# Leveraging hierarchical structures for genetic block interaction studies using the hierarchical transformer

Shiying Li[1], Shivam Arora[2], Redha Attaoua[1], Pavel Hamet[1], Johanne Tremblay[1], Alexander Bihlo[2], Bang Liu[3] and Guy A. Rutter[1,4,5*]

[1]Centre de Recherche du CHUM, and Faculty of Medicine, University of Montreal, QC, Canada

[2]Department of Mathematics and Statistics, Memorial University of Newfoundland, NL, Canada

[3]Département d'informatique et de recherche opérationnelle, Université de Montréal, QC, Canada

[4]Section of Cell Biology and Functional Genomics, Department of Metabolism, Diabetes and Reproduction, Imperial College of London, du Cane Road, London W120NN, United Kingdom

[5]Lee Kong Chian School of Medicine, Nan Yang Technological University, Singapore.

Address correspondence to:

Guy A. Rutter guy.rutter@umontreal.ca, or g.rutter@imperoial.ac.uk 514 890-8000, ext. 27081

Word count: 5558    Number of Figures: 5                    Number of Tables: 3

24    1. Abstract

25    Initially introduced in 1909 by William Bateson, classic epistasis (genetic variant

26    interaction) refers to the phenomenon that one variant prevents another variant from a

27    different locus from manifesting its effects. The potential effects of genetic variant

28    interactions on complex diseases have been recognized for the past decades.

29    Moreover, It has been studied and demonstrated that leveraging the combined SNP

30    effects within the genetic block can significantly increase calculation power, reducing

31    background noise, ultimately leading to novel epistasis discovery that the single SNP

32    statistical epistasis study might overlook. However, it is still an open question how we

33    can best combine gene structure representation modelling and interaction learning into

34    an end-to-end model for gene interaction searching. Here, in the current study, we

35    developed a neural genetic block interaction searching model that can effectively

36    process large SNP chip inputs and output the potential genetic block interaction

37    heatmap. Our model augments a previously published hierarchical transformer

38    architecture (Liu and Lapata, 2019) with the ability to model genetic blocks.  The

39    cross-block relationship mapping was achieved via a hierarchical attention mechanism

40    which allows the sharing of information regarding specific phenotypes, as opposed to

41    simple unsupervised dimensionality reduction methods e.g. PCA. Results on both

42    simulation and UK Biobank studies show our model brings substantial improvements

43    compared to traditional exhaustive searching and neural network methods.

44

45    2. Introduction

46    In the past decades, the genetic factors that contribute to the pathogenesis of complex

47    diseases have been extensively studied, and genome-wide association study (GWAS)

48    has played a fundamental role in unveiling novel risk alleles. The approach of the

49    GWAS framework assumes each SNP has an independent effect on phenotype and the

50    disease's statistical relevance was tested individually (Niel et al., 2015). However,

51    most complex disease variants identified so far confer relatively small increments in

52    risk, leading to many questions about how the remaining "missing heritability" can be

53    explained (Maher, 2008). For example, GWAS correlation identified >80 common

54    variants for type 2 diabetes with most of those associated with insulin secretion,

55    together, they only contribute ~ 10% of type 2 diabetes (T2D) heritability. T2D Low-

56    frequency and rare variants have also been identified, but their contribution towards

57    "missing heritability" is also limited (Stančáková and Laakso, 2016).  Indeed, with a

58    larger sample size and the advancement of sequencing techniques, GWAS will likely

59    continue to expand the number of novel complex disease genetic markers. However,

60    the current consensus underscores the growing recognition that the missing heritability

61    of complex diseases extends beyond the scope of singular genetic factors. Interactions

62    among two or more SNPs, a combinatorial effect known as epistasis, have been

63    proven (Turton et al., 2011) can at least partly explain the "missing heritability".

64

65    Finding the optimal interacting SNP combination for certain phenotypes, which

66    implies an exhaustive search of all possible cases, can be a challenging task. For

67    instance, in a dataset containing 500,000 SNPs, there are approximately 250 billion

68    possible pairwise SNP combinations. This immense number presents significant

69 challenges, not only in terms of computational hardware requirements but also in the

70 risk of losing true signals due to overcorrection for multiple comparisons resulting in a

71 reduction in statistical power. A problem we often refer to as the "curse of

72 dimensionality". To address this issue, several scalable statistical approaches have

73 been proposed in recent years; however, each comes with its own set of limitations.

74 Some methods (Cordell, 2002) only select "top SNPs" (the SNPs most correlated with

75 phenotypes) for epistasis searching, while ignoring the potential effects of

76 neighbouring SNPs. Indeed, over the past few decades, extensive research (Morris and

77 Kaplan, 2002; Zaykin et al., 2002; Chen et al., 2020) has highlighted the necessity and

78 efficiency of leveraging the combined effects of multiple SNPs within certain genetic

79 regions e.g. haplotype blocks, rather than focusing on individual SNPs in association

80 studies. Some have proposed to summarise multi-dimensional SNPs into one-

81 dimensional representations using unsupervised methods such as PCA (Li et al.,

82 2009). However, these methods overlook important phenotype information and

83 compress highly dependent SNPs into a single dimension, making it difficult to detect

84 signals within these units. In short, it is still an open question how we can best

85 combine genetic block representation learning and interaction modelling to an end-to-

86 end model to increase calculation power.

87

88 The rapid development of deep learning and artificial intelligence seems to hold

89 another promise for epistasis studies. Several studies (Pérez-Enciso and Zingaretti

90 2019; Cui et al., 2022) have suggested Deep Neural Networks (DNNs) can map the

91 flexible, both linear and non-linear relationships between SNPs and observed

92   phenotypes. However, DNNs, which are primarily designed for classification, with a

93   black-box nature that makes it challenging to interpret results, particularly in

94   identifying which SNPs are interacting. Many studies aim to bridge this gap between

95   interpretability and DNN structures using such as layerwise relevance (Mieth et al.,

96   2021), permutation testing (Cui et al., 2022), and more recently, transformer with

97   attention scores (Graça et al., 2024). While the potential application of transformer in

98   genome sequencing analysis has enjoyed renewed interest. Scanning through most of

99   the genetic transformer studies in recent years (Jubair et al., 2021; Reyes et al., 2022;

100  Zhou et al., 2022), the basic unit, as the "word" in natural language, is still single SNP

101  with little existing genetic structure e.g. haplotype block, the main focus of current

102  study, introduced within the attention block. The size of haplotype blocks can vary

103  and is often larger than the typical units analysed in natural language processing. This

104  variability poses a substantial challenge in encoding these differently-sized haplotype

105  blocks into the transformer encoder while preserving their biological significance

106  during this process, to ensure that the outputs of the attention scores are interpretable

107  and relevant.

108

109  In the current study, taking inspiration from the hierarchical transformer model (Liu

110  and Lapata, 2019), we proposed a novel haplotype block-haplotype block association

111  study workflow, Haplotype Block LSTM hierarchical Transformer (HB-LT). HB-LT

112  is constructed in a hierarchical manner which allows it to efficiently capture both

113  within and cross-haplotype relationships relevant to specific phenotypes. We

114  demonstrate with simulations that grouping SNPs into dimensionality-reduced

115 haplotype block structures significantly increases detection power for epistasis studies

116 compared to existing methods. Furthermore, by evaluating our model on the UK

117 Biobank dataset, we demonstrated its potential for real-world applications.

118

119 3. Methods

120 **3.1 Model description.**

121 The model in the current study is mainly inspired by Liu and Lapata (2019) and

122 several previous machine learning works (Chang et al., 2020; Cui et al., 2022 and

123 Graça et al., 2023) dedicated to epistasis studies. Our haplotype epistasis study system

124 is illustrated in Figure 1. The inputs of the model are the pre-organised haplotype

125 datasets and the associated phenotype of each individual, while the outputs are the

126 attention weights (epistasis) among potential candidates.

127

128 **3.1.1 Long short-term memory (LSTM) pre-selection.** For large dataset analysis,

129 applying pre-selection methods effectively reduces computational burdens and

130 enhances calculation efficiency. However, one of the key challenges in the current

131 study is the variation in haplotype block lengths, which can range from as few as 2 to

132 more than 100 SNPs. Here, we adopted a learning-based approach. A linear regression

133 model is applied to each haplotype individually, and its average root of mean square

134 error on the testing dataset is used as a score indicating whether it should be selected

135 as a phenotype-associated candidate. Haplotype blocks of the SNP dataset were first

136 constructed using the confidence interval method (Gabriel et al., 2002), we then use

137 recurrent neural network LSTM to represent each haplotype block. Let

138
$$\{u_{h1}, \cdots, u_{hi}\} = \mathrm{lstm}(\{w_{h1}, \cdots, w_{hi}\})$$

139 $w_h i$ are word embedding for tokens in each haplotype block, where $u_h i$ are updated

140 vectors for the token after LSTM.

141 An average-pooling is then used to obtain a fixed length representation and a linear

142 transformation yields the final representation of the haplotype block $s_i$.

143
$$s_i = \text{linear}\left(W \cdot \text{avgpool}\left(u_{h1}, \ldots, u_{hi}\right)\right)$$

144 All input haplotype blocks were pre-split into training and testing datasets. The model

145 is trained by minimising the root of mean square error of $s_i$ and the phenotype $y$. In

146 testing, the phenotype-associated haplotype block candidates were selected based on

147 the mean prediction score.

