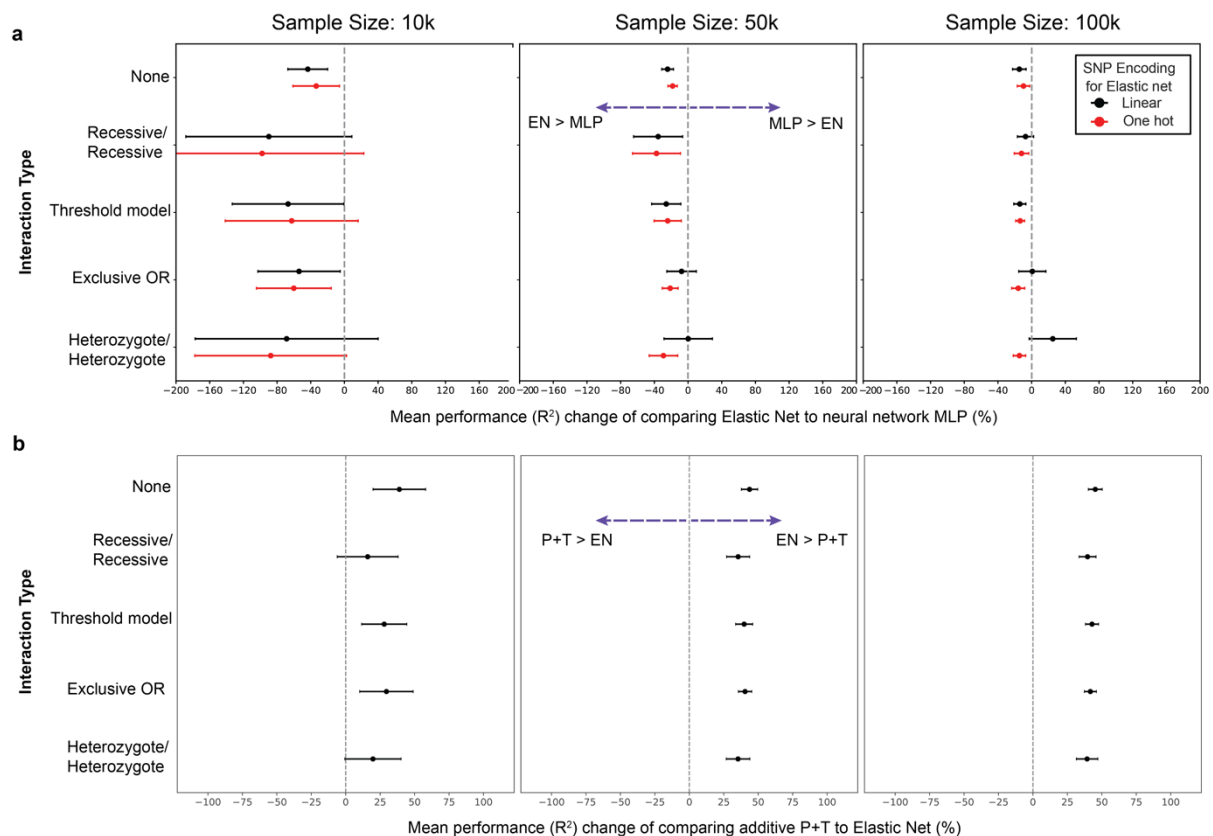
298

299 Finally, the two metrics: coefficient of determination (R^2) and Spearman correlation coefficient
300 (R_s), were used to measure the performance of each PGS method on the testing data of any
301 given phenotype setting as described above.

301

302 Results

303
304 In this study, we simulated genotype data of 100,000 individuals using the 1000G reference
305 panel, with which we further simulated 60 phenotypes of different settings, including (a) total
306 heritability (20, 50, or 80%), (b) GxG interaction type (HH, RR, XOR, or T)²⁶, and (c)
307 proportion of non-additive heritability (0, 20, 40, 60, or 80%). We then evaluated the
308 performance of three polygenic score methods, including a simple additive PGS method (the
309 pruning and thresholding), a linear machine learning method (elastic net) and a deep learning
310 method (multilayered perceptron), in predicting these simulated phenotypes using training data
311 of different sizes. Two different types of genotype encoding: additive dosage encoding and one
312 hot encoding, were also applied to investigate its impact on the performance of PGS methods.
313

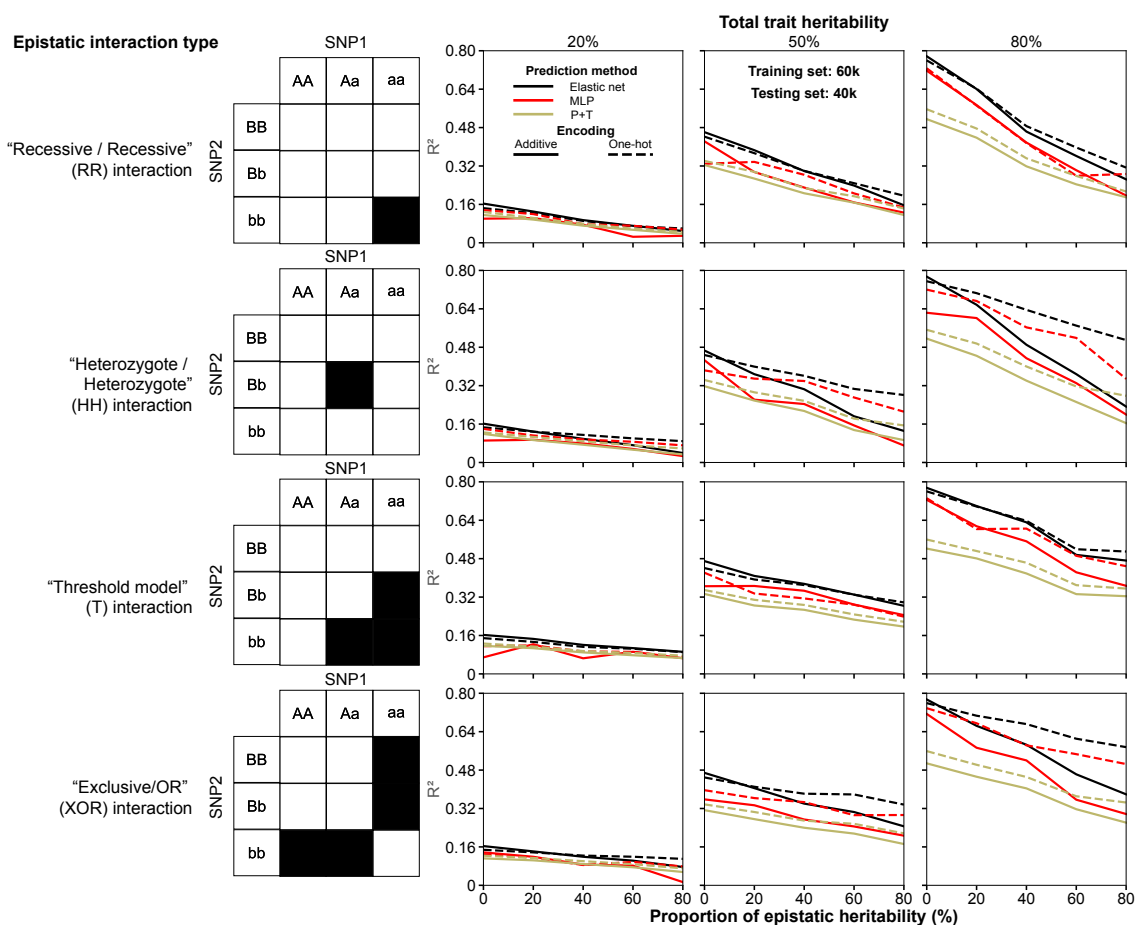


314
315 **Figure 3. Performance comparison of different PGS methods. a. Elastic net is more accurate than neural network MLP**
316 **in PGS development.** Each sub-plot shows the percentage change in the predictive performance (R^2) between Elastic net and
317 MLP at a given sample size. Simulated traits are grouped by interaction and variant encoding types, then we compare the
318 performance between Elastic net and MLP by using the mean and standard deviation of $(R_{MLP}^2 - R_{EN}^2)/R_{EN}^2$ in a selected trait
319 group, where R_{EN}^2 is the R^2 performance of EN on a trait. Note that both the one hot (red) and additively (black) encoded
320 elastic net PGS methods are compared against one hot encoded MLP. **b. Elastic net outperforms additive P+T method.**
321 Each sub-plot shows the percentage change in the predictive performance (R^2) of additive PGS method P+T and elastic net at
322 a given total sample size, which are measured using the mean and standard deviation of $(R_{EN}^2 - R_{P+T}^2)/R_{P+T}^2$ in each trait
323 group by GxG interaction type.

