

1 **Regional-specific calibration enables application of bioinformatic evidence for clinical**  
2 **classification of 5' cis-regulatory variants in Mendelian disease**

3

4 Rehan M. Villani<sup>1</sup>, Maddison E. McKenzie<sup>1</sup>, Aimee L. Davidson<sup>1</sup>, Amanda B. Spurdle<sup>1,2</sup>

5 1 Population Health Program, QIMR Berghofer Medical Research Institute, Brisbane, Australia

6 2 University of Queensland, Brisbane Australia

7

8 **Abstract**

9 To date, clinical genetic testing and approaches to classify genetic variants in Mendelian disease genes  
10 have focused heavily on exonic coding and intronic gene regions. This multi-step study was undertaken  
11 to provide an evidence base for selecting and applying bioinformatic approaches for use in clinical  
12 classification of 5' cis-regulatory region variants. Curated datasets of rare clinically reported disease-  
13 causing 5' cis-regulatory region variants, and variants from matched genomic regions in population  
14 controls, were used to calibrate six bioinformatic tools as predictors of variant pathogenicity.  
15 Likelihood ratio estimates were aligned to code weights following ClinGen recommendations for  
16 application of the American College of Medical Genetics (ACMG)/American Society of Molecular  
17 Pathology (AMP) classification scheme. Considering code assignment across all reference dataset  
18 variants, performance was best for CADD (81.2%) and REMM (81.5%). Optimized thresholds provided  
19 moderate evidence towards pathogenicity (CADD, REMM), and moderate (CADD) or supporting  
20 (REMM) evidence against pathogenicity. Both sensitivity and specificity of prediction were improved  
21 when further categorizing variants based on location in an EPDnew-defined promoter region.  
22 Combining predictions (CADD, REMM, and location in a promoter region) increased specificity at the  
23 expense of sensitivity. Importantly, the optimal CADD thresholds for assigning ACMG/AMP codes PP3  
24 ( $\geq 10$ ) and BP4 ( $\leq 8$ ) were vastly different to recommendations for protein-coding variants (PP3  $\geq 25.3$ ;  
25 BP4  $\leq 22.7$ ); CADD  $< 22.7$  would incorrectly assign BP4 for  $> 90\%$  of reported disease-causing cis-  
26 regulatory region variants. Our results demonstrate the need to consider a tiered approach and  
27 tailored score thresholds to optimize bioinformatic impact prediction for clinical classification of cis-  
28 regulatory region variants.

29 **Keywords:** ACMG/AMP; non-coding; cis-regulatory; promoter; variant classification; genomic variant;  
30 genome; humans; mutation; sequence analysis; DNA; genetic testing; genetics; human genetics

31

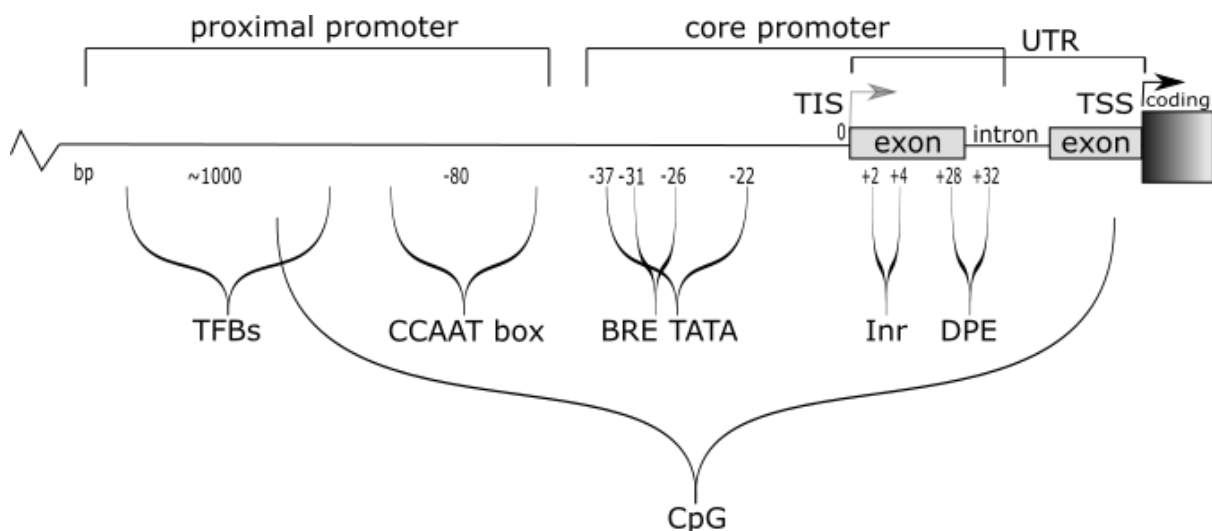
## 32 Introduction

33 Advances in genomic sequencing technology have led to dramatic improvements in diagnostic rates  
34 for inherited disease. Fundamental to these developments were the American College of Medical  
35 Genetics and Genomics and the Association for Molecular Pathology (ACMG/AMP) recommendations  
36 for clinical interpretation of genomic variants, which provided guidelines for classifying a given  
37 variant's potential role in disease <sup>(1)</sup>. To date, the overwhelming majority of variants classified as  
38 disease-causing are located in the protein-coding region of the genome <sup>(2)</sup>. The non-coding region,  
39 despite representing approximately 98% of the genome, remains largely unexplored as an explanation  
40 for Mendelian disease.