148

149 **3.1.2 Hierarchical transformer encoder.**

150 **3.1.2.1 Embedding.** Input SNPs are first represented by word embeddings. Let

151 $e_i \in \mathbb{R}^{d_e}$ represent the embedded dimensional vectors of the SNP $i$. Let $H$ denote the

152 haplotype where $H = \{e_i\}_{i=1}^{n}$. $n$ is the total number of SNPs in each haplotype block.

153 In our hierarchical haplotype transformer, each token (SNP) has two positions that

154 need to be considered, namely $i$, the position of the token (SNP) within the haplotype,

155 and $j$, the position of the haplotype block within the input sequence. We follow

156 (Vaswani et al., 2017) and use sine and cosine functions for calculating positional

157 embedding. These two positional embedding vectors were then concatenated and the

158 final input vector of each token (SNP) for the hierarchical haplotype transformer

159 model is: $x_i = e_i + e_p$.

160

161 **3.1.2.2 Local haplotype attention block.** The main aim of the local haplotype

162 attention block is to map the dynamic attention scores among SNPs within each

163     haplotype block. It contains several components, including multi-head attention, layer

164     normalisation and feed-forward. The number of these local attention blocks to be used

165     in the model will be decided by the researchers themselves. Let

166     $F = \{f_i\}_{i=1}^n, \quad f_i \in \mathbb{R}^{d_{model}}$ denote the features of SNPs within each haplotype input to

167     the local haplotype attention block. For the $i$ attention head, the query $Q^j = \{q_i^j\}_{i=1}^n$,

168     key $K^j = \{k_i^j\}_{i=1}^n$ and value $V^j = \{v_i^j\}_{i=1}^n$ of each SNP are calculated based on

169     $q_i^j = W_Q^j f_i$, $k_i^j = W_K^j f_i$ and $v_i^j = W_V^j f_i$ respectively. The linear projection learnable

170     parameters weight $W$ are matrices $W_i^Q \in \mathbb{R}^{d_{model} \times d_k}, W_i^K \in \mathbb{R}^{d_{model} \times d_k}, W_i^V \in \mathbb{R}^{d_{model} \times d_v}$.

171     $d_k = d_v = d_{model}/h$. The output of $j$ attention head will be: $\text{softmax}\left(\dfrac{QK^T}{\sqrt{d_k}}\right)V$. The

172     multi-head results were then concatenated and linear transferred with learnable weight

173     $W^O \in \mathbb{R}^{hd_v \times d_{model}}$ to get the final results. The feed-forward layer is composed of two

174     fully connected (FC) layers in inverse order with an activation function *Tahn* in

175     between.

176

177     **3.1.2.3 Inter-haplotype attention block.** To exchange information across different

178     haplotypes, an inter-haplotype attention block was used. To obtain a fixed-length

179     haplotype representation, a weighted, multi-head pooling was first used to represent

180     each haplotype. In each head, weight distributions over tokens (SNPs) are calculated

181     and different heads will encode haplotypes with different attention weights. Let

182     $x_i^{l-1} \in \mathbb{R}^d$ denote the output vector from the last layer of the local haplotype attention

183     block, which will be the input of the multi-head pooling layer. For haplotype $H_j$, for

184     head $z$, a linear transformation was first applied to convert the input vector into an

185     attention score $a_i^z$ and a value vector $b_i^z$ with the weight $W_a^z \in \mathbb{R}^{1 \times d}$ and $W_b^z \in \mathbb{R}^{d_{head} \times d}$.

186    The final output $d_{head} \times 1$ is a weighted sum representing haplotype $j$ in head $z$ where

187    $W_c^z \in \mathbb{R}^{d_{head} \times d_{head}}$.

$$\text{head}_j^z = \text{LayerNorm}\left( W_c^z \sum_{i=1}^{n} a_i^z b_i^z \right)$$

188

189    Similar to local haplotype attention blocks, inter-haplotype allows one haplotype to

190    attend to another to model the haplotype-haplotype dependencies.

191    $$q_j^z \& = W_q^z \text{head}_j^z$$

192    $$k_j^z = W_k^z \text{head}_j^z$$

193    $$v_j^z = W_v^z \text{head}_j^z$$

194    $$\text{context}_j^z = \sum_{j=1}^{m} \frac{\exp\left(q_j^{zT} k_{j'}^z\right)}{\sum_{o=1}^{m} \exp\left(q_j^{zT} k_o^z\right)} v_{j'}^z$$

195    Where $q_i^z, k_i^z \in \mathbb{R}^{d_{head} \times d_{head}}$ and $v_j^z \in \mathbb{R}^{d_{head} \times 1}$, the output $\text{context}_j^z$ was then flatted to

196    generate a vector with dimension $\mathbb{R}^{d_{head}}$. Finally, the different heads for each haplotype

197    were then concatenated and linear transformed $c_j = W_c[\text{context}_j^1; \dots; \text{context}_j^{n_{head}}]$

198    where $W_c \in \mathbb{R}^{d \times d}$ and $c_j \in \mathbb{R}^d$ will be added to the original token $i$ vector to update the

199    token. Figure 2 provides a schematic view of inter-haplotype block attention.

200

201    **3.2 Simulation dataset.**

202    The simulation datasets were achieved by re-sampling approaches with existing

203    genotype data as reference panels, thereby retaining allele frequency and LD patterns

204    (Wright et al., 2007). In the current study, the re-sampling based method Hapgen2 (Su

205    et al., 2011) was applied with 1000Genomes (Auton & Salcedo, 2015) as a reference

206    panel. In total, chromosome 20 of 1000 individuals was resampled and subjected to

207 the following analysis. The haplotype block was parsed based on the confidence

208 interval method (Gabriel et al., 2002) by *PLINK* (Purcell et al., 2007).

209

210 In our current study, we first simulate the expression of $h_i$ of each haplotype block $i$

211 according to a linear combination of all SNPs in the haplotype block.

$$h_i = \sum_{j=1}^{d_i} \alpha_{ij} x_{ij}$$

212

213 The phenotype was then simulated based on three epistasis models, which were

214 originally proposed by Burton et al (2007). More specifically, model 1 reflects an

215 epistasis model where the odds of disease increase multiplicatively within and

216 between 2-way disease markers. Using $h_i$ and $h_j$ to denote the expression of haplotype

217 $i$ and $j$, $\alpha$ and $\theta$ to denote the baseline and the factor of odd disease increase. Model1:

218 $$\text{odds}[h_i, h_j] = \alpha \times (1 + \theta_1)^{h_i} \times (1 + \theta_2)^{h_j} + \epsilon$$

219 In contrast, model 2 represents a disease model where the odds of disease only

220 increase unless both loci have at least one disease-associated allele,

221 $$\text{odds}[h_i, h_j] = \alpha \times (1 + \theta)^{(h_i) \times (h_j)} + \epsilon$$

222 Model 3 is similar to model 2, but renders a simpler threshold model as

223 $$\text{odds}[h_i, h_j] = \alpha \times (1 + \theta) + \epsilon$$

224

225 To simulate the simple epistasis model, two haplotype blocks will be randomly

226 selected each time and the phenotype model 1, model 2 or model 3 will be simulated

227 accordingly. To simulate the complex epistasis model, one dataset will contain

228 multiple epistasis from different SNP pairs. Here, we use the 'Combined Model

229    1+2+3' as a complex epistasis model which contains three epistasis from the previous

230    three basic models.

231

**3.3 Baseline models.**

233    To benchmark our model, we selected two state-of-the-art approaches for comparison

234    with the current proposed framework, matrixEpistasis (Zhu and Fang, 2017) and

235    GWAS_NN (Cui et al., 2022).  MatrixEpistasis represents a state-of-the-art method

236    for exhaustive epistasis searching. In contrast, GWAS_NN is one of the few methods

237    in the current field that tackles epistasis detection using neural networks, while also

238    providing an interpretation for the observed results. The GWAS_NN model first

239    learns the genetic block representations from all SNPs of a genetic block in a shallow

240    layer and then learns the complex relationships between genetic blocks in a deep

241    layer. These two baseline models exemplify the two main categories for epistasis

242    detection: exhaustive searching and machine learning. Both baseline models were

243    operated with the default settings unless indicated otherwise.

244

**3.4 Cohort description and statistical analysis.**

246    **3.4.1 Cohort description.** This research has been conducted using the UK Biobank

247    Resource. A material transfer agreement was signed with UK Biobank that covers

248    Research Tissue Bank (RTB) under projects 49731 and 59642.  The UK Biobank

249    study began in 2006 and, by 2010, had recruited over 500,000 participants from the

250    general UK population, aged 40 to 69 at the time of enrollment. The UK Biobank

251    genetic data contains genotypes for 488,377 participants. These were assayed using

252 two different yet similar assays, Applied Biosystems UK BiLEVE Axiom Array by

253 Affymetrix (Thermo Fisher Scientific) and Applied Biosystems UK Biobank Axiom

254 Array. More detail on the assay and quality control can be found in the UK Biobank

255 Genotyping and Quality Control

256 (https://biobank.ctsu.ox.ac.uk/crystal/crystal/docs/genotyping_qc.pdf.). Individuals

257 included in the current study from UK Biobank have T2D and are of European

258 descent. The UK The REC reference for UK Biobank is 16/NW/0274.