324 Performance of polygenic prediction methods across different settings

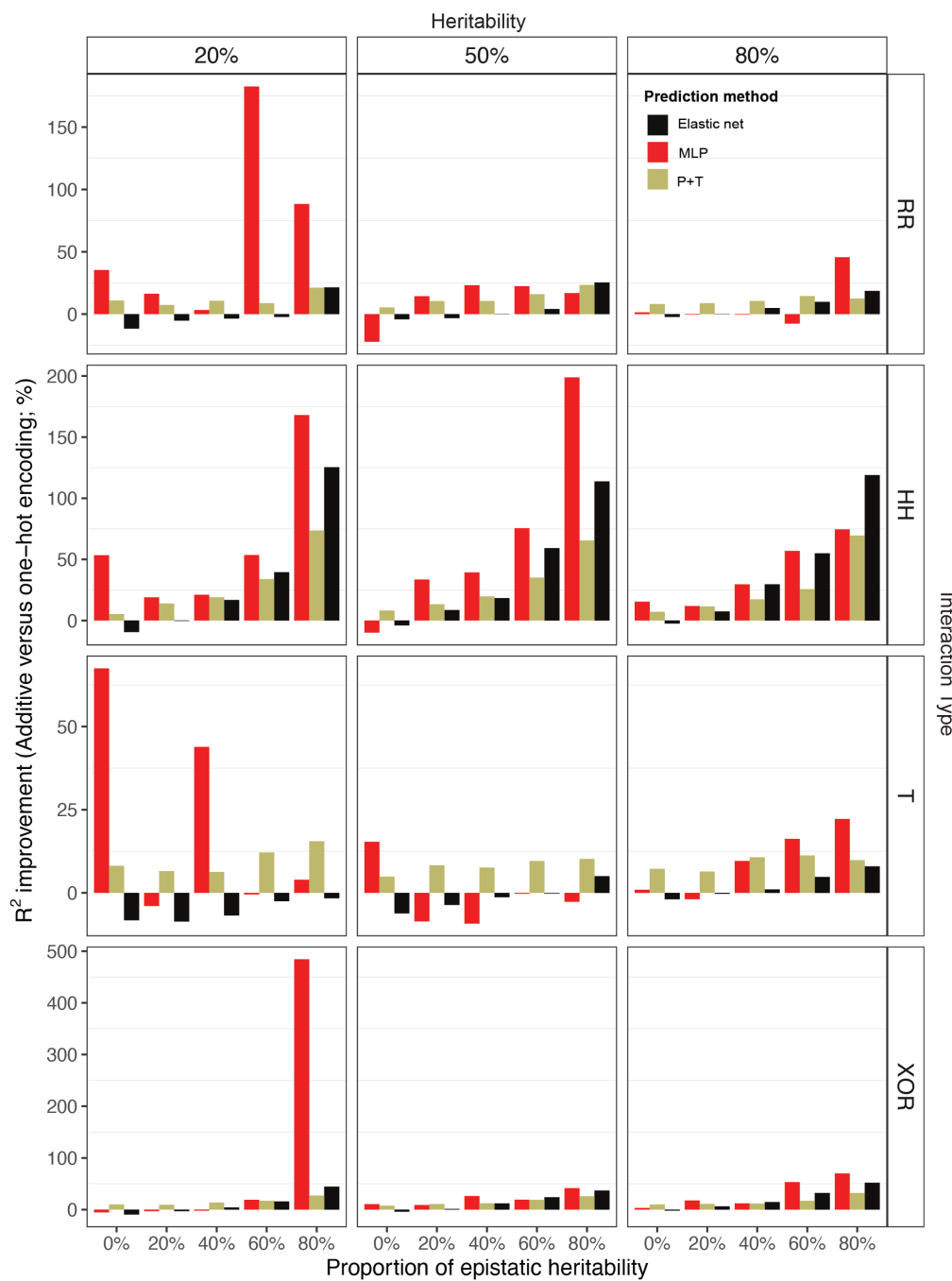
325
326
327 Across various simulation settings, elastic net performed consistently well across phenotypes
328 and training datasets when compared to the other methods (**Figures 3-4 and Supplementary**
329 **Figures 1-5**). Using the 10k simulated data (training sample size: 6k), EN improved R^2 by 22%
330 on average over the simple additive method (P+T) across traits of different GxG interaction

331 types and neural network model MLP underperformed EN by 65% in terms of R^2 (additive
 332 variant encoding for P+T and EN; **Figure 3**). However, at this smallest sample size applied,
 333 P+T slightly performed better than EN for a few traits with low heritability and high
 334 proportions of non-additive heritability, e.g. P+T improves over EN by 3% (R^2) for traits of
 335 20% heritability and 60% non-additive heritability (additive variant encoding; **Supplementary**
 336 **Figure 1** and **Supplementary Table 2**). As the sample size increased, EN continued to
 337 outperform both P+T and MLP for most traits. EN is commonly used with additive variant
 338 encoding ('dosage model'), and the MLP with one-hot encoding frequently outperformed this
 339 approach. In particular, for traits with HH interactions, MLP had R^2 25% greater than EN with
 340 additive variant encoding when sample size was 100k. However, using one hot encoding
 341 enabled EN to outperform MLP (**Figure 3a**).



342
 343 **Figure 4. Heritability and proportion of non-additive heritability affect prediction performances of different PGS**
 344 **methods (sample size: 100k).** Predictive performance (R^2) of all the three PGS methods with the two variant encoding types
 345 were compared for traits of different groups, where elastic net PGS are in black, neural network MLP in red and additive P+T
 346 in yellow; solid lines and dashed lines represent additively and one hot encoded models respectively. The plots detail each
 347 PGS method's performance for a given phenotype. Each row has the same underlying interaction model labelled by the tables
 348 and each column has the same total heritability as noted by the column title and within each sub-plot, the x-axis details the
 349 proportion of non-linear contribution for a given trait.

350



351

352 **Figure 5. R² performance improvement of PGS methods using the one-hot encoding over the additive encoding for**
 353 **traits of different groups (sample size: 100k).** Each row has the same interaction model as noted by the row title and each
 354 column has the same total heritability as noted by the column title and within each sub-plot, the x-axis shows the proportion
 355 of non-linear contribution for a given trait. The improvement is calculated using $\frac{R_{one-hot}^2 - R_{additive}^2}{R_{additive}^2}$, where $R_{one-hot}^2$ is the R²
 356 performance of a PGS method using one-hot encoding for a given trait.

357

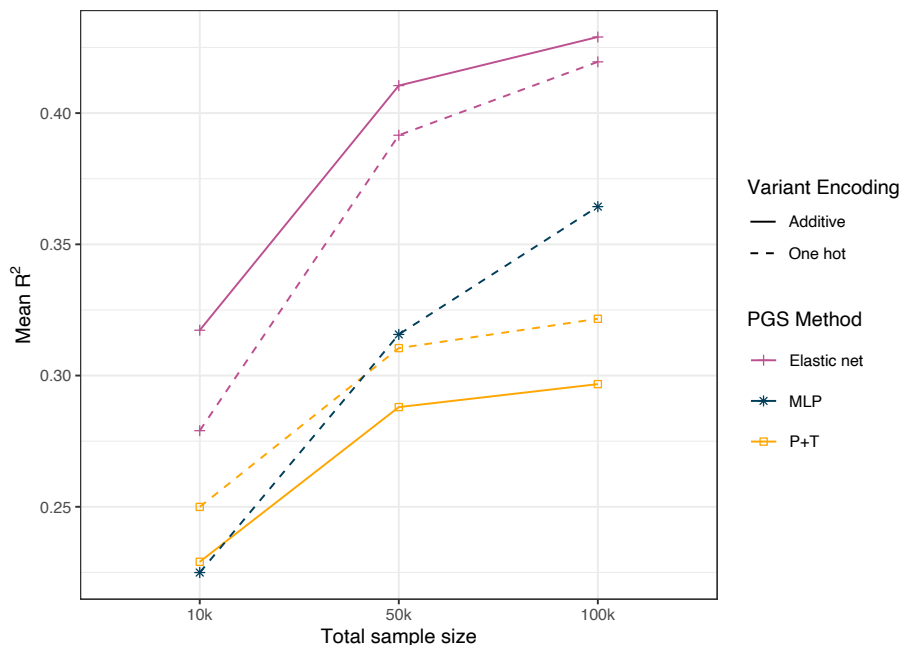
358 One hot encoding of variants frequently improved polygenic prediction

359 We further assessed how using one-hot encoding as variant inputs affect the accuracy of the
 360 three PGS methods. Overall, relative to additive encoding ('dosage model'), we found one-hot
 361 encoding improved the predictive accuracy (R²) for traits with non-additive variance on
 362 average by 42% for MLP, 14% for EN, and 20% for additive P+T PGS method (**Figures 3a,**
 363 **4-5 and Supplementary Figures 6-7**). However, for purely additive traits, one-hot encoding

364 resulted in a percentage change in R^2 performance of +14%, -9% and +8%, when using MLP,
365 EN and P+T respectively. This indicates that one hot encoding will consistently increase MLP
366 performance but may worsen EN's performance, depending on the genetic architecture (which
367 is rarely known *a priori*). Note that (i) the MLP using additive encoding were only tested on
368 data with total sample size of 100k due to its poor performance on smaller data sets and
369 extremely expensive training costs (both significant time and computational resource needed),
370 and (ii) performance gains at relatively low heritability (20%) can be susceptible to substantial
371 noise and should be interpreted with caution.

372 Sample size and relative performance of polygenic prediction method

373 We examined how increasing the sample size of training data results in an increased predictive
374 power for each PGS method (**Figure 6, Supplementary Figure 8**). One hot encoded MLP saw
375 the largest increase in predictive accuracy: increasing the sample size from 10k to 100k yielded
376 a 62% increase in mean R^2 scores for traits with 50% heritability and 20% or less GxG (0.22 to
377 0.36), compared to less than 50% for other PGS methods (**Figure 6**). Across all 60 simulated
378 traits, the MLP's mean R^2 improved by 113% (0.15 to 0.32), while other methods showed less
379 than 56% improvement (**Supplementary Figure 8**). Similarly, increasing the sample size from
380 50k to 100k resulted in the MLP achieving 15% mean R^2 improvement for traits with 50%
381 heritability and 20% or less GxG, and a 19% improvement across all simulated traits. In
382 contrast, other methods showed less than a 7% increase for both trait groups. In addition, one
383 hot encoding allowed EN to gain larger improvements when sample size increases compared
384 with using additive variant encoding. For example, increasing from 10k to 50k, the one hot
385 encoded EN had a mean R^2 increase of 42% (from 0.24 to 0.34), but the additively encoded
386 EN saw a smaller improvement of 29% across all simulated traits.



387
388
389 **Figure 6. The R^2 performance improvements of PGS methods by sample size increase.** This plot shows the mean R^2
390 performance of each PGS method (P+T, elastic net and MLP) with either of the two variant encoding types across these
391 simulated traits with 50% total heritability and either no or 20% non-linear contribution at given sample sizes (10k, 50k, 100k).