41 The non-coding sequence upstream of protein-coding genes, known as the cis-regulatory region, has  
42 important regulatory functions <sup>(3)</sup>. Variants in non-coding regions with known or suspected cis-  
43 regulatory function are thus high priority for investigating potential impact on gene function and  
44 disease predisposition.

45 Cis-regulatory regions contain a number of different functional domains (Figure 1), typically including:  
46 a core promoter which enables gene transcriptional output; a proximal promoter; and an upstream  
47 untranslated region (5' UTR). Further, the 5' UTR may contain introns that modulate gene output e.g.  
48 expression level, spatial or temporal modifications. Within these domains are identifiable cis-  
49 regulatory sequence motifs. The cis-regulatory region domains can contain: promoter motifs required  
50 for transcription initiation such as a TATA box; downstream promoter element (DPE); initiator element  
51 (Inr); or motif ten element (MTE) <sup>(4-6)</sup>. Additionally, domains such as CpG islands, CCAAT regions,  
52 regions of open chromatin and various epigenetic markers, can convey regulatory function and are  
53 enriched in promoter regions <sup>(7-11)</sup>. Finally, a diverse range of transcription factor binding motifs enable  
54 temporal and spatial gene modulation <sup>(12)</sup>. The variable composition of domains and motifs in the cis-  
55 regulatory region upstream of a gene dictate its expression and behavior. Thus, while not directly  
56 encoding protein sequence, the cis-regulatory regions proximal to the protein-coding gene sequence  
57 are crucial for normal biological function.

58



59

60 *Figure 1. Overview of features associated with cis-regulatory regions.*

61 Promoter regions contain a variety of motifs, with substantial diversity in features present and relative  
62 locations of motifs between genes. The cis-regulatory region includes the main regulatory regions  
63 upstream (5') to the translation start site (TSS), including the core and proximal promoter/s, and  
64 encompassing any untranslated regulatory introns and exons, and the transcription initiation site (TIS).  
65 The functional components of the core promoter may include a Beta recognition element (BRE), a TATA  
66 box or variation thereof (TATA), an initiator sequence (termed Inr), and/or a downstream promoter

67 element (DPE). These sequence features are generally observed within -40 bp to +40 bp of the TSS (+1  
68 position within the Inr). The cis-regulatory region also contains a variety of transcription modulating  
69 regions, including CCAATT box sites and transcription factor binding sites (TFBs) which regulate  
70 temporal and spatial control of gene expression. Not represented in this image are a number of  
71 additional cis-regulatory elements, including the E box, X box and GC sites. Cis-regulatory elements are  
72 observed in various combinations, sometimes with multiple instances of a functional element, and  
73 also are not all found concurrently.

74 Variants in cis-regulatory regions cause inherited disease through impacting gene function, generally  
75 via altering gene regulation. Disease-causing variants have been observed across the range of cis-  
76 regulatory motifs, and these variants have generally been reported to cause phenotypes similar to  
77 pathogenic variants within the associated protein-coding regions <sup>(13)</sup>. Some previously reported  
78 examples include: variants upstream of *PTEN* that reduce promoter activation causing Cowden  
79 syndrome <sup>(14)</sup>; variants in the TATA box recognition sites of *HBB* and *HBD* that alter transcription  
80 initiation by TATA Binding Protein and are reported as causal for  $\beta$ - and  $\delta$ -thalassemia <sup>(15)</sup>; deletions  
81 upstream of the TSS in the *APC* gene promoters 1A or 1B identified as causal for Familial Adenoma  
82 Polyposis <sup>(16)</sup>; and variants in the 5' UTR of *MLH1* that reduce transcription that have been reported to  
83 cause hereditary non-polyposis colorectal cancer (also known as Lynch Syndrome) <sup>(17)</sup>. Despite such  
84 examples establishing precedence, cis-regulatory region variants are not routinely examined in the  
85 clinical diagnostic setting <sup>(2)</sup>.

86  
87 Ellingford et al. <sup>(2)</sup> recently published recommendations to support the interpretation of non-coding  
88 variants in alignment with the ACMG/AMP variant classification guidelines <sup>(1)</sup>. These recommendations  
89 included a general description regarding use of bioinformatic prediction tools for non-coding variant  
90 interpretation, with reference to several tools that might be used to predict variant impact on splicing,  
91 or deleteriousness of other categories of non-coding region variants. The authors specifically  
92 highlighted the importance of accurately annotated true positive pathogenic variants for training, and  
93 cautioned against over-interpretation of output from genome-wide predictors.