259

260 **3.4.2 Candidate genes selection.** To test the potential application of HB-LT in a real-

261 world scenario, we applied HB-LT to pre-selected glycated haemoglobin associated

262 genes. The candidate genes were first extracted from the DisGeNET database (Piñero,

263 et al., 2015). In total, fourteen genes were extracted and their coordinates

264 (chromosome, gene start position, and gene end position) were obtained using the

265 BioMart Project martview tool (Supplemental material). Next, SNPs located in each

266 gene ± 10 kbps were extracted in *PLINK* (Purcell et al., 2007). The thresholds set for

267 quality control including, imputation quality, Hardy-Weinberg equilibrium,

268 genotyping missing data across individuals, and genotyping missing rate were 0.8, 10-

269 10, 0.05, and 0.05 respectively.

270

271 **3.4.3 Covariants pre-filtering.** The datasets were subjected to a PCA-based covariant

272 pre-filtering stage to reduce the confounding effects before they feed into the HB-LT

273 for potential epistasis signal mining. Four covariants, including sex, age of diabetes

274 diagnosis, diabetes duration and population genetic structure were standardised

13

275 (*Scikit-learn* package) and subjected to PCA (*Scikit-learn* package) to reduce into a 2-

276 dimensional space. A $1 \times 1$ square was then applied to locate the most densely

277 populated area and individuals within this area were selected and subjected to the

278 following analysis.

279

280 **3.5 Software support.**

281 We conducted model-building and statistical analysis mainly using Python 3.9

282 (https://www.python.org/) and additional packages including *Pandas* (2.2.2), *Numpy*

283 (1.26.4), *PyTorch* (2.3.1), and *Tensorflow* (2.16.1). Other software includes R 4.2.2

284 (https://www.r-project.org/), Hapgen2

285 (https://mathgen.stats.ox.ac.uk/genetics_software/hapgen/), PLINK 1.9

286 (https://www.cog-genomics.org/plink/) and BioMart Project martview

287 (https://mart.ensembl.org/). The figures in the current study were drawn by *Matplotlib*

288 (3.8.2) and *Plotly* (5.22.0).

289

290 **3.6 Data availability statement.**

291 UK Biobank data are available to registered investigators under approved applications

292 (http://www.ukbiobank.ac.uk). Other relevant data are available from the

293 corresponding author upon request. The source code will be available shortly after the

294 deposition.

295

296 4. Results

297 **Long short-term memory (LSTM) selects potential phenotype-associated signals**

298 **as stage 1 of the current model.**

299    To pre-select potential phenotype-associated candidates from pools of haplotype

300    blocks in the human genome while maintaining a feasible computational burden, it is

301    essential to implement efficient filtering techniques. These techniques should be

302    capable of learning the representations of haplotype blocks by considering all the

303    SNPs within each block to the phenotype. Long short-term memory (LSTM) is a

304    recurrent neural network that is capable of learning long-range dependency and can

305    process sequences with variable lengths. It was widely used in datasets that process

306    "sequential" properties, such as natural language translation, before the introduction of

307    the Transformer model (Vaswani et al., 2017). Nevertheless, LSTMs still demonstrate

308    several advantages, especially for small datasets such as haplotype blocks, making

309    them a potentially effective technique for the filtering stage.

310

311    In the current study, we evaluated three epistasis models (model 1, model 2, and

312    model 3) as described by Burton et al. (2007), using chromosome 20 data from 1,000

313    individuals. Single nucleotide polymorphisms (SNPs) were pre-organized into

314    haplotype blocks using *PLINK*. Detailed descriptions of the dataset simulation and

315    haplotype block parsing methods can be found in the Methods section. Each haplotype

316    was tested individually by LSTM (n=10) and the root of mean square error was

317    recorded each time. Figure 3 shows the LSTM performance of chromosome 20 with 2

318    random haplotype blocks selected as epistasis signals of model 1, model 2 and model

319    3 (from top to bottom) at one record. Multiple valleys can be observed in all three

320    plots, indicating the presence of real epistasis signals. To quantify the performance of

321    the model across multiple runs with different haplotype block sizes, we repeated this

322   process 10 times and recorded the ROC AUC each time. As shown in the figure, the

323   LSTM can distinguish the epistasis signals from the background noise, achieving an

324   area under the curve (AUC) close to 1 in all three simple epistasis models.

325   Additionally, we conducted tests using the complex epistasis model to evaluate the

326   LSTM's capability in distinguishing signals. This complex epistasis model

327   encapsulates the combined effects of model 1, model 2, and model 3, which feature

328   multiple epistatic interactions from different haplotype block pairs. Employing

329   assessment criteria similar to those used for the simple epistasis model, Figure 3

330   shows that the LSTM demonstrates robust performance in accurately identifying all

331   haplotype blocks which contain complex epistasis signals. In short, in both simple and

332   complex epistasis models, the LSTM is an effective tool for selecting potential

333   candidates from large haplotype pools, significantly reducing the computational

334   burden and increasing the calculation power for subsequent hierarchical transformer

335   analysis.

336

337   **The hierarchical transformer encoder maps the haplotype block interactions**

338   **with the complex epistasis model as stage 2 of the current study.**

339   After selecting the potential phenotype-associated haplotype block candidates, we

340   applied a hierarchical transformer encoder and continued with the simulation datasets

341   to assess its ability for epistasis signal detections. The hierarchical transformer

342   encoder is a modified version of the hierarchical transformer that was originally

343   proposed by Liu and Lapata (2019). The potential epistasis signals were quantified

344   and visualised using the attention weights, which served as the main output of the

345    current hierarchical transformer encoder. "Attention," first introduced by Vaswani et

346    al. (2017), is a mechanism in the transformer neural network that enables the model to

347    dynamically weigh the importance of each element in an input sequence relative to the

348    others. The attention weights, the main output of our current model, are computed

349    using a scaled dot-product that quantifies this "relatedness" between pairs of haplotype

350    blocks and SNPs (within haplotype block). Several studies (Ahmed, Aly and Liu,

351    2024; Graça et al., 2024) have demonstrated these attention weights that learned by

352    the model, at least partially, can be interpreted as the epistasis interactive scores

353    between genetic regions.

354

355    In the hierarchical transformer encoder, attention weights are initially mapped

356    between each SNP within each haplotype block. Subsequently, a fixed-length

357    representation of each haplotype is generated using multi-head pooling. Multi-head

358    attention weights are then calculated between haplotypes. This updated information is

359    incorporated into the original SNP embedding and processed through the feed-forward

360    layer. For further details, please refer to the Methods section. Figure 4a illustrates the

361    Root Mean Square Error (RMSE) and loss for both training and testing datasets, split

362    in a ratio of 0.8:0.2, for a single record. The plots demonstrate a smooth training

363    trajectory with no significant discrepancies between the training and testing datasets,

364    indicating stable model performance and effective training without overfitting.

365    Additionally, this process was repeated 10 times, each iteration using different pre-

366    selected simulated signals (Figure 4b). The model consistently demonstrated a robust

367    performance in phenotype predictions with mean RMSE lower than 0.012 for both

368    training and testing datasets.

369

370    More importantly, to evaluate whether the simulated complex epistasis signals

371    (model1+2+3) can be at least partially captured by the cross-haplotype attention

372    weights, which are the primary output of our current HB-LT model, we compared

373    these weights with the ground truth matrix. We then plotted the ROC AUC curve

374    (n=10), as shown in Figure 4c. The model achieved an average AUC of 0.83. In short,

375    after the LSTM pre-selected the potential phenotype-associated haplotype blocks, the

376    hierarchical transformer encoder demonstrated reasonable performance in

377    distinguishing potential interaction/epistasis signals between haplotype blocks.

378

379    **HB-LT outperforms baseline models for both simple and complex epistasis**

380    **models.**

381    Traditional exhaustive epistasis searching methods are often facing a challenge

382    referred to as the "curse of dimensionality". The SNPs that are involved in the

383    epistatic interactions might have low minor allele frequencies, however, the variants

384    to be tested can be huge. As a result, detecting these interactions becomes challenging

385    due to the reduced statistical power and the increased likelihood of both type 1 and

386    type 2 errors. In the current study, we selected matrixEpistasis (Zhu and Fang, 2017)

387    as one of the representatives for an exhaustive epistasis searching method against our

388    current HB-LT model.

389    In contrast, we also tested two deep learning neural network baseline models,

390    GWAS_NN (Cui et al., 2022) and LSTM with the vanilla transformer encoder, against

391    our current HB-LT model. In GWAS_NN, the long sequence of SNPs was initially

392    divided into different genetic blocks (SNPs layer). Fully connected multilayer

393    perceptrons (MLPs) were then used to learn a fixed representation for each genetic

394    block $(g_1, \ldots, g_m)$. Another set of MLPs was subsequently trained to learn the epistasis

395    between these genetic blocks. Additionally, we also tested the LSTM with the vanilla

396    transformer encoder against our HB-LT model to assess the potential advantages of

397    the hierarchical transformer structure compared to the vanilla single SNP transformer

398    structure for epistasis studies.

399

400    A total of 6,156 haplotype blocks (60,501 SNPs) from chromosome 20 in the

401    simulation datasets were analysed using three baseline models and the HB-LT model.

402    These models were evaluated under both simple and complex epistasis conditions

403    with recorded ROC AUC (n=10). Overall, all four models showed robust performance

404    in identifying epistasis signals. MatrixEpistasis demonstrated a stable yet

405    comparatively low performance, with an average of around 0.73 across both simple

406    and complex epistasis models. This outcome likely reflects the reduced statistical

407    power of exhaustive searching methods when handling large SNP datasets. In

408    contrast, HB-LT exhibited the highest performance across both simple and complex

409    models, although this advantage was less pronounced when dealing with complex

410    models. The reason HB-LT has a larger ROC AUC than baseline models could mainly

411    be because HB-LT employs multi-head pooling to utilise all SNPs for representing

412    haplotype blocks in a supervised manner. This approach not only reduces

413    dimensionality compared to considering each SNP individually but also provides more

414    informative representations than selecting a single SNP or using unsupervised

415    methods such as PCA.