392
393

394 **Elastic net PGS compared to P+T**

395

396 Consistent with previous studies, our results showed that elastic net outperformed the additive
397 P+T method for almost every trait under different settings (**Figure 4, Supplementary Figures**
398 **1-5**). For example, EN consistently outperformed P+T for all the traits (mean R^2 improvement:
399 40%) when a larger sample size (50k and 100k) was used in training; even with the smallest
400 sample size (10k), there were only 10 traits (out of 60) that had a lower R^2 score with EN.

401

402 We further examined how differently EN and P+T methods estimate linear effect sizes of
403 variants in PGS development for traits of different settings. Overall, our results showed the
404 outperformance of EN can be reflected in its better estimation of linear effect sizes of causal
405 variants of the trait (both additively encoded and sample size = 100k; **Supplementary Figures**
406 **9-12**). For example, for purely additive traits, the Spearman correlation (R_s) between the true
407 linear effect size (from simulation) and the EN-estimated effect size was about 25% higher on
408 average when compared with using P+T method (**Supplementary Figure 10**). When total
409 heritability of the trait decreases, EN maintained its accuracy in estimating variant effect sizes
410 ($R_s = 0.71, 0.73$ and 0.72 for traits with total heritability of 80%, 50% and 20% respectively),
411 but P+T saw a significant performance drop ($R_s = 0.63, 0.56$ and 0.54 for traits with total
412 heritability of 80%, 50% and 20% respectively) (**Supplementary Figure 10**). We also found
413 EN was able to better capture true linear effect sizes of variants for traits that are controlled by
414 very high or low proportions of GxG interactions (**Supplementary Figures 9, 11-12**). For
415 example, the R_s between the true linear effect size and the estimated effect size from EN was
416 43% higher on average for traits with either 20% or 80% of non-linear heritability (80% total
417 heritability; **Supplementary Figures 9, 11**) than additive P+T, but this improvement decreased
418 to 9% for traits with 50% of non-linear heritability (**Supplementary Figure 12**).

419

420

421 **Discussion**

422

423 In this work, we simulated realistic genotype-to-phenotype data with varying key parameters
424 then compared the performance of deep learning using a multi-layer perceptron to two other
425 common methods (i.e. elastic net and P+T, pruning and thresholding). These key parameters
426 included GxG interaction types, total trait heritability, proportions of non-linear heritability
427 (i.e. from GxG), genotype encoding (additive and one hot encoding), and different sample sizes
428 for training data. Our results showed that traits with low GxG heritability were best predicted
429 by EN but when the proportion of GxG increases, one hot encoding allowed EN to outperform
430 MLP. Our results also showed that the MLP performed considerably better with an increased
431 sample size in training, and as the total trait heritability increases, the relative performance of
432 MLP in comparison with the linear PGS methods increased. Our results suggest that as the size
433 of the training dataset increases substantially beyond 100k toward a million or more
434 individuals, neural network models may achieve equal or better performance as linear PGS
435 methods. However, currently the computational and financial expense of training even an MLP
436 to UK Biobank data is out of reach for the vast majority of academic groups. As costs come
437 down, neural network models, such as MLP, could be useful for the prediction of highly
438 heritable, substantially GxG phenotypes (e.g. some autoimmune diseases) in massive-scale
439 biobanks, e.g. those of millions of individuals.

440

441 We found that EN was better at capturing the true linear effect sizes present in the causal
442 variants involved in GxG, indicating EN can better predict traits with GxG even when no
443 interaction terms are explicitly defined in the model. When individual-level genotypes are not

444 available, lasso and related models can be run on GWAS summary statistics using tools such
445 as LDpred^{34,35}, Lassosum³⁶, PRS-CS³⁷ and SBayesRC³⁸. Studies have also shown some traits
446 may benefit from PGS methods (e.g. EN) that are based on individual-level genetic data in the
447 current era of large-scale biobanks such as blood cell traits³⁹, and ensemble PGS methods, that
448 combine both summary-level (or PGS previously developed in external cohorts) and
449 individual-level data, can result in improved PGS for phenotypes such as coronary heart
450 disease⁴⁰. Nonetheless, our findings support the use of and continued access to individual-level
451 data by *bona fide* researchers so that optimal PGS can be constructed using the most advanced
452 methods.

453

454 **Limitations**

455 Whilst the MLP model used in our work is relatively standard and unspecialized, domain
456 knowledge, such as total heritability and GxG interaction type, can be utilised to further
457 optimize neural network architectures. These factors may make neural networks become better
458 predictive models in polygenic prediction of certain traits. Our study utilized simulated
459 phenotypes involving statistical GxG, with a fixed proportion of variance attributed to epistasis.
460 However, in real data analyses, the presence of statistical epistasis can be confounded by LD.
461 In such cases, untyped causal variants may be jointly tagged by SNPs included in the dataset,
462 potentially manifesting as statistical epistasis⁴¹. To differentiate between these joint tagging
463 effects and true epistasis, dedicated methods are required.²² In this study, we only included
464 GxG interaction types that were enumerated by one of the previous studies²⁶, a variation of
465 diverse pairwise or two-loci interactions. However, it is possible that higher orders of
466 interactions, not considered in this study, could be present within the human genome. For
467 instance, higher order of interactions have been reported in genes affecting several non-human
468 traits, such as chicken body weight⁴² or colony morphology in yeast⁴³. If such interactions exist
469 in humans as well, it is conceivable that neural networks or other complex prediction methods
470 would be more favourable in their polygenic prediction. Finally, we could not justify the costs
471 (both financial and carbon emissions) of simulating data and training neural networks to
472 datasets substantially greater than 100k individuals. We believe such an approach may be
473 justifiable for real data for select autoimmune diseases where substantial GxG is likely (e.g.
474 type 1 diabetes); however, given the paucity of autoimmune cases in existing biobanks a
475 concerted effort would be needed to assemble and harmonise individual-level data for neural
476 network training.

477

478 **Conclusion**

479 In summary, this work provides a detailed assessment of neural network models for predicting
480 traits in diverse genetic architectures, in comparison with two commonly used linear PGS
481 methods. It gives general insights into the application of deep learning methods in polygenic
482 prediction, and provides guidance for the selection of optimal PGS methods, variant encoding
483 approach, and training sample size when developing PGS for a target trait. Investigations into
484 customised neural network models, that utilise the genetic architecture of a target trait, may
485 represent a promising future for deep learning in polygenic prediction.

486

487 Carbon impact of this study

488
489 Based in Victoria, Australia, the computational methods used in this study had an estimated
490 carbon footprint of 2,973 kgCO₂, which is equivalent to 3,130 tree months. This was estimated
491 using calculated using green-algorithms.org v1.0⁴⁴.
492

493 Code availability

494
495 The original codes used to simulate phenotypes of various genetic architectures are available
496 at <https://github.com/JasonGrealey/Simulations>. The codes of using the three methods (P+T,
497 EN and MLP) to develop PGS are available at [https://github.com/xuyu-cam/Deep-learning-](https://github.com/xuyu-cam/Deep-learning-for-genetic-prediction-of-complex-traits)
498 [for-genetic-prediction-of-complex-traits](https://github.com/xuyu-cam/Deep-learning-for-genetic-prediction-of-complex-traits).
499

500 Conflicts of Interest

501 M.I. is a trustee of the Public Health Genomics (PHG) Foundation, a member of the Scientific
502 Advisory Board of Open Targets, and has research collaborations with AstraZeneca,
503 Nightingale Health and Pfizer which are unrelated to this study.
504

505 Acknowledgements

506 This study was supported by the Victorian Government's Operational Infrastructure Support
507 (OIS) program. JG was supported by a La Trobe University Postgraduate Research Scholarship
508 jointly funded by the Baker Heart and Diabetes Institute and a La Trobe University Full-Fee
509 Research Scholarship.
510

511 The support of the UK Economic and Social Research Council (ESRC) is gratefully
512 acknowledged (ES/T013192/1). This work was supported by core funding from: the UK
513 Medical Research Council (MR/L003120/1), the British Heart Foundation (RG/13/13/30194;
514 RG/18/13/33946) and the NIHR Cambridge Biomedical Research Centre (BRC-1215-20014).
515 This work was also supported by Health Data Research UK, which is funded by the UK
516 Medical Research Council, Engineering and Physical Sciences Research Council, Economic
517 and Social Research Council, Department of Health and Social Care (England), Chief Scientist
518 Office of the Scottish Government Health and Social Care Directorates, Health and Social Care
519 Research and Development Division (Welsh Government), Public Health Agency (Northern
520 Ireland), British Heart Foundation, and Wellcome. This study was supported by the Victorian
521 Government's Operational Infrastructure Support (OIS) program. The views expressed are
522 those of the authors and not necessarily those of the NHS, the NIHR or the Department of
523 Health and Social Care. M.I. was supported by the Munz Chair of Cardiovascular Prediction
524 and Prevention.
525