94  
95 Another important consideration for regulatory region variant effect prediction is how to prioritize,  
96 compare and select bioinformatic tool/s for both calibration and ongoing use. There are numerous  
97 tools with potential relevance for impact prediction of non-coding variants (see Table S1 for examples).  
98 For ease of application in a variant curation setting, ideally bioinformatics tool/s should be: current  
99 and maintained; easy to use (if possible, even for those without coding skills); publicly available  
100 without cost; and capable of batch variant annotation. While many previous studies have compared  
101 tool performance in the process of assessing a new bioinformatic tool for non-coding regions, we  
102 identified relatively few apparently impartial reviews of computational tools that predict impact on  
103 function for non-coding variants <sup>(2, 18-25)</sup>. Of the latter, only one study <sup>(24)</sup> identified optimal thresholds  
104 for predicting impact of non-coding variants, and reported tool sensitivity and specificity using these  
105 thresholds. While sensitivity and specificity are key factors in selecting which tool/s may be used to  
106 predict variant pathogenicity, formal calibration of a tool using known pathogenic and benign variants  
107 is required to determine the appropriate evidence weight for application in clinical variant  
108 classification <sup>(26)</sup>. Bayesian modelling of the ACMG/AMP variant classification guidelines has provided  
109 a framework on how to assign evidence weights based on likelihood ratio (LR) towards pathogenicity  
110 <sup>(27)</sup>. Recently, a ClinGen computational subgroup (Pejaver et al. 2022) used this approach to define  
111 score thresholds for bioinformatic prediction evidence weighting for missense variants <sup>(28)</sup>. In this study,  
112 only four of the 13 tools assessed were potentially applicable for non-coding variants: two  
113 conservation/constraint scores (GERP and PhyloP) and two meta-predictors (CADD and BayesDel) <sup>(29-  
114 32)</sup>. Given that the mechanisms underlying cis-regulatory region function are quite different to those  
115 for protein-coding regions, we hypothesized that the score thresholds and evidence weights derived  
116 for missense variant impact cannot be assumed to be applicable for cis-regulatory region variants.

117 We applied multiple quality control and filtering steps to publicly accessible information to generate  
118 refined reference datasets of reported disease-causing variants (representing pathogenic variants) and  
119 region-matched control variants, exclusively located in cis-regulatory regions. These reference datasets  
120 were used to compare performance of, and evaluate evidence weight for, scores from six bioinformatic  
121 prediction tools and promoter region annotation. The results from this calibration study demonstrate  
122 the need to consider a tiered approach with tailored score thresholds to optimize impact prediction  
123 for clinical classification of cis-regulatory region variants.

## 124 **Methods**

### 125 **Scoping analyses relating to variant effect and location in non-coding regions**

126 A preliminary dataset of non-coding variants was sourced from ncVarDB <sup>(20)</sup>, and annotated using VEP  
127 (online GUI version 109, 28 March 2023) to obtain the Ensembl molecular consequence and CADD  
128 PHRED scores (v1.6). CADD score profiles for benign and pathogenic variants, categorized as defined  
129 by ncVarDB, were compared using a density plot. Non-coding variants were then grouped by Ensembl  
130 molecular consequence, with splicing-related molecular consequences were collapsed into a single  
131 'splicing' group, and variants with no molecular consequence annotated were collapsed into a single  
132 group 'other'. The CADD PHRED score of benign and pathogenic variants was compared for each  
133 molecular consequence group using bar graph visualization.

### 134 **Sourcing reported disease-causing cis-regulatory variants**

135 Source data included large-scale studies and variant databases <sup>(20, 33-35)</sup>, and smaller research  
136 publications (published up to February 2023) reporting Mendelian disease-causing regulatory region  
137 variants identified in the clinical setting (clinically reported and/or patient-identified). Variants that  
138 were annotated by the original source as 5' UTR, upstream or regulatory region variants were selected  
139 to generate a combined dataset of 962 variant records (Table S2). 2 variants  
140 (NC\_000001.11:g.11023351G>A, NC\_000014.9:g.75958692G>A) were excluded based on literature  
141 reporting their location as 3' UTR (though the original source annotation as 5'UTR).

142 These reported disease-causing variants (hereafter also referred to as disease variants) were  
143 investigated for literature and functional evidence via the following approaches: ClinVar (collected  
144 November 2022)<sup>(36)</sup>; dbSNP; LitVar search using rsID and/or variant location; and Google search (online  
145 search completed 6 December 2022) for variant MANE transcript associations, HGVS  
146 nomenclature/dbSNP identifiers, gene and alternate gene references, and promoter-related  
147 information. PMIDs were recorded for all publications that appeared to capture evidence specific to  
148 the variant (Table S3). After removal of duplicates, 576 unique cis-regulatory region variants remained.

### 149 **Identifying cis-regulatory regions of interest, 5kb upstream regions of MANE transcripts**

150 The translation and transcription start sites for all MANE\_Select and MANE\_Plus\_Clinical transcripts  
151 were collected using BioMart (Ensembl) <sup>(37)</sup>. The region start was determined as genomic location 5kb  
152 upstream of the transcription start site for the positive strand or 5kb downstream of the transcription  
153 start site for the negative strand. The last nucleotide 5' to the translation start site in positive strand  
154 or first nucleotide 3' to the translation start site in the negative strand was designated the region end  
155 location. A 'region of interest' input BED file <sup>(38)</sup> was then created to match the relevant genes for the  
156 cis-regulatory reported disease variants.

### 157 **Population variant frequency/conservation correlation**

158 Variants located within the regions of interest (MANE genes) were selected from gnomAD v3.0 VCF  
159 files <sup>(39)</sup>. Maximum population allele frequency (maxAF) was calculated for 314,817 variants by  
160 selecting variants based on the highest alternative AF from (non-founder) populations (Non-Finnish  
161 European, South-Asian, African-American/African ancestry, Latino, East Asian). Precomputed GERP  
162 (version homo\_sapiens GRCh38, downloaded 02 February 2023) and phyloP 100V GRCh38 vertebrate

163 (version hg38.100way.phyloP100way 2015-05-11, downloaded 02 February 2023) scores were  
164 sourced from UCSC and gnomAD variants annotated with precomputed scores via VEP (See  
165 Supplementary Methods). GERP and phyloP 100V scores were obtained for single nucleotide variants  
166 only. The correlation between maxAF and GERP, and maxAF and phyloP 100V, was investigated via  
167 scatterplot with linear regression and Spearman's correlation coefficient.