416

417

418    **Interaction discovery in a diabetes glycated haemoglobin study.**

419    To inspect how trustworthy our proposed framework is in a real-world scenario,

420    experiments on real-world cohort, UK BioBank for glycated haemoglobin (HbA1c)

421    are performed. Glycated haemoglobin (HbA1c) is the most common biomarker used

422    to monitor glucose control in diabetes patients (WHO, 2011), which reflects the

423    glycemic load ~ 3 months and traits such as hemoglobinopathies and alteration in

424    intracellular glucose metabolism (Nathan, Turgeon and Regan, 2007). HbA1c levels

425    are influenced by both environmental and genetic factors. Research studies (Snieder et

426    al., 2001; Meigs et al., 2002) estimating the heritability of HbA1c in non-diabetic

427    individuals report a range from 27% to 62%, providing strong evidence of a

428    significant genetic component in HbA1c variability. Previous genome-wide

429    association studies (GWAS) (Wheeler et al., 2017) have identified more than 100

430    genetic variants to be associated with HbA1c. In this study, we re-examined the

431    HbA1c-associated loci in the UK Biobank cohort to explore potential novel epistasis

432    signals, both within and across haplotype blocks. To limit the potential confounding

433    factors in our current study, patients were pre-filtered based on 4 covariants (age of

434    diabetes diagnosis; diabetes duration; sex and population genetic structure) in UK

435    Biobank to select individuals with similar characteristics and eliminate the potential

436    outliers. The clinical features of individuals from the UK Biobank used in the current

437    study are shown in Table 1. The chosen datasets have 1277 individuals, with in total

438    of 14 genes with 10kbp flanking regions added to both ends, 74 haplotype blocks and

439    1821 SNPs. The haplotype blocks were parsed based on the confidence interval

440    method (Gabriel et al., 2002). More details regarding individual covariants pre-

441    selection and haplotype block parsing can be referred to the Methods section.

442

443    **Cross-haplotype block epistasis**

444    Unlike the simulated datasets, the interacting haplotype blocks are not known in the

445    real datasets. One of the common approaches to validate the proposed framework's

446    prediction and interpretation is to find previous works that report genes and epistatic

447    relationships on the related disease. Afterwards, the objective is to map the framework

448    haplotype block outputs to the target genes. In the UK Biobank, using HB-LT, we

449    observed 7 pairwise interaction candidates. All interaction candidates of the HbA1c

450    phenotype and their corresponding attention scores were listed in the first two

451    columns of Table 2. The threshold is set as attention scores higher than 0.1. There are

452    no standard ways to choose the attention threshold. For the future studies, researchers

453    can set the threshold wherever they think that fits their hypothesis.

454    We then investigate whether the interaction candidates discovered by HB-LT can be

455    detected by other methods. Similarly to the simulation dataset section, we trained

456    GWAS_NN and recorded the interaction scores for each interaction candidate in the

457    third column in Table 2. Not all signals detected by HB-LT can also be mapped out by

458    GWAS_NN as significant. This could potentially highlight that HB-LT can detect

459    novel signals which can be overlooked by other methods. Finally, We then check if

460    any of the gene interactions have already been recorded in the previous studies, shown

461    in the last column of Table 2.

462

463    **Within-haplotype SNP epistasis**

464    One of the key advantages of HB-LT compared to other multi-dimensional reduction

465    methods is that it trains each SNP individually before pooling them into a fixed

466    haplotype representation. This approach preserves the haplotype structure during the

467    training process, allowing us to monitor not only potential cross-haplotype

468    interactions but also SNP interactions within each haplotype block. Indeed, it is

469    believed that SNPs within a functional region have a higher chance to interact with

470    each other and influence the phenotype (Ma, Clark and Keinan, 2013). We listed all

471    the within-haplotype SNP interaction candidates, gene names and attention scores in

472    the first three columns of Table 3. By setting the threshold $> 0.05$, there are 16 within

473    haplotype block pairwise SNP interactions observed. Similarly, we compared the

474    within-haplotype pairwise SNP interaction results obtained from HB-LT with those

475    identified using the previously published exhaustive search method, MatrixEpistasis

476    and recorded the *p-value* in the last column of Table 3.

477

478    In summary, leveraging pre-selected HbA1c-associated genes from the UK Biobank,

479    we explore the potential real-world application of HB-LT. Future studies are essential

480    to statistically and biologically validate our current findings. Additionally, further

481  investigation is warranted to assess the feasibility of HB-LT in hypothesis-free, large-

482  scale genetic datasets.

483

484  5. Discussion

485  The current existing approaches for detecting interactions in genetic study face several

486  challenges including: (i) to reduce the "curse of dimensionality", genes are typically

487  represented as the most important SNPs while ignoring the potential effects of

488  neighbouring SNPs. (ii) only restrictive forms of interactions are considered. (iii)

489  while scalable methods like PCA have been proposed to reduce the multidimensional

490  SNP data into a one-dimensional representation, these techniques often neglect

491  important phenotype-related information. Additionally, such methods make it difficult

492  to analyse the internal structure of each condensed multi-SNP unit in relation to the

493  phenotype. Indeed, there is a need for a framework that integrates genetic block

494  representation learning with the modelling of both intra- and inter-block interactions

495  within an end-to-end model. Here, we proposed a deep learning genetic block

496  detection method, Haplotype Block LSTM hierarchical Transformer (HB-LT). HB-LT

497  can hierarchically encode genetic SNP sequences. In HB-LT, each SNPs in

498  relationship to its surrounding SNPs within each genetic block were learned by a

499  multi-attention head. Next, a pooling method is applied to get a fixed representation of

500  each genetic block, cross genetic block relationships via an attention method were

501  then mapped as opposed to concatenating dimensional condensed genetic units into

502  flat sequences and then fed into the model. This approach enables the model to

503  dynamically learn richer structural dependencies among SNPs within each genetic

504    block and effectively incorporate these insights into the inter-genetic block layer. In

505    the experimental work, the results obtained from HB-LT were compared with both

506    exhaustive searching (MatrixEpistasis) and deep learning methods (GWAS_NN). The

507    current study shows that HB-LT exhibits a better performance for epistasis detections

508    in both simple and complex epistasis models. Moreover, HB-LT was also tested on

509    HbA1c in the UK Biobank to assess its application in a real-world scenario.

510

511    The current proposed deep learning framework may have many attractive features, but

512    it also has several shortcomings. First, it should be noted, that although the potential

513    applications of attention weights in transformer as an indicator of epistasis have been

514    studied and demonstrated in recent years (Reyes et al., 2022; Graça et al., 2024), these

515    weights cannot be directly interpreted as measurements of epistasis levels in genetics.

516    While these weights can highlight regions of interest in relation to the phenotype, the

517    focus of attention mechanisms is on capturing token dependencies and relevant

518    patterns in the dataset to improve outcome predictions, not necessarily to isolate or

519    quantify specific genetic interactions. Moreover, the intricate relationship between

520    statistical and biological epistasis adds an additional layer of complexity (Moore and

521    Williams, 2005; Ebbert, Ridge, and Kauwe, 2015). The disparity between these two

522    models of epistasis often obscures the biological relevance and implications of

523    statistical findings, making it challenging to draw clear, meaningful conclusions about

524    the underlying genetic mechanisms. Indeed, we view our proposed HB-LT framework

525    as a tool for mining and filtering large datasets. Regions of interest identified by HB-

526    LT should be further investigated and validated through more targeted statistical

527     methods and potential complementary biological experiments. Second, the potential

528     covariants, such as age and population genetic structure are not directly incorporated

529     in the current model, instead, a PCA-based pre-filtering stage was applied to minimise

530     the confounding effects. By implementing this approach, we minimised the risk of

531     extraneous factors contaminating the HB-LT model outputs, making the results more

532     straightforward to interpret. However, this adjustment means that the HB-LT model

533     will not be applied to the full dataset size, potentially leading to a loss of information.

534     Additionally, this change may introduce new complexities for future users. The

535     challenge of how to incorporate potential covariates into the model remains an open

536     question that needs to be addressed. Finally, the potential of whole-genome

537     hypothesis-free epistasis studies to significantly enhance outcome prediction has been

538     a topic of debate and scepticism for decades. Several studies (González-Camacho et

539     al., 2012; Mäki-Tanila and Hill, 2014 and Wei et al., 2014) have demonstrated that

540     despite the possible biological ubiquity of epistasis, the total genetic variance of

541     polygenic traits is likely largely to be explained by the "additive top SNPs model".

542     However, other studies (Dudley and Johnson, 2009; Hu et al., 2011; Álvarez-Castro et

543     al., 2012; Wang et al., 2012) conducted showed that the inclusion of epistatic effect

544     networks for prediction improved prediction over the use of additive effects only.

545     Indeed, we do not intend to propose HB-LT as a replacement for the "top SNPs

546     approach". This study is not aimed at comparing epistasis and non-epistasis models

547     for outcome prediction. Instead, HB-LT serves as a complementary tool designed to

548     uncover interactions that might otherwise be overlooked, potentially revealing novel

549     genetic structures underlying complex disease development and leading to the

550     discovery of new markers.