526

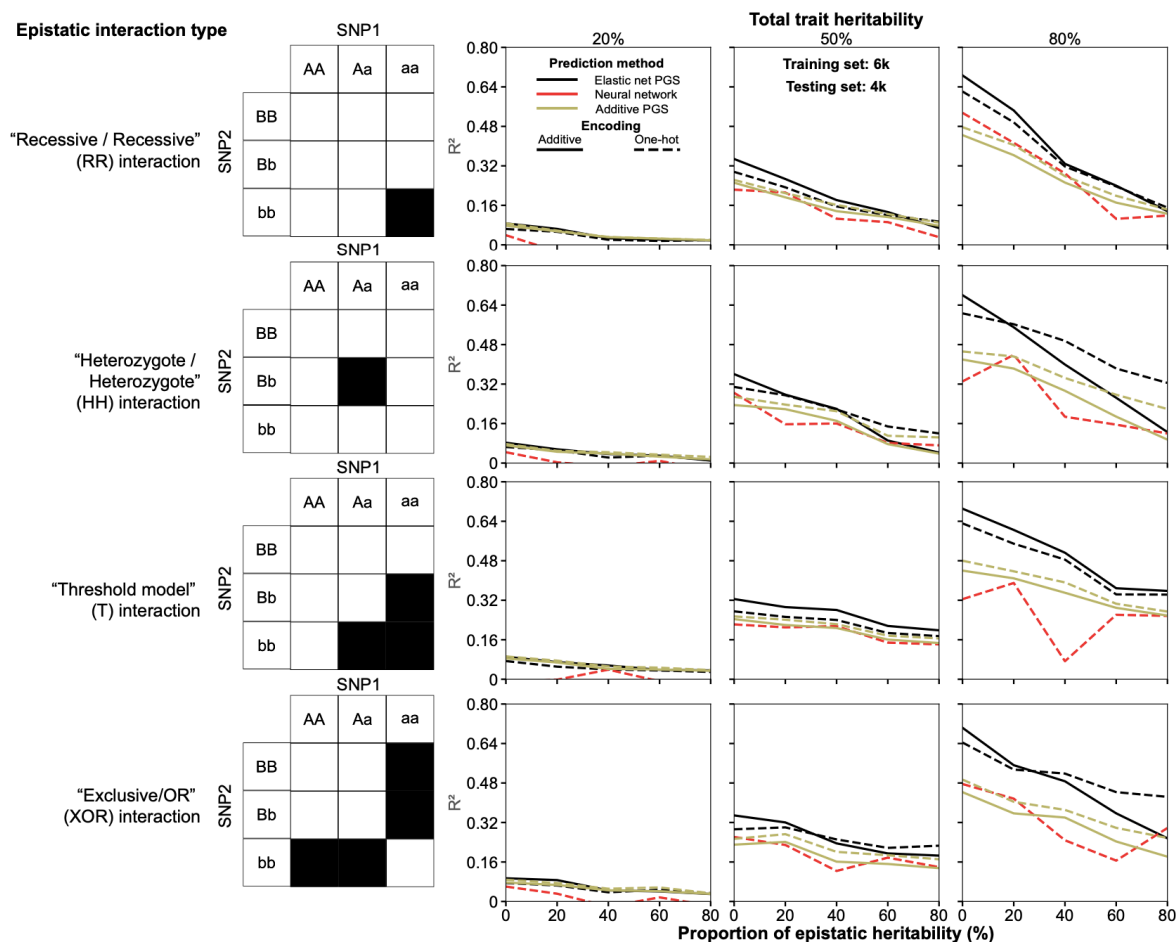
527 References

- 528 1. Lambert, S. A., Abraham, G. & Inouye, M. Towards clinical utility of polygenic risk scores.
529 *Hum. Mol. Genet.* **28**(R2), R133–R142 (2019).
- 530 2. Gibson, G. On the utilization of polygenic risk scores for therapeutic targeting. *PLOS Genet.*
531 **15**, e1008060 (2019).
- 532 3. Ritchie, S. C. *et al.* Integrative analysis of the plasma proteome and polygenic risk of
533 cardiometabolic diseases. *Nat. Metab.* **3**, 1476–1483 (2021).
- 534 4. Bellot, P., de Los Campos, G. & Pérez-Enciso, M. Can Deep Learning Improve Genomic
535 Prediction of Complex Human Traits? *Genetics* **210**, 809–819 (2018).
- 536 5. Xu, Y. *et al.* Machine learning optimized polygenic scores for blood cell traits identify sex-
537 specific trajectories and genetic correlations with disease. *Cell Genomics* **2**, 100086 (2022).
- 538 6. Badré, A., Zhang, L., Muchero, W., Reynolds, J. C. & Pan, C. Deep neural network improves
539 the estimation of polygenic risk scores for breast cancer. *J. Hum. Genet.* **66**, 359–369 (2020).
- 540 7. Sigurdsson, A. I. *et al.* Deep integrative models for large-scale human genomics. *Nucleic
541 Acids Res.* **51**, e67–e67 (2023).
- 542 8. Luo, X., Kang, X. & Schönhuth, A. Predicting the prevalence of complex genetic diseases
543 from individual genotype profiles using capsule networks. *Nat. Mach. Intell.* **5**, 114–125
544 (2023).
- 545 9. Motsinger-Reif, A. A., Dudek, S. M., Hahn, L. W. & Ritchie, M. D. Comparison of
546 approaches for machine-learning optimization of neural networks for detecting gene-gene
547 interactions in genetic epidemiology. *Genet. Epidemiol.* **32**, 325–340 (2008).
- 548 10. Zhou, X. *et al.* Deep learning-based polygenic risk analysis for Alzheimer’s disease prediction.
549 *Commun. Med.* **3**, 1–20 (2023).
- 550 11. Kim, S. bin, Kang, J. H., Cheon, M. J., Kim, D. J. & Lee, B. C. Stacked neural network for
551 predicting polygenic risk score. *Sci. Rep.* **14**, 1–15 (2024).
- 552 12. Hornik, K., Stinchcombe, M. & White, H. Multilayer feedforward networks are universal
553 approximators. *Neural Networks* **2**, 359–366 (1989).
- 554 13. Bengio, Y., Goodfellow, I. J. & Courville, A. *Deep Learning*. (Massachusetts, USA: MIT
555 press, 2017).
- 556 14. Abdellaoui, A., Yengo, L., Verweij, K. J. H. & Visscher, P. M. 15 years of GWAS discovery:
557 Realizing the promise. *Am. J. Hum. Genet.* **110**, 179–194 (2023).
- 558 15. Dudbridge, F. & Wray, N. R. Power and Predictive Accuracy of Polygenic Risk Scores. *PLOS
559 Genet.* **9**, e1003348 (2013).
- 560 16. Hill, W. G., Goddard, M. E. & Visscher, P. M. Data and Theory Point to Mainly Additive
561 Genetic Variance for Complex Traits. *PLOS Genet.* **4**, e1000008 (2008).
- 562 17. Sharp, S. A. *et al.* Development and Standardization of an Improved Type 1 Diabetes Genetic
563 Risk Score for Use in Newborn Screening and Incident Diagnosis. *Diabetes Care* **42**, 200–207
564 (2019).
- 565 18. Privé, F., Aschard, H. & Blum, M. G. B. Efficient Implementation of Penalized Regression for
566 Genetic Risk Prediction. *Genetics* **212**, 65–74 (2019).
- 567 19. Elgart, M. *et al.* Non-linear machine learning models incorporating SNPs and PRS improve
568 polygenic prediction in diverse human populations. *Commun. Biol.* **5**, 1–12 (2022).
- 569 20. Jiajie Peng, A. *et al.* A Deep Learning-based Genome-wide Polygenic Risk Score for Common
570 Diseases Identifies Individuals with Risk. *medRxiv* 2021.11.17.21265352 (2021)
571 doi:10.1101/2021.11.17.21265352.
- 572 21. van Hilten, A. *et al.* GenNet framework: interpretable deep learning for predicting phenotypes
573 from genetic data. *Commun. Biol.* 2021 41 **4**, 1–9 (2021).
- 574 22. Kelemen, M. *et al.* Performance of deep-learning based approaches to improve polygenic
575 scores. *medRxiv* (2024).
- 576 23. Su, Z., Marchini, J. & Donnelly, P. HAPGEN2: simulation of multiple disease SNPs.
577 *Bioinformatics* **27**, 2304–2305 (2011).
- 578 24. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68 (2015).
- 579 25. Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based

- 580 Linkage Analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 581 26. Li, W. & Reich, J. A Complete Enumeration and Classification of Two-Locus Disease
582 Models. *Hum. Hered.* **50**, 334–349 (2000).
- 583 27. Meyer, H. V. & Birney, E. PhenotypeSimulator: A comprehensive framework for simulating
584 multi-trait, multi-locus genotype to phenotype relationships. *Bioinformatics* **34**, 2951–2956
585 (2018).
- 586 28. Pedregos, F. et al. Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**, 2825–
587 2830 (2011).
- 588 29. Privé, F., Vilhjálmsson, B. J., Aschard, H. & Blum, M. G. B. Making the Most of Clumping
589 and Thresholding for Polygenic Scores. *Am. J. Hum. Genet.* **105**, 1213–1221 (2019).
- 590 30. Choi, S. W. & O'Reilly, P. F. PRSice-2: Polygenic Risk Score software for biobank-scale data.
591 *Gigascience* **8**, 1–6 (2019).
- 592 31. Abraham, G., Kowalczyk, A., Zobel, J. & Inouye, M. SparSNP: Fast and memory-efficient
593 analysis of all SNPs for phenotype prediction. *BMC Bioinformatics* **13**, 1–8 (2012).
- 594 32. Abadi, M. et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed
595 Systems. (2016) doi:10.48550/arxiv.1603.04467.
- 596 33. Retrieved from <http://github.com/autonomio/talos>. Autonomio Talos [Computer software].
597 (2020).
- 598 34. Privé, F., Arbel, J. & Vilhjálmsson, B. J. LDpred2: better, faster, stronger. *Bioinformatics* **36**,
599 5424–5431 (2021).
- 600 35. Vilhjálmsson, B. J. et al. Modeling Linkage Disequilibrium Increases Accuracy of Polygenic
601 Risk Scores. *Am. J. Hum. Genet.* **97**, 576–592 (2015).
- 602 36. Mak, T. S. H., Porsch, R. M., Choi, S. W., Zhou, X. & Sham, P. C. Polygenic scores via
603 penalized regression on summary statistics. *Genet. Epidemiol.* **41**, 469–480 (2017).
- 604 37. Ge, T., Chen, C. Y., Ni, Y., Feng, Y. C. A. & Smoller, J. W. Polygenic prediction via Bayesian
605 regression and continuous shrinkage priors. *Nat. Commun.* **10**, 1–10 (2019).
- 606 38. Zheng, Z. et al. Leveraging functional genomic annotations and genome coverage to improve
607 polygenic prediction of complex traits within and between ancestries. *Nat. Genet.* **56**, 767–777
608 (2024).
- 609 39. Plagnol, V. Polygenic score development in the era of large-scale biobanks. *Cell Genomics* **2**,
610 100088 (2022).
- 611 40. Inouye, M. et al. Genomic Risk Prediction of Coronary Artery Disease in 480,000 Adults:
612 Implications for Primary Prevention. *J. Am. Coll. Cardiol.* **72**, 1883–1893 (2018).
- 613 41. Wood, A. R. et al. Another explanation for apparent epistasis. *Nature* **514**, E3–E5 (2014).
- 614 42. Pettersson, M., Besnier, F., Siegel, P. B. & Carlborg, Ö. Replication and Explorations of High-
615 Order Epistasis Using a Large Advanced Intercross Line Pedigree. *PLOS Genet.* **7**, e1002180
616 (2011).
- 617 43. Taylor, M. B. & Ehrenreich, I. M. Genetic Interactions Involving Five or More Genes
618 Contribute to a Complex Trait in Yeast. *PLOS Genet.* **10**, e1004324 (2014).
- 619 44. Lannelongue, L., Grealey, J. & Inouye, M. Green Algorithms: Quantifying the Carbon
620 Footprint of Computation. *Adv. Sci.* **8**, 2100707 (2021).