168 gnomAD variants were binned into seven groups by maxAF: [1] 0 -0.00001, [2] >0.00001-0.00002,  
169 [3] >0.00002-0.0001, [4] >0.0001-0.001, [5] >0.001-0.01, [6] >0.01-0.1, and [7] >0.1-1.

170 Summary statistics of both phyloP 100V (Table S4) and GERP (Table S5) were calculated for each of  
171 these bins, including mean and standard deviation of each maxAF bin.

172 Based on conservation metrics for the different bins, the maxAF bin [3] >0.00002-0.0001 was used as  
173 source for a presumed benign control variant set. To enable analysis via web-based annotation  
174 platforms (e.g. CADD web annotation recommends limiting to around 10 000 variants), a subset of the  
175 maxAF bin [3] was created by random selection of 10% of the total 127,868 variants. This formed the  
176 initial control variant set comprising 12,788 variants (Table S6), hereafter also referred to as control  
177 variants.

## 178 **Compilation of reference datasets**

179 To select refined reference datasets of cis-regulatory region variants for calibrating bioinformatic tools,  
180 the 576 reported disease-causing variants (Table S3) and 12,788 control variants (gnomAD population  
181 variants with maxAF >0.00002-0.0001) (Table S6) were filtered to remove variants with potential to  
182 confound the analysis. Variants were excluded from the reference datasets if they: were predicted to  
183 alter an amino acid in any transcript; overlapped with the coding region of the MANE transcript  
184 (including introns between coding exons); were predicted to alter splicing by max SpliceAI delta score  
185  $\geq 0.2$  (of which a subset had published evidence for impact on splicing, Table S7); had VEP-annotated  
186 ClinVar classification in opposition with the reference dataset grouping (i.e. disease variants with a  
187 benign classification or control variants with a pathogenic classification); were GWAS-identified with  
188 experimental evidence supporting causality for common disease<sup>(40)</sup>; or had ambiguity concerning their  
189 role in disease from a broad literature search (e.g. some variants were reported *in cis* with a second  
190 potentially causal variant). Annotations relating to all variant exclusions are shown in Table S8, and  
191 detailed description of variant exclusion methods is provided in the Supplemental methods.

## 192 **Selection of bioinformatic impact prediction tools**

193 A literature search identified 269 bioinformatic tools with potential application for non-coding variant  
194 classification (Table S1). To prioritize tools for further clinical evaluation and calibration, we selected a  
195 subset of six tools previously evaluated as highly performing in Wang et al. 2022. In addition, the  
196 EPDnew database of promoter regions (version H. sapiens 006, GRCh38)<sup>(41)</sup> was selected as a source  
197 of experimental and computationally derived promoter locations.

## 198 **Bioinformatic tool score collection and variant annotation**

199 Variant annotations from multiple sources were combined in R (version 4.2.3), further information on  
200 tools and datasets can be found in Supplemental Methods and Table S9, and a full collation of  
201 reference dataset variant annotations can be found in Table S10.

202 In summary, the annotations were as follows. VEP (online GUI version 109, 28 March 2023) was used  
203 to source RefSeq transcripts, consequence (from Ensembl), MANE\_Select, MANE\_Plus\_Clinical, Amino  
204 Acids, CLIN\_SIG (ClinVar classification) annotations. Custom VEP command line version 99.2<sup>(42)</sup> was  
205 used to collect GERP, vertebrate phyloP 100V, LINSIGHT, and Eigen annotations. CADD (v1.6), FATHMM-  
206 MKL, FATHMM-XF and REMM (V0.4) annotations were obtained via web GUI. In addition, promoter-  
207 associated "sub"-annotations were extracted from the CADD results table (web GUI sourced), including;  
208 CpG, percent CpG in a window of +/- 75bp (default: 0.02); GC percent, percent GC in a window of +/-

209 75bp (default: 0.42); RemapOverlapTF, Regulatory region map number of different transcription  
210 factors binding (default: - 0.5); and Encode Regulatory Region Feature annotations. Control and disease  
211 variant scores were compared, using default CADD score settings. For Encode regulatory features,  
212 regulatory feature overlap categories were lumped to calculate the dataset proportion overlap with  
213 any regulatory feature. Splicing prediction analysis was performed using SpliceAI <sup>(43)</sup>. Based on  
214 calibration results reported previously <sup>(44)</sup>, maximum delta score of  $\geq 0.2$  was considered as  
215 bioinformatic evidence for predicted impact on splicing. Annotation of variant location within a  
216 promoter region was extracted via locational overlap with EPDnew (version H. sapiens 006, GRCh38)  
217 using the R GenomicRanges (version 1.50.2).

218 All annotations were performed on GRCh38, except LINSIGHT and FATHMM-MKL, for which variant  
219 GRCh37 positions were determined using web-based UCSC LiftOver tool, annotations collected using  
220 GRCh37 positions and returned to the corresponding GRCh38 locations.

## 221 **Statistics and bioinformatic tool calibration**

222 Figures were generated in R (version 4.2.3)/R Studio (2023.06.01), Microsoft excel and/or Inkscape  
223 (0.92). Statistical analyses were performed using R/R Studio, including linear regression, Spearman's  
224 correlation, Wilcoxon rank sum tests, Chi-Square tests and summary statistics.