551

552     6. Conclusion

553     For the past decades, Genome-wide Association Studies (GWAS) have successfully

554     identified thousands of risk alleles for complex diseases. Despite this, it usually failed

555     to capture the statistical epistasis i.e. interaction between SNPs, which is

556     acknowledged as a fundamental factor for understanding complex disease genetic

557     pathways. Traditional epistasis searching tools often suffer from computational burden

558     and lack of calculation power. Moreover, it has become increasingly recognized that

559     leveraging the combined effects of SNP groups within specific genetic blocks for

560     epistasis searching can yield greater phenotypic variance than focusing on individual

561     SNPs alone. However, to our best knowledge, there is yet a framework that has been

562     proposed that can incorporate this hierarchical genetic structure to form an end-to-end

563     model for epistasis searching. Here, we proposed HB-LT, which takes advantage of

564     the haplotype block structure existing in the genome to reduce the dimensionality of

565     SNP features and increase statistical power. Haplotype block is not the only way to

566     measure SNP dependencies and grouping, the impact of different methods for

567     epistasis study should be investigated in the future.

568

569     7. Acknowledgement

572    16/0005485) and NIH-NIDDK (R01DK135268) project grants, a CIHR-JDRF Team

573    grant (CIHR-IRSC TDP-186358 and JDRF 4-SRA-2023-1182-S-N), CRCHUM start-

574    up funds, and an Innovation Canada John R. Evans Leader Award (CFI 42649). AB

575    was supported by the Canada Research Chairs program.

576

577    **Duality of Interest:** G.R. has received grant funding from, and is a consultant for,

578    Sun Pharmaceuticals Inc. All other authors declare that there are no relationships or

579    activities that might bias, or be perceived to bias, their work.

580

581    **Contribution Statement:** The current project was designed by S.L. and S.A. under

582    the supervision of G.R. and with professional support from P.H., J.T., A.B. and B.L.

583    The current model (HB-LT) was designed by S.L. and S.A. The source code was

584    written by S.A. and S.L. The model training and testing were performed by S.L. and

585    S.A. The real-world cohort genetic datasets were organized, maintained and queried

586    by R.A. The manuscript was written by S.L.

587

588    8. References

589    Ahmed, F. S., Aly, S. and Liu, X. (2024). EPI-Trans: an effective transformer-based

590    deep learning model for enhancer promoter interaction prediction. BMC

591    Bioinformatics, 25 (1), p.216.

592

593    Álvarez-Castro, J. M. et al. (2012). Modelling of genetic interactions improves

594    prediction of hybrid patterns--a case study in domestic fowl. Genetics Research, 94

595    (5), pp.255–266.

596

597    Auton, A. and Salcedo, T. (2015). The 1000 Genomes Project. Assessing Rare

598    Variation in Complex Traits, pp. 71–85. doi:10.1007/978-1-4939-2824-8_6.

599

600    Bayat, A. et al. (2021). Fast and accurate exhaustive higher-order epistasis search with

601    BitEpi. Scientific Reports, 11 (1), p.15923.

602

603    Burton, P.R. et al. (2007). Genome-wide association study of 14,000 cases of seven

604    common diseases and 3,000 shared controls. Nature, 447(7145), pp.661–678.

605    doi:10.1038/nature05911.

606

607    Chen, H. et al. (2020). A fast-linear mixed model for genome-wide haplotype

608    association analysis: application to agronomic traits in maize. BMC Genomics, 21 (1),

609    p.151.

610

611    Chen, Y. et al. (2023). Molecular language models: RNNs or transformer? Briefings

612    in Functional Genomics, 22 (4), pp.392–400.

613

614    Cordell, H. J. (2002). Epistasis: what it means, what it doesn't mean, and statistical

615    methods to detect it in humans. Human Molecular Genetics, 11 (20), pp.2463–2468.

616

617    Cui, T. et al. (2022). Gene-gene interaction detection with deep learning.

618    Communications Biology, 5 (1), p.1238.

619

620    Cui, T. et al. (2022). Gene–gene interaction detection with deep learning.

621    Communications Biology, 5(1). doi:10.1038/s42003-022-04186-y.

622

623    De Franco, E. et al. (2020). Update of variants identified in the pancreatic β-cell

624    KATP channel genes KCNJ11 and ABCC8 in individuals with congenital

625    hyperinsulinism and diabetes. Human Mutation, 41 (5), pp.884–905.

626

627    Dudley, J. W. and Johnson, G. R. (2009). Epistatic models improve prediction of

628    performance in corn. Crop Science, 49 (4), pp.1533–1533.

629

630    Ebbert, M. T. W., Ridge, P. G. and Kauwe, J. S. K. (2015). Bridging the gap between

631    statistical and biological epistasis in Alzheimer's disease. BioMed Research

632    International, 2015, p.870123.

633

634    Gabriel, S.B. et al. (2002). The structure of haplotype blocks in the human genome.

635    Science, 296(5576), pp.2225–2229. doi:10.1126/science.1069424.

636

637    González-Camacho, J. M. et al. (2012). Genome-enabled prediction of genetic values

638    using radial basis function neural networks. TAG. Theoretical and Applied Genetics,

639    125 (4), pp.759–771.

640

641    Graça, M. et al. (2023). Interpreting High Order Epistasis Using Sparse Transformers.

642    [Online]. Available at: https://ieeexplore.ieee.org/document/10183767.

643

644    Graça, M. et al. (2024). Distributed transformer for high order epistasis detection in

645    large-scale datasets. Scientific Reports, 14 (1), p.14579.

646

647    Hill, W.G. and Robertson, A. (1968). Linkage disequilibrium in finite populations.

648    Theoretical and Applied Genetics, 38(6), pp.226–231. doi:10.1007/bf01245622.

649

650    Hu, Z. et al. (2011). Genomic value prediction for quantitative traits under the

651    epistatic model. BMC Genetics, 12, p.15.

652

653    Jubair, S. et al. (2021). GPTransformer: A Transformer-Based Deep Learning Method

654    for Predicting Fusarium Related Traits in Barley. Frontiers in Plant Science, 12,

655    p.761402.

656

657    Li, J. et al. (2009). Identification of gene-gene interaction using principal components.

658    BMC Proceedings, 3 Suppl 7 (Suppl 7), p.S78.

659

660     Liu, Y. and Lapata, M. (2019). Hierarchical Transformers for Multi-Document

661     Summarization. arXiv [cs.CL]. arXiv [Online]. Available at:

662     http://arxiv.org/abs/1905.13164.

663

664     Ma, L., Clark, A. G. and Keinan, A. (2013). Gene-based testing of interactions in

665     association studies of quantitative traits. PLoS Genetics, 9 (2), p.e1003321.

666

667     Maher, B. (2008). Personal genomes: The case of the missing heritability. Nature, 456

668     (7218), pp.18–21.

669

670     Mäki-Tanila, A. and Hill, W. G. (2014). Influence of gene interaction on complex trait

671     variation with multilocus models. Genetics, 198 (1), pp.355–367.

672

673     Meigs, J. B. et al. (2002). A genome-wide scan for loci linked to plasma levels of

674     glucose and HbA(1c) in a community-based sample of Caucasian pedigrees: The

675     Framingham Offspring Study. Diabetes, 51 (3), pp.833–840.

676

677     Mieth, B. et al. (2021). DeepCOMBI: explainable artificial intelligence for the

678     analysis and discovery in genome-wide association studies. NAR Genomics and

679     Bioinformatics, 3 (3), p.lqab065.

680

681    Moore, J. H. and Williams, S. M. (2005). Traversing the conceptual divide between

682    biological and statistical epistasis: systems biology and a more modern synthesis.

683    BioEssays, 27 (6), pp.637–646.

684

685    Morris, R. W. and Kaplan, N. L. (2002). On the advantage of haplotype analysis in the

686    presence of multiple disease susceptibility alleles. Genetic Epidemiology, 23 (3),

687    pp.221–233.

688

689    Nathan, D. M., Turgeon, H. and Regan, S. (2007). Relationship between glycated

690    haemoglobin levels and mean glucose levels over time. Diabetologia, 50 (11),

691    pp.2239–2244.

692

693    Niel, C. et al. (2015). A survey about methods dedicated to epistasis detection.

694    Frontiers in Genetics, 6, p.285.

695

696    Oishi, Y. and Manabe, I. (2018). Krüppel-Like Factors in Metabolic Homeostasis and

697    Cardiometabolic Disease. Frontiers in Cardiovascular Medicine, 5, p.69.

698

699    Pérez-Enciso, M. and Zingaretti, L. M. (2019). A Guide for Using Deep Learning for

700    Complex Trait Genomic Prediction. Genes, 10 (7). [Online]. Available at:

701    doi:10.3390/genes10070553.

702

703    Piñero, J. et al. (2015). DisGeNET: a discovery platform for the dynamical

704    exploration of human diseases and their genes. Database: the journal of biological

705    databases and curation, 2015 (0), p.bav028.

706

707    Purcell, S. et al. (2007). PLINK: a tool set for whole-genome association and

708    population-based linkage analyses. The American Journal of Human Genetics, 81 (3),

709    pp.559–575.

710

711    Reyes, D. M. et al. (2022). Genomics transformer for diagnosing Parkinson's disease.

712    *IEEE-EMBS International Conference on Biomedical and Health Informatics. IEEE-*

713    *EMBS International Conference on Biomedical and Health Informatics,* 2022.