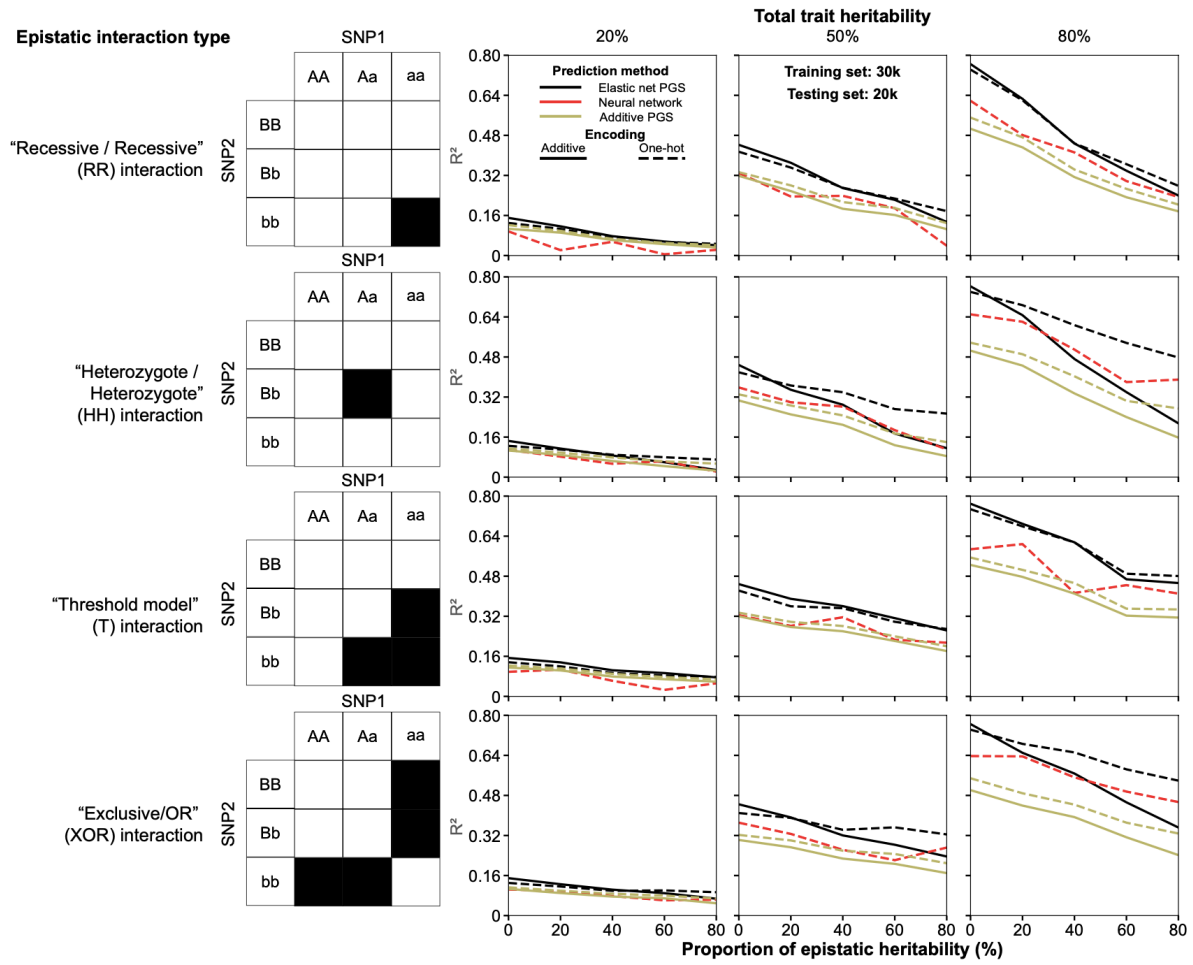
621
622

623 Supplementary Figures



624
 625 **Supplementary figure 1. Predictive performances (R^2) of all PGS methods with a training size of 6k and testing set of**
 626 **4k samples.** Elastic net PGS method is in black, neural network MLP in red and additive P+T PGS in yellow. Solid lines and
 627 dashed lines represent additively and one hot encoded models respectively. The plots detail each PGS Method's performance
 628 for a given phenotype. Each row has a different underlying interaction model labelled by the tables; each column has the same
 629 total heritability as noted by the column title and within each plot. The x axis details the proportion of epistatic contribution
 630 for a given trait.

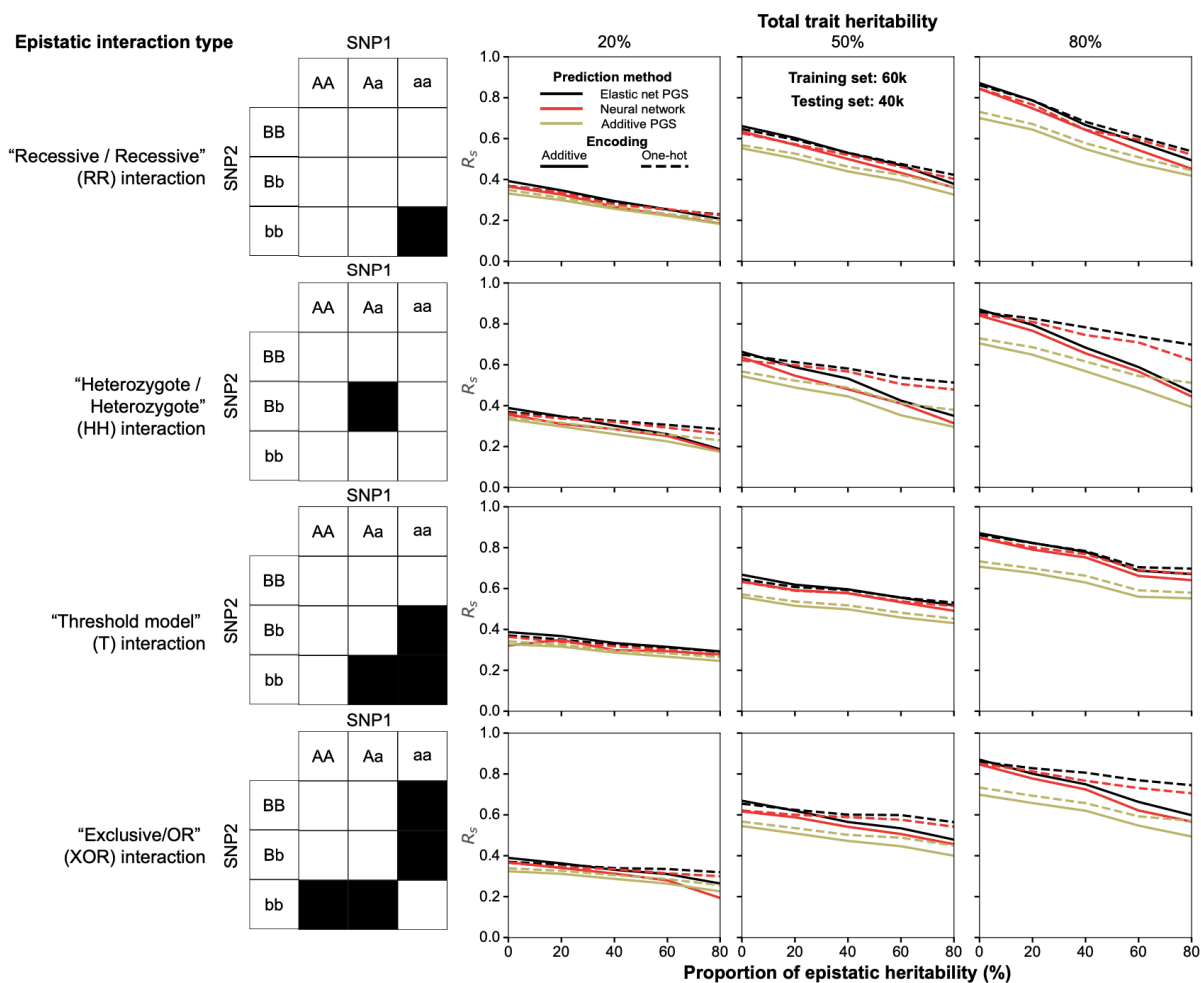
631



632

633 **Supplementary figure 2. Predictive performances (R^2) of all PGS methods with a training size of 30k and testing set**
 634 **of 20k samples.**