225 The overall bioinformatic tool evaluation was performed by: (i) allocating score categories; (ii)  
226 calculating the score category LR and resulting evidence category/strength; (iii) evaluating the  
227 combined performance of the categories. An online LR calculation tool developed and applied for  
228 calibration of splicing prediction tool thresholds <sup>(44)</sup> was used to determine the area under the curve  
229 (AUC), Youden's index and the score threshold corresponding to the Youden's index. Upper and lower  
230 thresholds defining score categories were determined using the score defined by Youden's index to  
231 designate the central point for an uncertain zone comprised of approximately 10% of variants.  
232 Sensitivity, specificity and accuracy of the scored variants were determined based on the defined score  
233 categories. LRs were estimated for the different bioinformatic score categories by comparison of the  
234 proportions observed for control and reported disease-causing variants, as described previously <sup>(45)</sup>.  
235 ACMG/AMP criteria weights were assigned based on LR, following published LR range/threshold  
236 recommendations <sup>(27)</sup>.

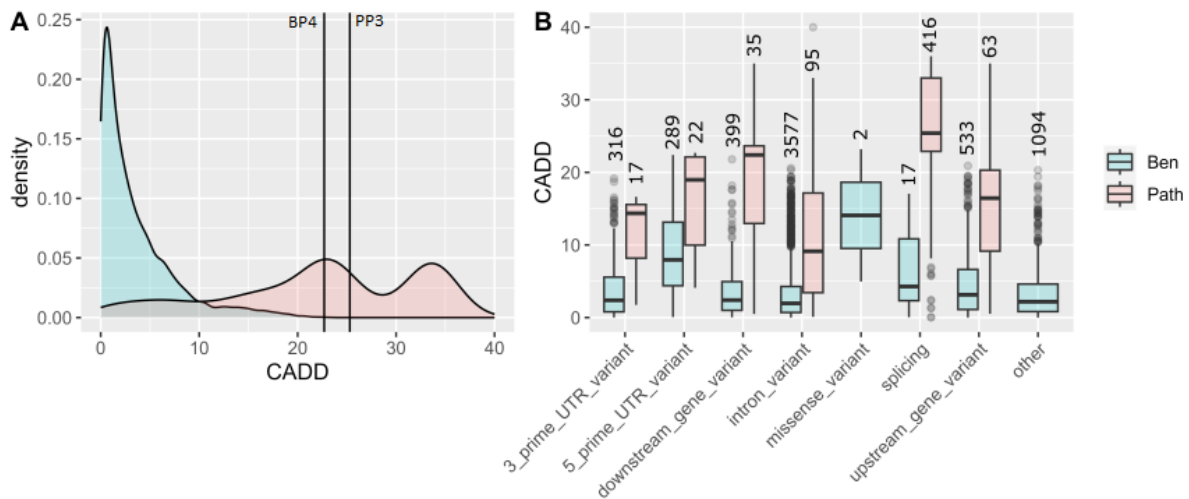
237 The evaluation of selected bioinformatics impact predictor tools was then adjusted to include all  
238 variants (including unscored and uninformative variants combined, referred to as the undetermined  
239 group), with these whole reference set evaluation results designated as the clinical performance. For  
240 clinical performance comparisons, the overall score category alignment with the reference set  
241 experimental group was determined as correct, incorrect and undetermined. Correct referred to the  
242 variant scoring in a category providing at least supporting evidence consistent with the reference set  
243 status, either towards pathogenicity for disease variants (true positive) or against pathogenicity for  
244 control variants (true negative). Incorrect referred to the variant scoring in a category providing  
245 evidence inconsistent with reference set status, either against pathogenicity for disease variants (false  
246 negative) or towards pathogenicity for control variants (false positive). Undetermined referred to both  
247 any variant that did not score (no score) and/or variants for which the respective tool did not reach an  
248 LR corresponding to at least supporting evidence either towards or against pathogenicity  
249 (uninformative).

## 250 **Results**

### 251 **Non-coding variants have distinct impact prediction score profiles**

252 Non-coding regions contain motifs with a variety of functions, therefore it should be anticipated that  
253 non-coding variants can cause impact on gene expression/function via a broad set of mechanisms. We  
254 hypothesized that this may present as large variability in variant impact prediction scores depending  
255 on specific non-coding region features, with implications for selection of appropriate thresholds for

256 assigning evidence towards or against pathogenicity. The bioinformatic prediction tool CADD was  
 257 selected for preliminary analysis based on its common use as a variant impact predictor in clinical  
 258 settings <sup>(28)</sup>, and CADD (CADD\_PHRED) scores for different Ensembl variant consequence categories  
 259 were compared for variants in the publicly accessible dataset of non-coding variants ncVarDB <sup>(20, 31)</sup>.  
 260 Variants in different locational-based consequence categories showed different CADD score profiles  
 261 (Figure 2), indicating need to consider non-coding variant location category when calibrating  
 262 bioinformatic tools for predicting clinical impact for this broad group of variants.  
 263



264  
 265 **Figure 2. CADD scores of ncVarDB variants separated by Ensembl consequence category.**  
 266 A) Density plot of CADD scores (CADD PHRED) for ncVarDB variants, comparing benign (Ben; n=7228)  
 267 and pathogenic (Path; n=721) variants. Categorization as benign or pathogenic is as per ncVarDB.  
 268 Vertical lines indicate missense variant thresholds <sup>(28)</sup>, BP4 indicates the category in which variants  
 269 would meet at least supporting level of evidence for benignity and PP3 indicating the category in  
 270 which variants would meet at least supporting level of evidence for pathogenicity according to the  
 271 missense calibrated thresholds <sup>(28)</sup>. B) CADD scores median (box center line), 25<sup>th</sup> and 75<sup>th</sup> percentiles  
 272 (upper and lower box boundaries respectively), the inter-quartile range (whisker line), with outlier  
 273 points plotted individually (dots) comparing benign and pathogenic variants, separated by Ensembl  
 274 consequence; ‘splicing’ includes grouped splicing type consequences, and ‘other’ includes those  
 275 variants that did not annotate with a molecular consequence. The number of variants is indicated  
 276 above each group.