714    [Online]. Available at: doi:10.1109/bhi56158.2022.9926815.

715

716    Snieder, H. et al. (2001). HbA(1c) levels are genetically determined even in type 1

717    diabetes: evidence from healthy and diabetic twins. Diabetes, 50 (12), pp.2858–2863.

718

719    Stančáková, A. and Laakso, M. (2016). Genetics of Type 2 Diabetes. Endocrine

720    Development, 31, pp.203–220.

721

722    Su, Z., Marchini, J. and Donnelly, P. (2011). Hapgen2: Simulation of multiple disease

723    snps. Bioinformatics, 27(16), pp.2304–2305. doi:10.1093/bioinformatics/btr341.

724

725    Turton, J. C. et al. (2011). Investigating statistical epistasis in complex disorders.

726    Journal of Alzheimer's Disease, 25 (4), pp.635–644.

727

728    Vaswani, A. et al. (2017). Attention Is All You Need. arXiv [cs.CL]. arXiv [Online].

729    Available at: http://arxiv.org/abs/1706.03762.

730

731    Wang, D. et al. (2012). Prediction of genetic values of quantitative traits with epistatic

732    effects in plant breeding populations. Heredity, 109 (5), pp.313–319.

733

734    Wei, W.-H., Hemani, G. and Haley, C. S. (2014). Detecting epistasis in human

735    complex traits. Nature Reviews Genetics, 15 (11), pp.722–733.

736

737    Wheeler, E. et al. (2017). Impact of common genetic determinants of Hemoglobin

738    A1c on type 2 diabetes risk and diagnosis in ancestrally diverse populations: A

739    transethnic genome-wide meta-analysis. PLoS Medicine, 14 (9), p.e1002383.

740

741    World Health Organization. (2011). Glycated haemoglobin (HbA1c) for the diagnosis

742    of diabetes.

743

744    Wright, F.A. et al. (2007). Simulating Association Studies: A data-based resampling

745    method for candidate regions or whole genome scans. Bioinformatics, 23(19),

746    pp.2581–2588. doi:10.1093/bioinformatics/btm386.

747

748    Zaykin, D. V. et al. (2002). Testing association of statistically inferred haplotypes

749    with discrete and continuous traits in samples of unrelated individuals. Human

750    Heredity, 53 (2), pp.79–91.

751

752    Zhou, J. et al. (2022). Deep learning predicts DNA methylation regulatory variants in

753    the human brain and elucidates the genetics of psychiatric disorders. Proceedings of

754    the National Academy of Sciences of the United States of America, 119 (34),

755    p.e2206069119.

756

757    Zhu, S. and Fang, G. (2018). Matrixepistasis: Ultrafast, exhaustive epistasis scan for

758    quantitative traits with covariate adjustment. Bioinformatics, 34(14), pp.2341–2348.

759    doi:10.1093/bioinformatics/bty094.

760

761    Table 1. Characteristics of the included individuals of the cohorts.

| Clinical Features | UK Biobank  (n=1277)  (Mean+-SD) |
| --- | --- |
| Age of diabetes diagnosis | 56.3+-3.0 |
| Duration (year) | 4.0+-2.4 |
| Diastolic blood pressure | 83.0+-5.1 |
| Systolic blood pressure | 142.0+-8.8 |

| Sex (Male: Female) | 482:795 |
|---|---|

762

763    Table 2. Cross-haplotype blocks attention in UK Biobank by two independent

764    methods.

| Phenotype | Genes | Haplo Block* | HB-LT Attention Scores (UK Biobank) | GWAS_NN Interaction Scores (UK Biobank) | References |
|---|---|---|---|---|---|
| HbA1c | LOC112268412, KLF11, CYS1; GCK, LOC105375257, YKT6 | H2, H11 | 0.23 | 0.0012 | Oishi and Manabe, 2018 |
| | LOC112268412, KLF11, | H2, H15 | 0.22 | 0.17 | NA |

| | | | | | |
|---|---|---|---|---|---|
| | CYS1; BLK | | | | |
| | BLK, LOC105379241; NSMCE2 | H23, H30 | 0.13 | 0.09 | NA |
| | NSMCE2 | H30, H31 | 0.30 | 0.23 | |
| | CEL; KCNJ11 | H34, H41 | 0.13 | 0.19 | NA |
| | CEL; NCR3LG1, KCNJ11 | H34, H40 | 0.12 | 0.19 | NA |
| | ABCC8; KCNJ11 | H40, H42 | 0.12 | 0.008 | De Franco et al., 2020 |

765     * Full length of each Haplotype block seen in Supplemental material.

766

767     Table 3. Within haplotype block SNP attention in UK Biobank by both HB-LT and

768     MatrixEpistasis.

| Phenotype | Genes | Pairwise | Attention | P-value (UK BioBank) |
|---|---|---|---|---|

| | | SNPs | Scores (UK BioBank) HB-LT | MatrixEpistasis |
|---|---|---|---|---|
| HbA1c | NA | rs6756950, rs7420169 | 0.11 | 0.00037 |
| | HNF4A | rs736820, rs736824 | 0.11 | 0.00046 |
| | GCK | rs2908285, rs118180640 | 0.09 | 0.00023 |
| | YKT6, GCK | rs1814253, rs118180640 | 0.07 | 0.00163 |
| | GCK, NA | rs118180640, rs758983 | 0.07 | 0.00091 |
| | NEUROD1, NA | rs12053195, rs6756950 | 0.06 | 0.052 |

|  | NA, NEUROD1 | rs6756950, rs12052558 | 0.05 | 0.002 |
|---|---|---|---|---|
|  | NA, NEUROD1 | rs6756950, rs16867467 | 0.05 | 0.0002 |

769

770

771    **Figure Legend**

772    **Figure 1. Pipeline of our within and cross haplotype block epistasis detection**

773    **method (HB-LT).** The haplotype dataset was first selected by Long short-term

774    memory (LSTM) to obtain the potential phenotype-associated haplotype blocks

775    (HBs). These candidates are then fed into the Hierarchical transformer encoder to

776    obtain the within and cross-haplotype block attention weights, which could be

777    subjected to the following analysis for potential epistasis signal discoveries.

778

779    **Figure 2.  The overall pipeline of the hierarchical haplotype transformer**

780    **encoder.**

781

782    **Figure 3. Long short-term memory (LSTM) selects potential phenotype-**

783    **associated haplotype block candidates on both simple and complex epistasis**

784    **models using simulated datasets.** Each haplotype block was independently trained

785    and tested by the LSTM model. The root mean square error (RMSE) of the testing

786    datasets was plotted along with the mean and standard deviation (SD). To evaluate the

787    model's performance with different haplotype block lengths, the process was repeated

788    (n=10). For each iteration, the area under the curve (AUC) of the model's predictions

789    was plotted against the ground truth.

790

791    **Figure 4. The hierarchical transformer encoder distinguishes epistasis signals in**

792    **a complex epistasis model using simulated datasets.** Pre-selected haplotype block

793    candidates identified by the long short-term memory (LSTM) model were
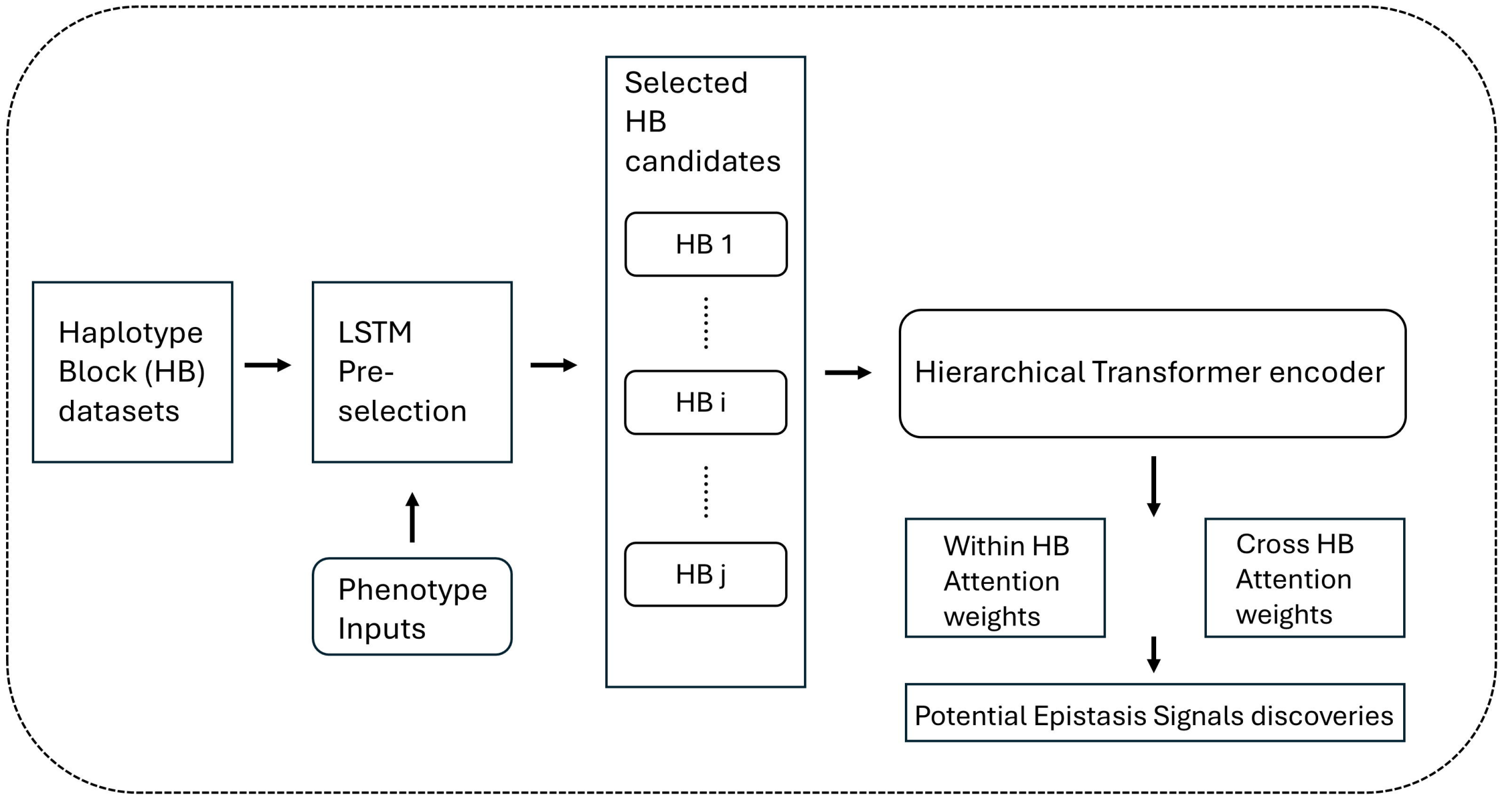
794    subsequently processed by the hierarchical transformer encoder. (a) Training and