635



636

637

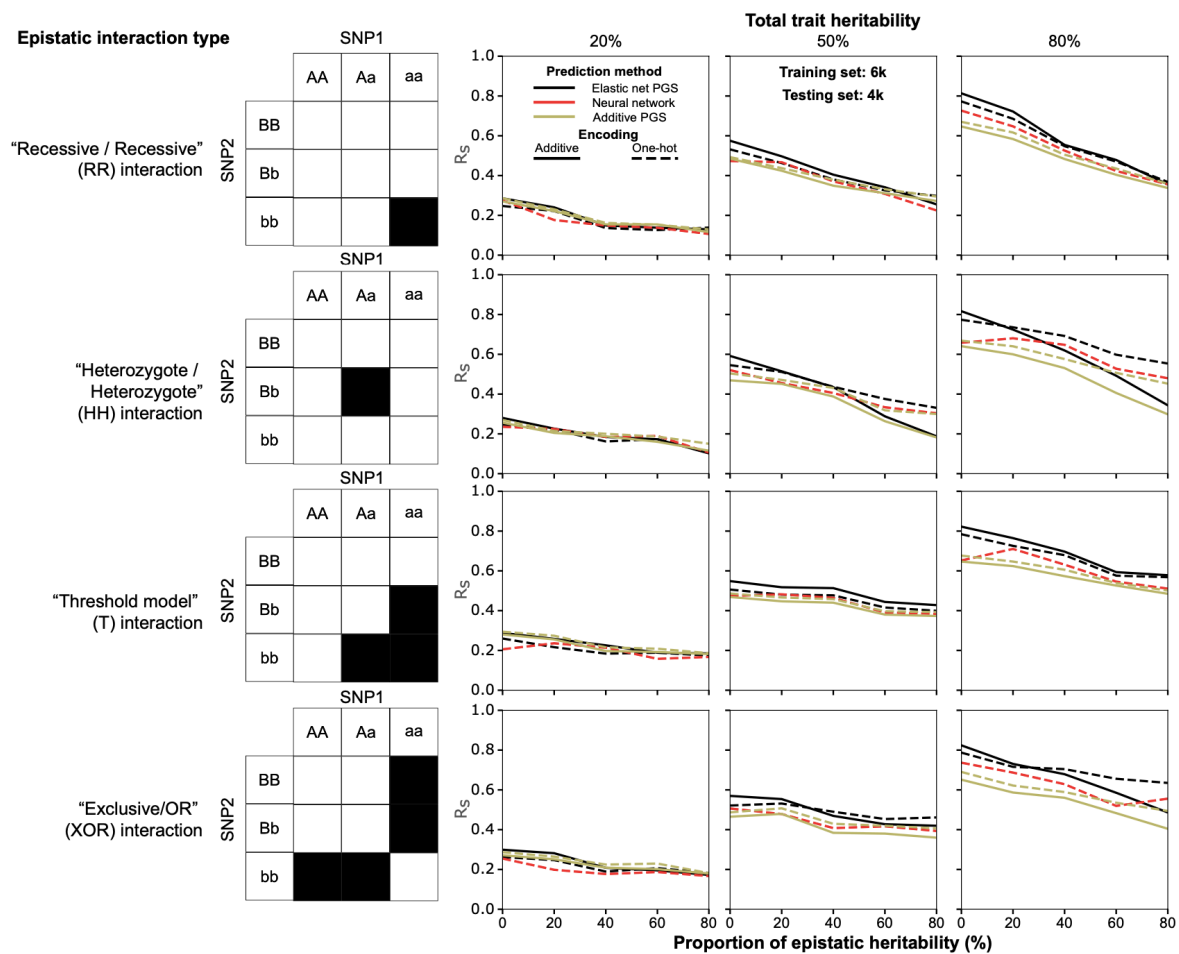
638

Supplementary figure 3. Predictive performances (Spearman R_s) of all PGS methods with a training size of 60k and testing set of 40k samples.

639

640

641



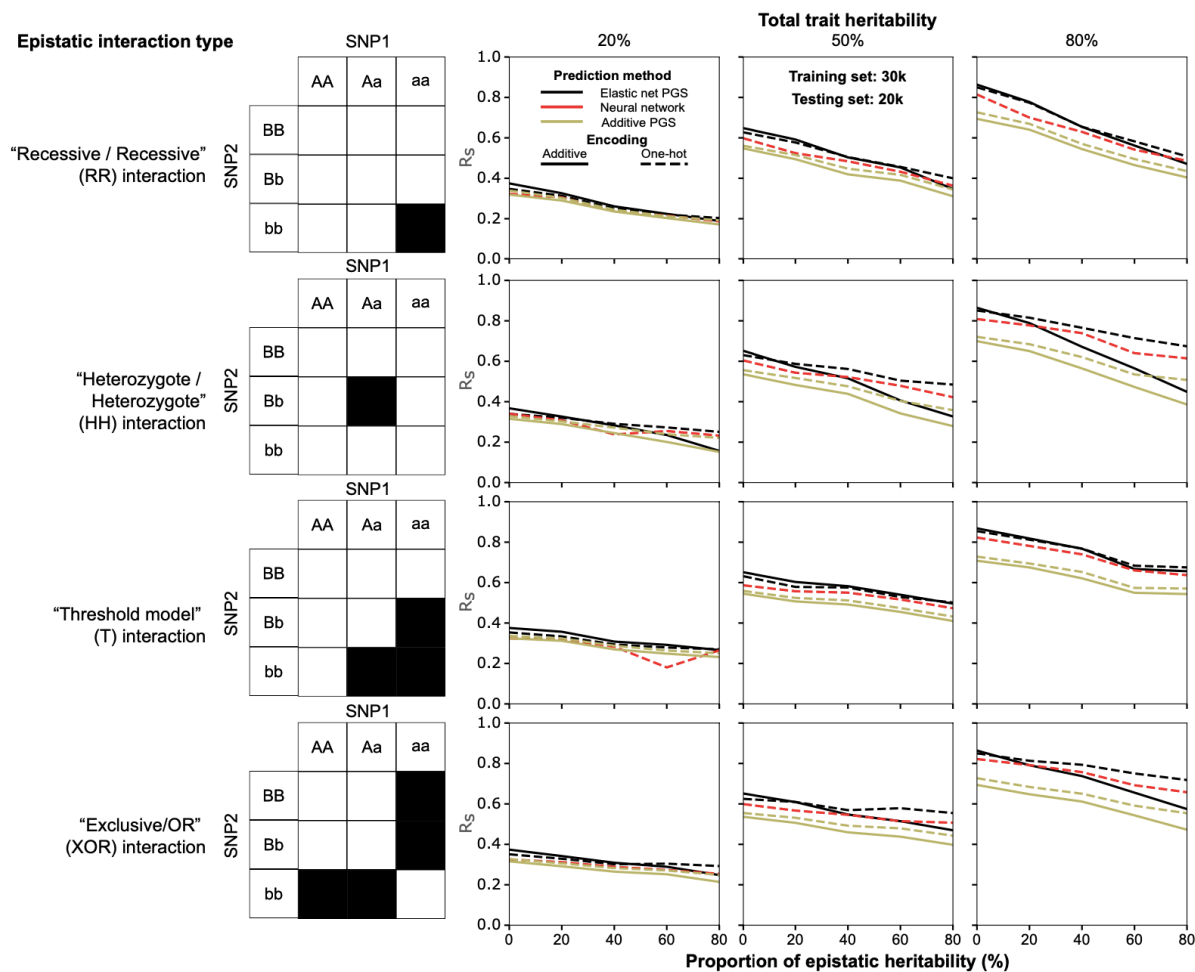
642

643 **Supplementary figure 4. Predictive performances (Spearman R_s) of all PGS methods with a training size of 6k and**
 644 **testing set of 4k samples.**

645

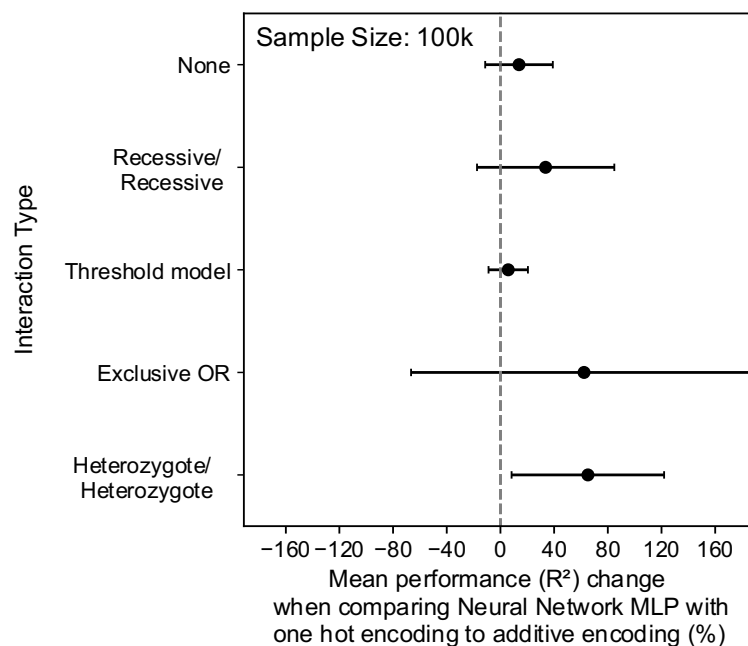
646

647



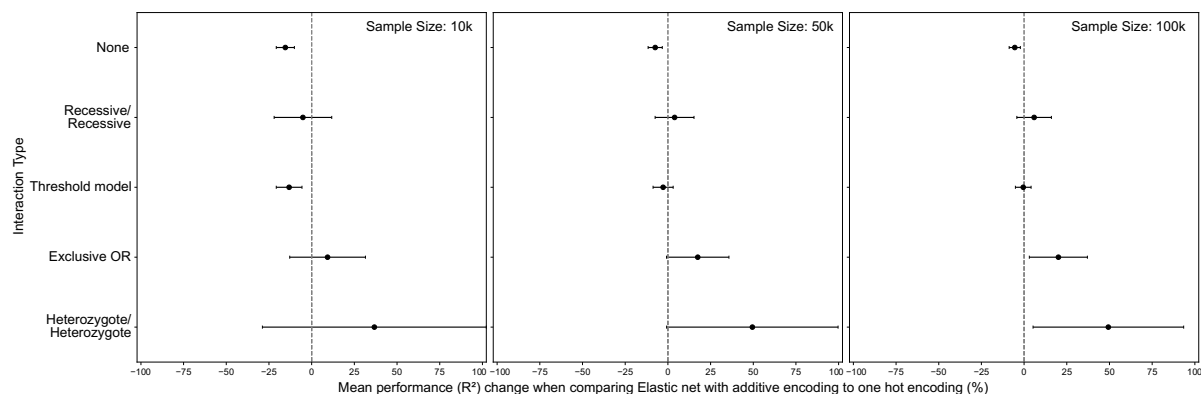
648
649 **Supplementary figure 5. Predictive performances (Spearman R_s) of all PGS methods with a training size of 30k and**
650 **testing set of 20k samples.**

651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669



670
671 **Supplementary figure 5. Neural network MLP predicts more accurately with one hot encoded variants.** The plot details
672 the percentage change in the predictive performance (R²) of neural network MLP method with one hot encoding over that of
673 using additive encoding at the total sample size of 100k, which are measured using the mean and standard deviation of
674 $(R_{one-hot}^2 - R_{additive}^2)/R_{additive}^2$ ($R_{one-hot}^2$ and $R_{additive}^2$ are R² performance of MLP using one hot encoding and additive
675 encoding respectively for a trait) in each trait group by GxG interaction type or with no interactions.