277 Our subsequent analysis focused on the cis-regulatory region of the genome, since it is a relatively  
 278 well-studied non-coding region with a number of recognizable motifs. For the purposes of this study,  
 279 the cis-regulatory region was defined as the 5kb sequence upstream of the translation start site to the  
 280 transcription start site (the nucleotide 5’ to the ATG start site). This region spans the 5’ UTR, any  
 281 untranslated introns, the core promoter and the proximal promoter, and is an area generally  
 282 understood to regulate transcription of the neighboring gene (and would include the 5’ UTR and  
 283 upstream gene region as defined in the ncVarDB <sup>(20)</sup>).

#### 284 **Identifying a set of disease-causing cis-regulatory region variants**

285 We collated a set of 576 unique variants located in a cis-regulatory region, as defined by their source  
 286 dataset, and reported as Mendelian disease-causing (see Methods, Table S3). The variant list included  
 287 mostly single nucleotide variants (SNVs) (536, 92%), but also other small insertion/deletion variants  
 288 (44, 7.6%), including insertion, deletion and small multi-substitution variants. These 576 variants were  
 289 located in the cis-regulatory region of 317 genes (or 1523 RefSeq transcripts), meaning some variants

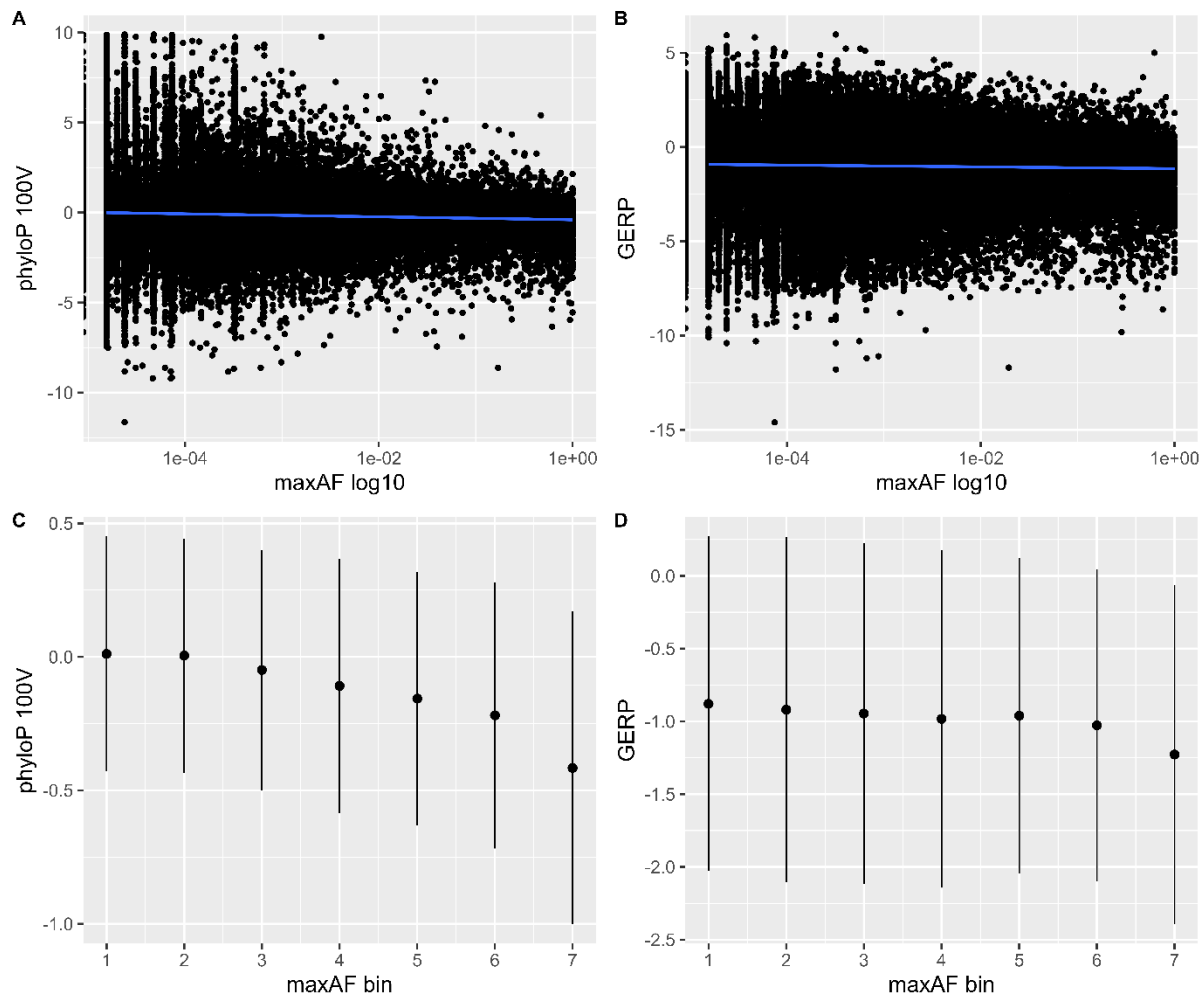
290 were in proximity to, and therefore potentially functionally relevant for, multiple genes and/or  
291 transcripts. We then: (i) selected the clinically relevant gene, the gene reported to be causal for the  
292 clinically reported condition; (ii) identified the MANE transcript for that gene; and (iii) selected the  
293 variant annotations relevant to this MANE transcript. For 14 variants where a MANE transcript was not  
294 available, the relevant transcript was identified based on the original clinical report of the variant  
295 (identified via publication, or ClinVar submission). When considering the clinically relevant gene and  
296 transcript, the 576 variants annotated to 193 genes. The revised gene and transcript-based annotation  
297 identified a range of Ensembl consequences across the 576 reported disease cis-regulatory variants,  
298 including 231 5' UTR variants, 310 upstream gene variants and 11 intron variants but also 15 splicing  
299 variants, one stop-gained variant, one start-lost, five frameshift and two missense variants (Figure S1).  
300 This observation raises the importance of considering multiple alternative mechanisms for the impact  
301 of potential "regulatory region" variants, and also highlights the need to ensure calibrations are  
302 performed on a verified set of solely cis-regulatory region variants.