795    testing loss, and root mean square error (RMSE) for a single record. (b) Box plot

796    showing the RMSE for training and testing datasets (n=10). (c) The area under the

797    curve (AUC) for the complex epistasis signals predicted by the hierarchical

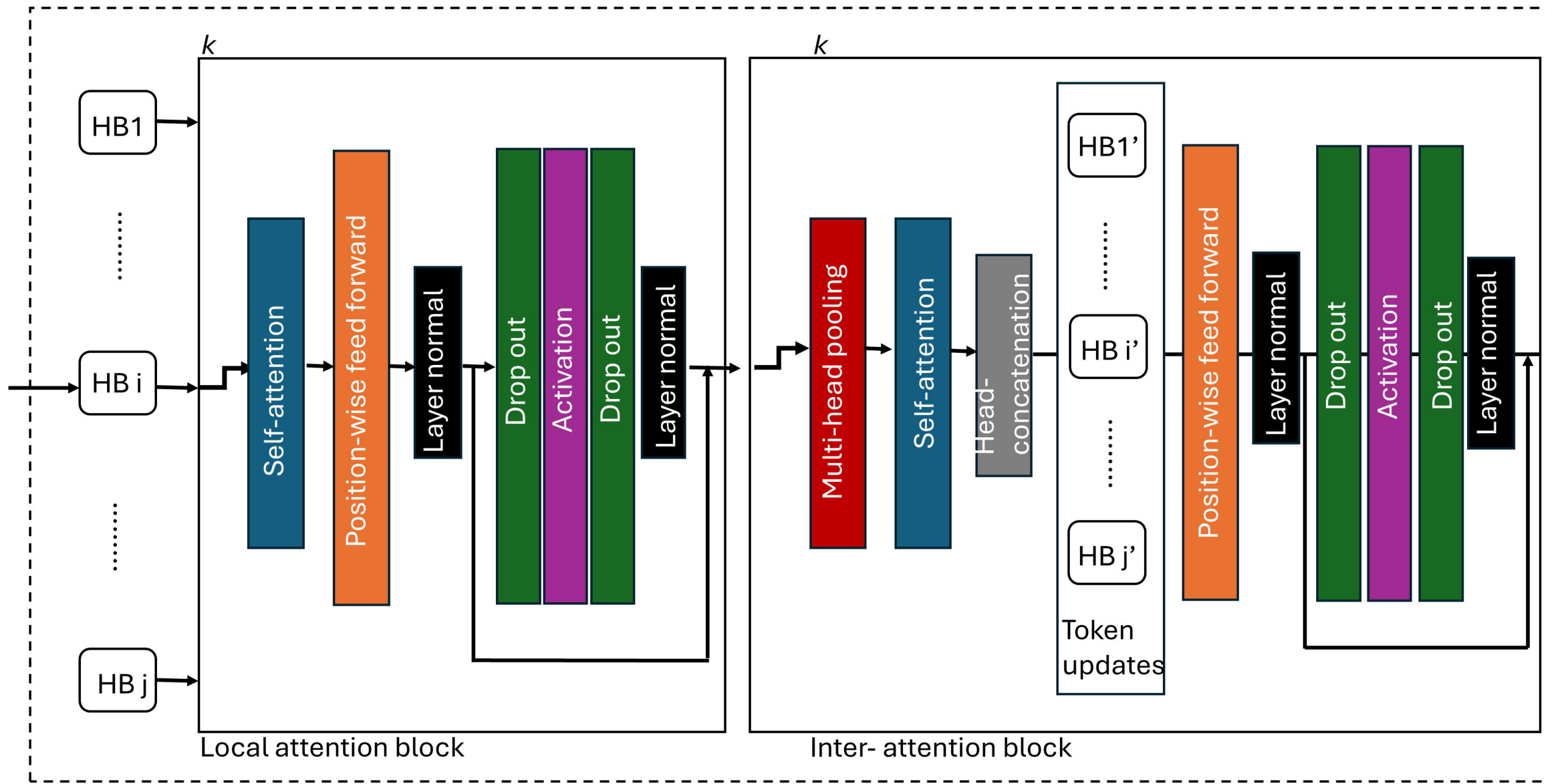798    transformer encoder, compared to the ground truth.

799

800    **Figure 5. The comparison of HB-LT with state-of-art Epistasis searching**

801    **methods, MatrixEpistasis (Zhu and Fang, 2017) and GWAS_NN (Cui et al.,**

802    **2022) in both simple and complex epistasis models with n=10.**
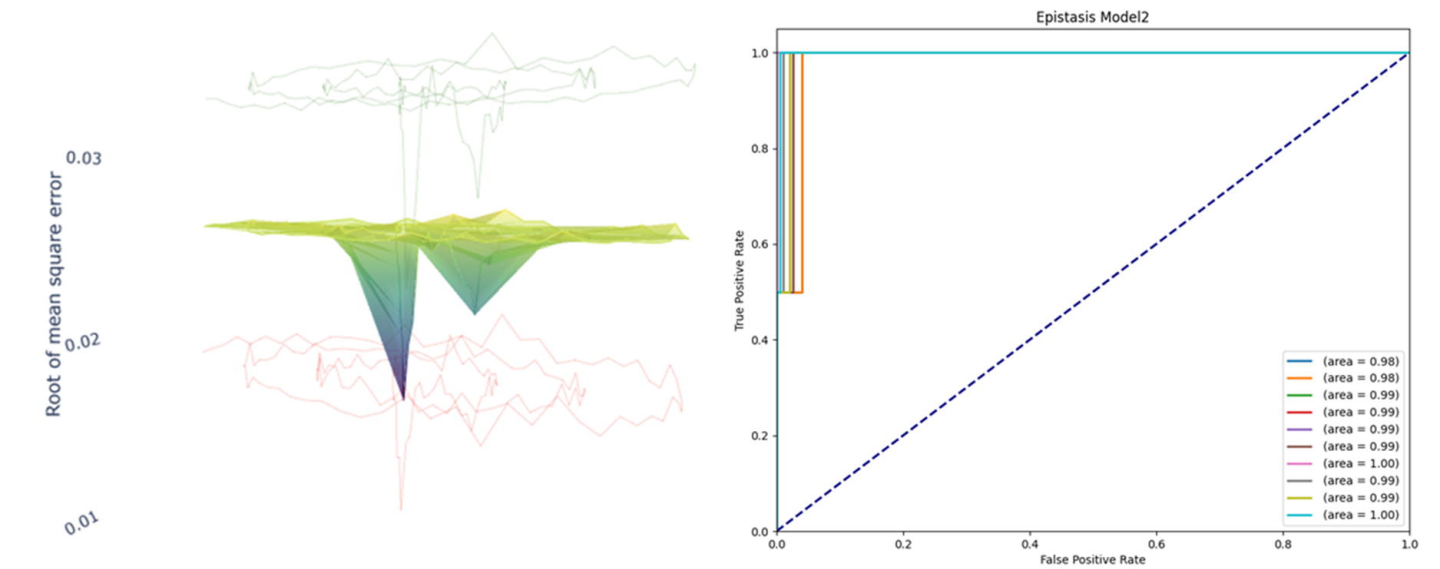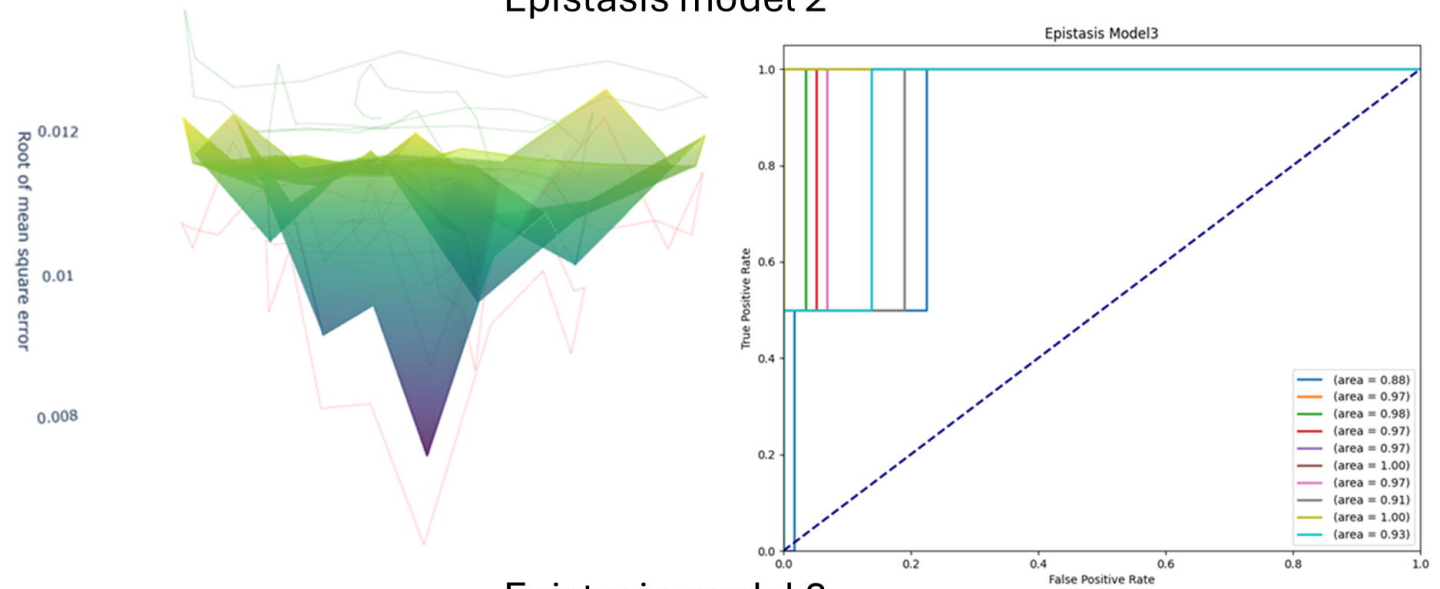
803