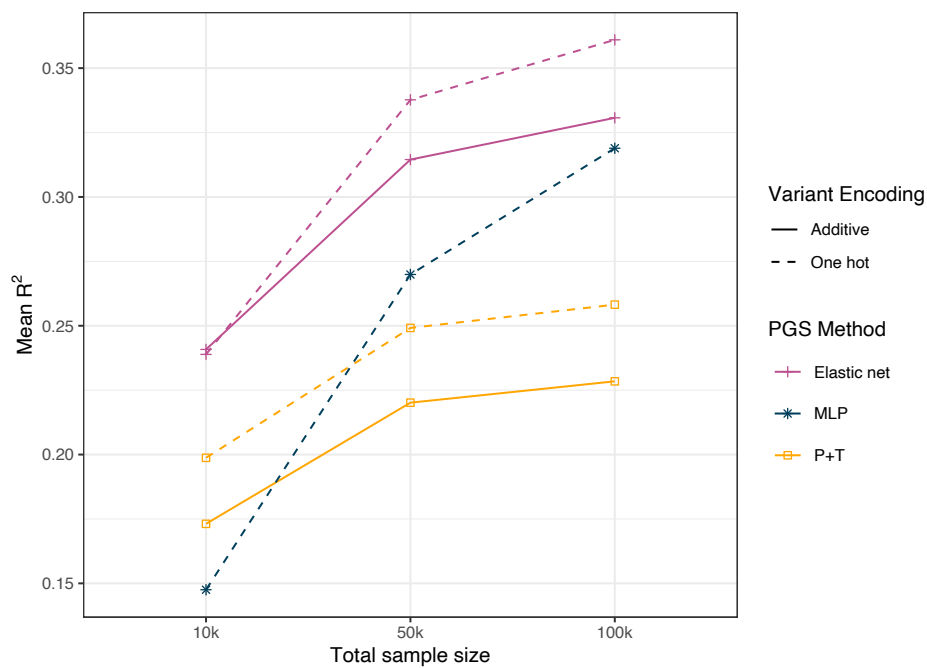
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695



696
697
698
699

Supplementary figure 7. One hot encoding allows elastic net to better predict non-linear traits. The plot details the percentage change in the predictive performance (R^2) when comparing additively encoded elastic net to one hot encoded elastic net by groups of traits with different interaction types in each of the sample sizes used in this study.

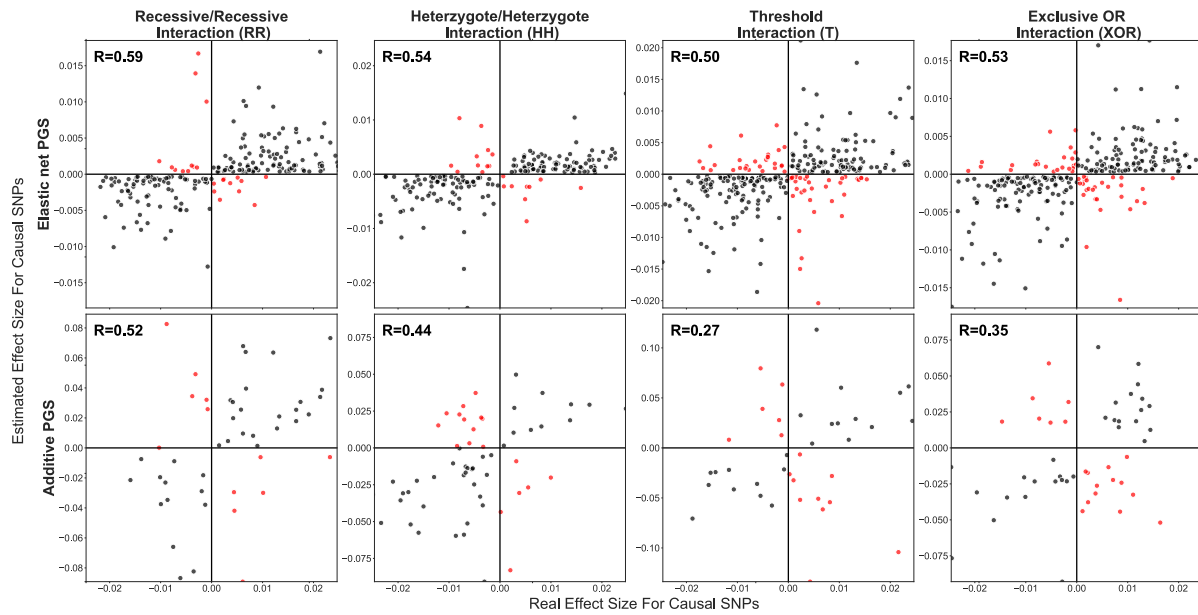
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732



733

734 **Supplementary figure 8. The R^2 performance improvements of PGS methods by sample size increase.** This plot shows
735 the mean R^2 performance of each PGS method (P+T, elastic net and MLP) with either of the two variant encoding types across
736 all the 60 simulated traits at given sample sizes (10k, 50k, 100k).

737

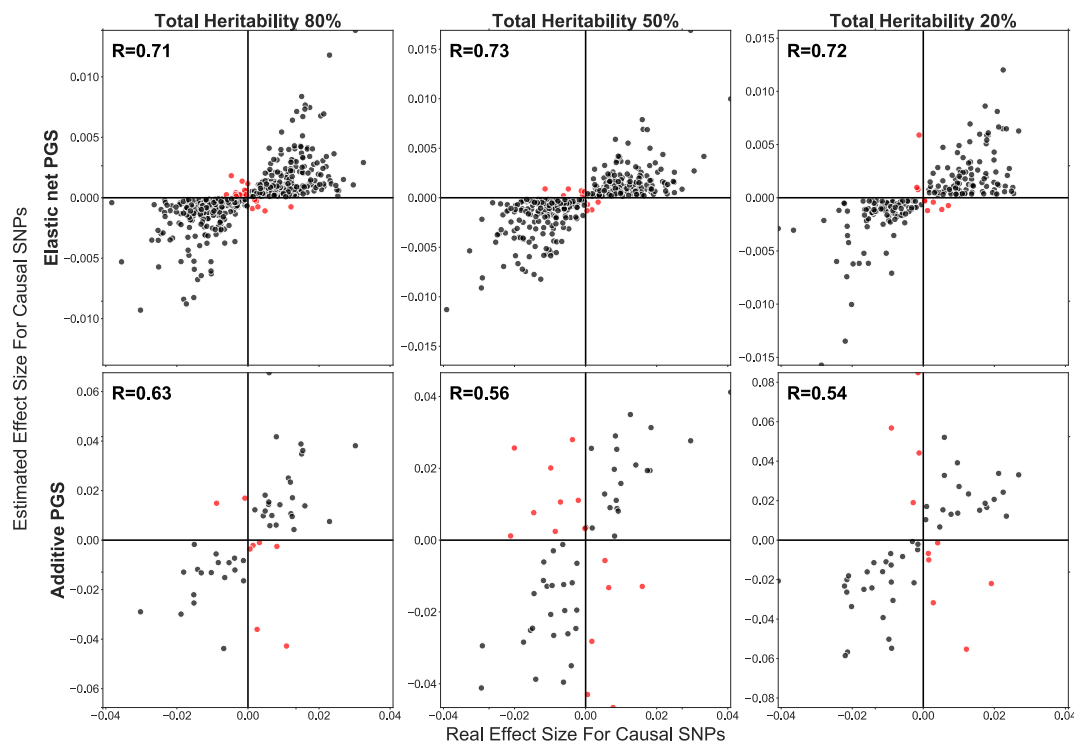


738

739 **Supplementary figure 9. Elastic-net better estimates effect sizes in highly non-additive traits than additive P+T method**
740 **(80% total heritability and 80% GxG).** Each sub-plot corresponds to a single trait under control of different types of
741 interactions, in which the trait has the total heritability of 80% and 80% of heritability is explained by GxG. The total sample
742 size used is 100k. The *x*-axis presents the real linear effects of causal variants and the estimated effect sizes are displayed on
743 the *y* axis. The upper row shows non-zero effect sizes estimated by Elastic net for causal variants; bottom row shows effect
744 sizes of causal variants estimated by P+T. Columns are separated by the GxG interaction type present in the trait (i.e. XOR,
745 RR, HH and T). Points in the off diagonal are coloured in red. The Spearman correlation between the two effect sizes across
746 all the variants in each plot is labelled at the top left. Note that for clarity any effect size estimated to be exactly zero is removed
747 from the plot.