### 303 **In cis-regulatory regions, population maximum allele frequency correlates with conservation**

304 Variant observation and frequency in large population datasets, such as gnomAD, is used to inform  
305 variant pathogenicity<sup>(46)</sup>. Following the ACMG/AMP classification guidelines, absence from gnomAD is  
306 considered evidence for pathogenicity (code PM2), while presence in gnomAD at a frequency higher  
307 than expected for disease prevalence provides evidence against pathogenicity (codes BA1, BS1)<sup>(1)</sup> .  
308 Previous studies have used variants with high population allele frequency (e.g. a population frequency  
309 greater than 5%) as presumed benign controls for bioinformatic tool development and evaluation<sup>(20,</sup>  
310<sup>34)</sup>. We hypothesized that high variant frequency will correlate with lower conservation, and since  
311 conservation is a key component of many bioinformatic prediction tools, selecting very common  
312 variants as controls could potentially confound tool calibration.

313 Analysis of gnomAD variants located within the cis-regulatory regions of interest (i.e. matched to those  
314 for the reported disease variants) showed an inverse correlation between maxAF and conservation  
315 score (Figure 3). Evidence for the negative correlation remained after grouping variants into seven  
316 maxAF bins, with lower conservation scores for variants observed when maxAF>0.0001 (Figure 3).  
317 Based on this information, variants with a gnomAD maxAF of >0.00002 to  $\leq$  0.0001 were considered  
318 suitable for inclusion as a control group, as this bin showed minimal conservation skewing but  
319 importantly remained within maxAF levels defined as evidence against pathogenicity from the ClinGen  
320 ENIGMA *BRCA1* and *BRCA2* Variant Curation Expert Panel VCEP (See CSpecs,  
321 <https://clinicalgenome.org/affiliation/50087/>; BS1\_Supporting may be applied for MAF >0.00002 to  
322  $\leq$  0.0001).





323

324 *Figure 3. Correlation of population allele frequency with conservation scores.*

325 The maxAF for 314,817 single nucleotide variants across 193 cis-regulatory gene regions was selected  
 326 from non-founder gnomAD populations and binned into maxAF groups. (A) Correlation scatterplot  
 327 comparing maxAF to phyloP 100V. Spearman's rank correlation coefficients were determined as rho-  
 328 value -0.05492315,  $S = 3.9247e+15$ , p-value  $< 2.2e-16$ . (B) Correlation scatterplot comparing maxAF to  
 329 GERP scores. Spearman's rank correlation coefficients were determined as rho value -0.01679457,  $S =$   
 330  $4.0028e+15$ , p-value  $< 2.2e-16$ . MaxAF was then grouped into the following maxAF bins: [1] 0-0.00001,  
 331 [2]  $>0.00001-0.00002$ , [3]  $>0.00002-0.0001$ , [4]  $>0.0001-0.001$ , [5]  $>0.001-0.01$ , [6]  $>0.01-0.1$ , and  
 332 [7]  $>0.1 - 1$ . Conservation measures were compared for variants in different maxAF bins, using (C)  
 333 mean phyloP 100V scores (line indicates 25-75% IQR) and (D) mean GERP scores (line indicates 25-75%  
 334 IQR)). The number of values per bin, and other summary characteristics, are indicated in Table S4  
 335 (phyloP 100V) and Table S5 (GERP).