804

HB-LT

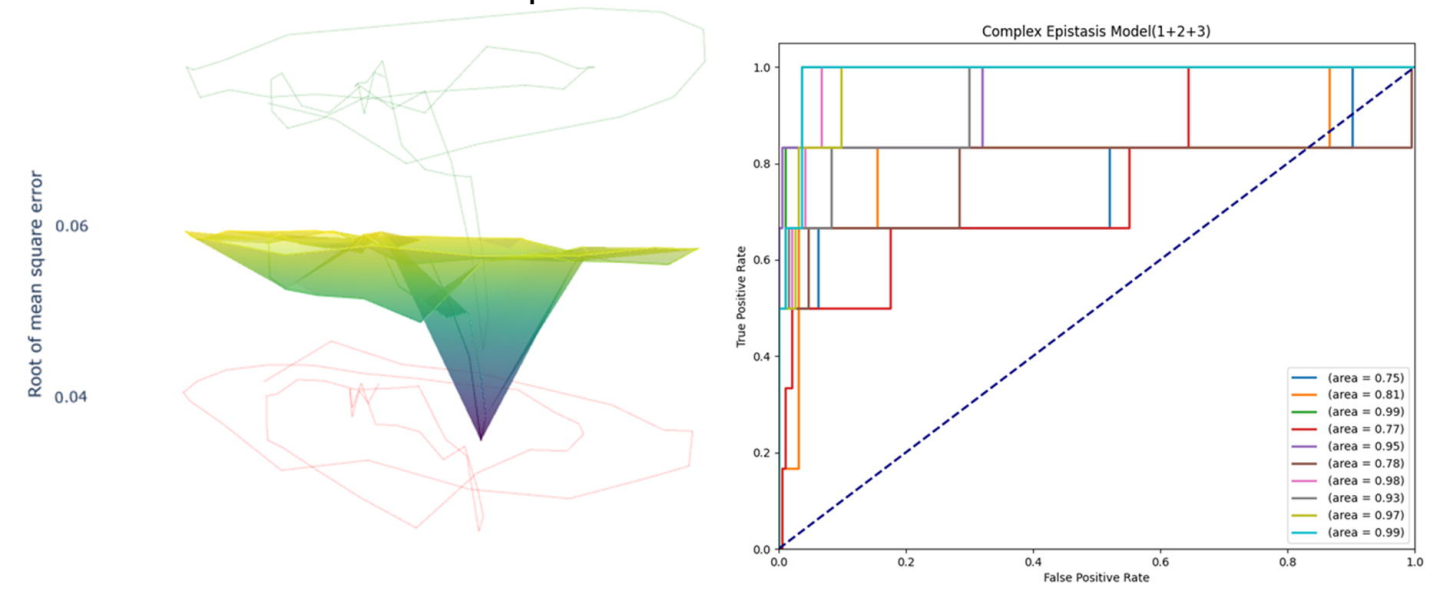Hierarchical haplotype transformer encoder

Epistasis model 1

Epistasis model 2

Epistasis model 3

Complex Epistasis model (model1+2+3)

a.

Training and Validation Loss Over Epochs

— Training Loss
— Validation Loss

Loss

Epoch

Training and Validation RMSE Over Epochs

— Training RMSE
— Validation RMSE

RMSE

Epoch

b.

Root of mean square error

RMSE

Training RMSE                    Validation RMSE

c.

Receiver Operating Characteristic (ROC) Curve

True Positive Rate

False Positive Rate

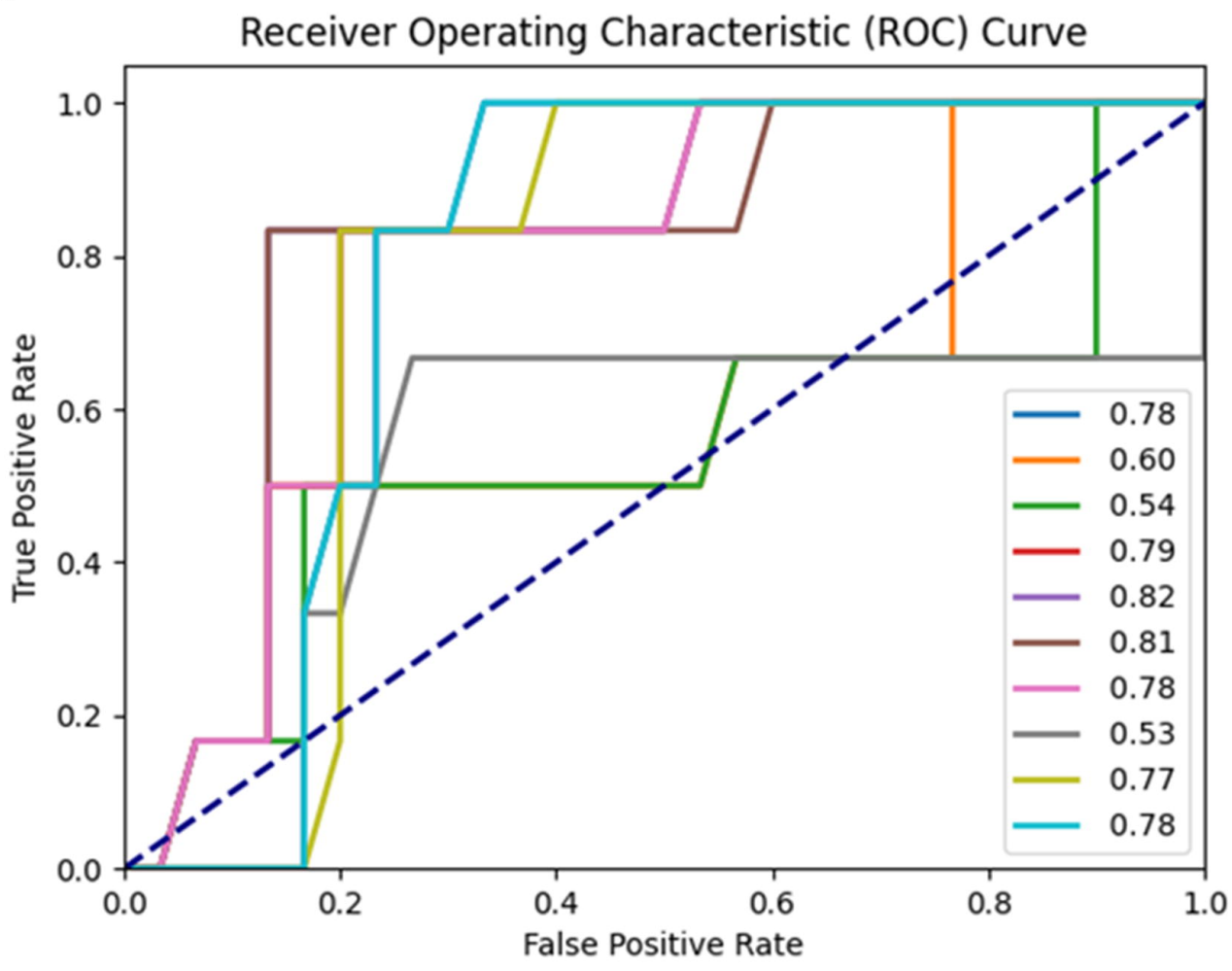| | |
|---|---|
| — | 0.78 |
| — | 0.60 |
| — | 0.54 |
| — | 0.79 |
| — | 0.82 |
| — | 0.81 |
| — | 0.78 |
| — | 0.53 |
| — | 0.77 |
| — | 0.78 |

ROC AUC of HB-LT and baseline models