748

749



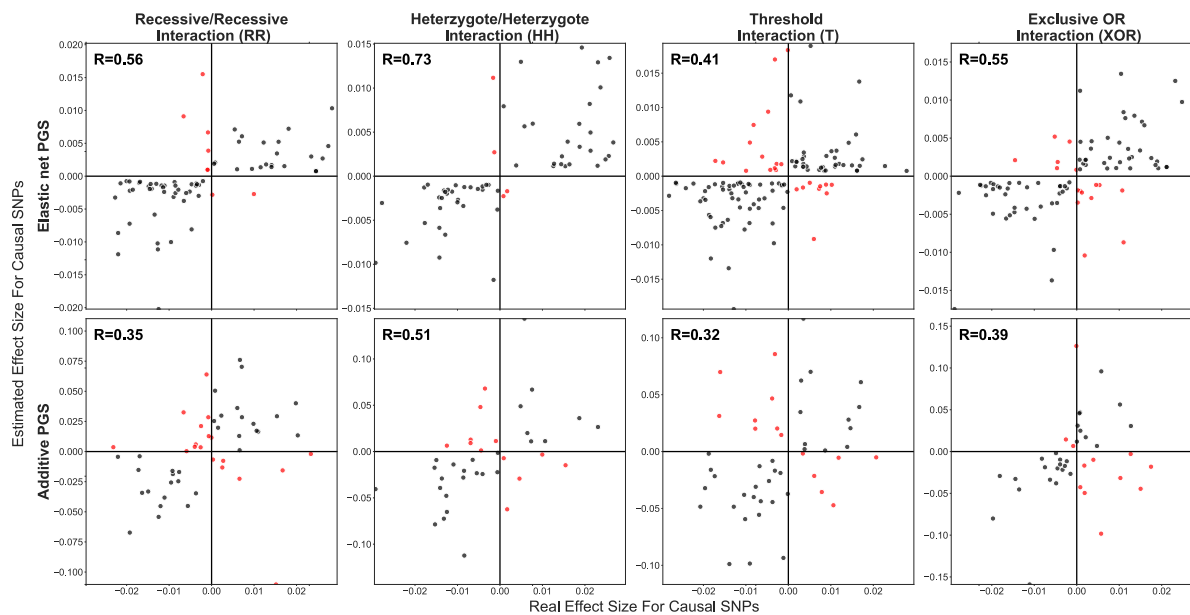
750

751 **Supplementary figure 10. Elastic-net better learns effect sizes of variants in linear traits than additive P+T.** Each sub-
752 plot corresponds to a single trait with no interactions (i.e. purely additive traits). The total sample size used is 100k. The x-
753 axis presents the real linear effects for each causal variant (from estimation) and the y axis shows the estimated effect sizes
754 using EN or P+T. The upper row shows non-zero effect sizes of causal variants estimated by Elastic net for traits with 80%,
755 50%, and 20% total heritability from left to right respectively; bottom row shows effect sizes of causal variants estimated by
756 P+T. Points in the off diagonal are coloured in red. The Spearman correlation between the two effect sizes across all the
757 variants in each plot is labelled at the top left. Note that for clarity any effect sizes estimated to be exactly zero are removed
758 from the plot.

759

760

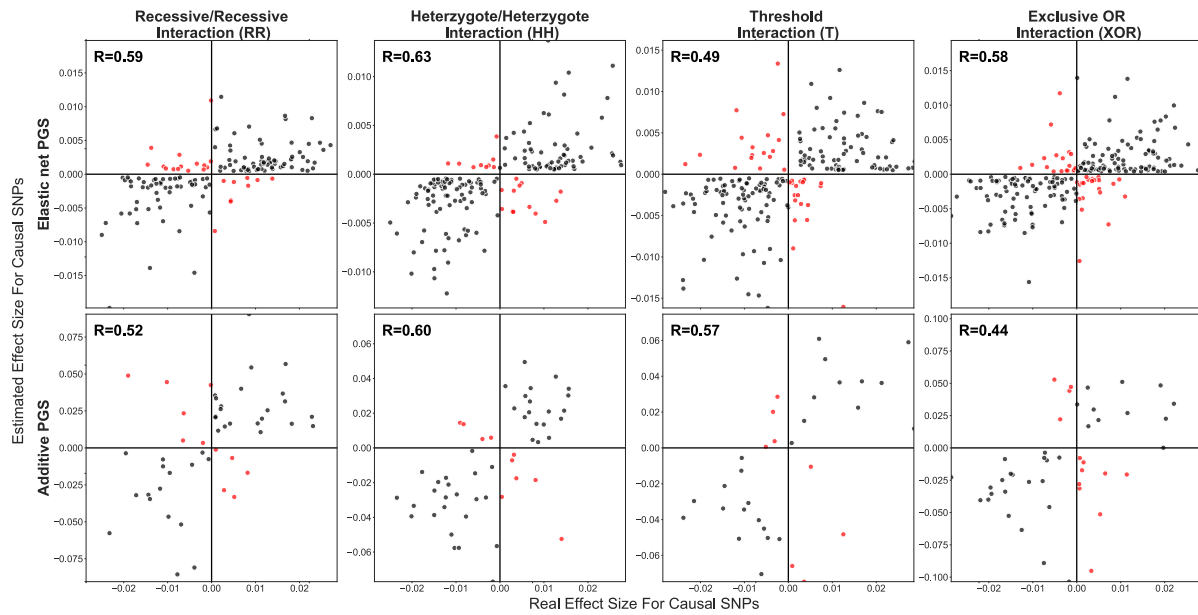
761



762
763
764
765
766
767
768
769
770

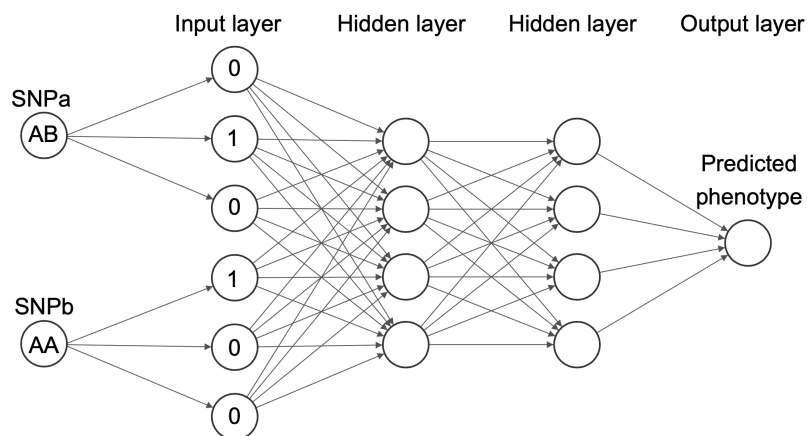
Supplementary figure 11. Elastic-net better estimates effect sizes in highly non-additive traits than additive P+T (80% total heritability and 20% GxG). Each sub-plot corresponds to a single trait under control of different types of interactions, in which the trait has the total heritability of 80% and 20% of heritability is explained by GxG. The total sample size used is 100k. The x-axis presents the real linear effects of causal variants and the estimated effect sizes are displayed on the y-axis. The upper row shows non-zero effect sizes estimated by Elastic net for causal variants; bottom row shows effect sizes of causal variants estimated by P+T. Columns are separated by the interaction type present in the trait (i.e. XOR, RR, HH and T). Points in the off diagonal are coloured in red. The Spearman correlation between the two effect sizes across all the variants in each plot is labelled at the top left. Note that for clarity any effect size estimated to be exactly zero is removed from the plot.

771
772
773
774
775
776
777
778
779
780
781
782
783
784
785



786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811

Supplementary figure 12. Elastic-net better estimates effect sizes in highly non-additive traits than additive P+T (80% total heritability and 50% GxG interaction). Each sub-plot corresponds to a single trait under control of different types of interactions, in which the trait has the total heritability of 80% and 50% of heritability is explained by GxG. The total sample size used is 100k. The x-axis presents the real linear effects of causal variants and the estimated effect sizes are displayed on the y-axis. The upper row shows non-zero effect sizes estimated by Elastic net for causal variants; bottom row shows effect sizes of causal variants estimated by P+T. Columns are separated by the interaction type present in the trait (i.e. XOR, RR, HH and T). Points in the off diagonal are coloured in red. The Spearman correlation between the two effect sizes across all the variants in each plot is labelled at the top left. Note that for clarity any effect size estimated to be exactly zero is removed from the plot.

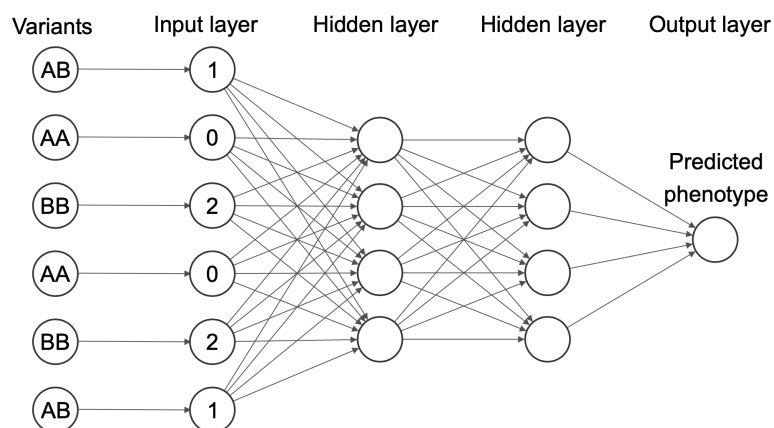


812

813 **Supplementary figure 13. One hot encoded Neural network (NN) schematic.** The above is an example of the MLP network
814 structure used in this study. This NN has two input variants, two hidden layers and an output layer where the phenotype is
815 predicted. The input variants are firstly encoded into their genotype classes through one hot encoding, then these are passed
816 through the network's hidden layers and finally the phenotype is predicted. The NNs in this study had 100,455 variants as
817 input.

818

819



820

821 **Supplementary figure 14. Additively encoded Neural network (NN) schematic.** The plot is an example of the MLP network
822 structure used in this study. This NN has six input variants, two hidden layers and an output layer where the phenotype is
823 predicted. The input variants are encoded by counting the number of affect alleles, then these are passed through the network's
824 hidden layers and finally the phenotype is predicted. The NNs in this study had 100,455 variants as input.

825

826