336

### 337 Selection of reference variants for tool calibration

338 To select reference datasets of cis-regulatory region variants for calibrating bioinformatic tools, the  
 339 576 reported disease variants (Table S3) and 12,788 control variants (gnomAD population variants with  
 340 maxAF  $>0.00002-0.0001$ ) (Table S4) were combined and filtered to remove variants with potential to  
 341 confound the analysis. A substantial proportion of variants were excluded after applying filters:  
 342 1221/13,364 (or 9.14%) were located between the transcriptional start and transcriptional end of a  
 343 MANE transcript (including introns between coding exons); 210/13,364 or 1.57%) were predicted to  
 344 alter an amino acid (when considering any protein-coding transcripts); 52 variants, 26 each disease

345 and control, were predicted to alter splicing by max SpliceAI score  $\geq 0.2$ , of which 18 disease variants  
 346 had published evidence for impact on splicing (Table S7); 26 disease variants had a maxAF  $>0.01$  and  
 347 would be considered too common to be disease-causing; 61 variants had a reported ClinVar  
 348 classification in opposition with the reference dataset grouping (53/576 or 9.20% of disease variants  
 349 and 8/12,788 or 0.06% of control variants); one GWAS-identified variant had experimental evidence  
 350 supporting it as functionally causal for common disease <sup>(40)</sup>; and a broad literature search identified  
 351 another 19 variants where manual review identified ambiguity in disease causality (see Methods).  
 352 Details relating to all variant exclusions are shown in Table S8 and Supplemental Methods.

353 As summarized in Table 1, after application of these filters a combined cis-regulatory region reference  
 354 dataset consisting of 445 reported disease variants (representing “pathogenic” reference variants) and  
 355 9,505 control variants (representing “benign” reference variants) was compiled. This combined cis-  
 356 regulatory region reference dataset included 8,872 SNVs and 1,078 indels.

357 *Table 1. Filters applied to select reference datasets.*

	Reported disease- causing variants	gnomAD control variants
<b>STARTING set<sup>1</sup></b>	<b>576</b>	<b>12,788</b>
<b>Coding regions (between start and end of MANE transcript)</b>	47	1,174
<b>Coding (amino acid annotation, any transcript)</b>	35	175
<b>Disease-associated identified by GWAS <sup>(40)</sup></b>	1	-
<b>Clinical association unclear in literature</b>	19	-
<b>gnomAD maxAF <math>&gt;0.01</math></b>	26	-
<b>Reported disease variants with ClinVar benign classification</b>	53	-
<b>Control variants with ClinVar pathogenic classification</b>	-	8
<b>Predicted spliceogenic variants (SpliceAI max delta <math>\geq 0.2</math>)</b>	26	26
<b>Additional control variants excluded based on location in regions associated with excluded reported disease variants</b>	-	2,650
<b>FINAL refined reference dataset</b>	<b>445</b>	<b>9,505</b>

358 <sup>1</sup> Some variants met more than one criterion for exclusion.

359

### 360 **Calibration of bioinformatic tools for predicting pathogenicity of cis-regulatory region variants**

361 We selected six variant impact prediction score tools for calibration, based on their relatively high  
 362 evaluation performance in Wang 2022, CADD, REMM, FATHMM-MKL, FATHMM-XF, Eigen and LINSIGHT  
 363 <sup>(24, 31, 35, 47-50)</sup>. The distribution of prediction scores differed between control and reported disease  
 364 variants for all six bioinformatic tools analyzed (Figure 4). To determine the clinical utility of the impact  
 365 prediction scores for variant curation against ACMG/AMP recommendations for evidence weighting,  
 366 we calibrated each of these six tools using the cis-regulatory region variant reference datasets. Variant  
 367 scores were categorized into three groups based on an optimal score threshold as defined by the  
 368 Youden’s index, an upper and lower threshold were then designated to capture an intermediate,  
 369 uninformative group of approximately 10% of the variants (Table 2). Distribution of score ranges for  
 370 the reference datasets and the determined thresholds for each bioinformatic prediction score are  
 371 shown in Figure 4.

372 Considering the calibrated score categories of successfully scored variants only, all of the tools showed  
373 sensitivity above 74% and specificity above 67%. Accuracy was highest for FATHMM-MKL (87.28%,  
374 7,738/8,866) and FATHMM-XF (88.33 %, 7,309/8,275) compared to REMM (81.47%, 8,106/9,950) and  
375 CADD (81.22%, 8,081/9,949) (Table 2). However, there was a relatively large proportion of unscored  
376 variants for FATHMM-MKL (10.89%), FATHMM-XF (16.88%) and Eigen (19.84%) which do not score  
377 indels (Table 2), and an extremely high number of unscored variants for LINSIGHT (86.11%) which relies  
378 on dbSNP ID for annotation (and not genomic location/allele).

379 To evaluate bioinformatic tool performance consistent with application in a clinical diagnostic setting,  
380 sensitivity, specificity and accuracy were adjusted against a baseline of all reference set variants  
381 (scored and unscored). This evaluation revealed highest clinical accuracy for CADD (81.22%,  
382 8,081/9,950 variants) and REMM (81.47%, 8,106/9,950) (Table 2). To determine the strength of  
383 evidence provided by the calibrated score categories, the likelihood ratio (LR) associated with each  
384 score category was then calculated (Table 2). LRs estimated for the optimal score category groups are  
385 shown graphically in Figure 4. The LRs indicate that all six tools can be used to provide at least  
386 supporting evidence towards and against pathogenicity for cis-regulatory region variants. Evidence  
387 towards pathogenicity reached moderate level (LR >4.3) for CADD, FATHMM-MKL, FATHMM-XF, Eigen  
388 and REMM, and supporting level (LR >2.08) for LINSIGHT (Table 2). Evidence against pathogenicity  
389 reached moderate level (LR <0.23) for CADD, Eigen, FATHMM-MKL and LINSIGHT, and supporting level  
390 (LR <0.48) for REMM and FATHMM-XF (Table 2).

391

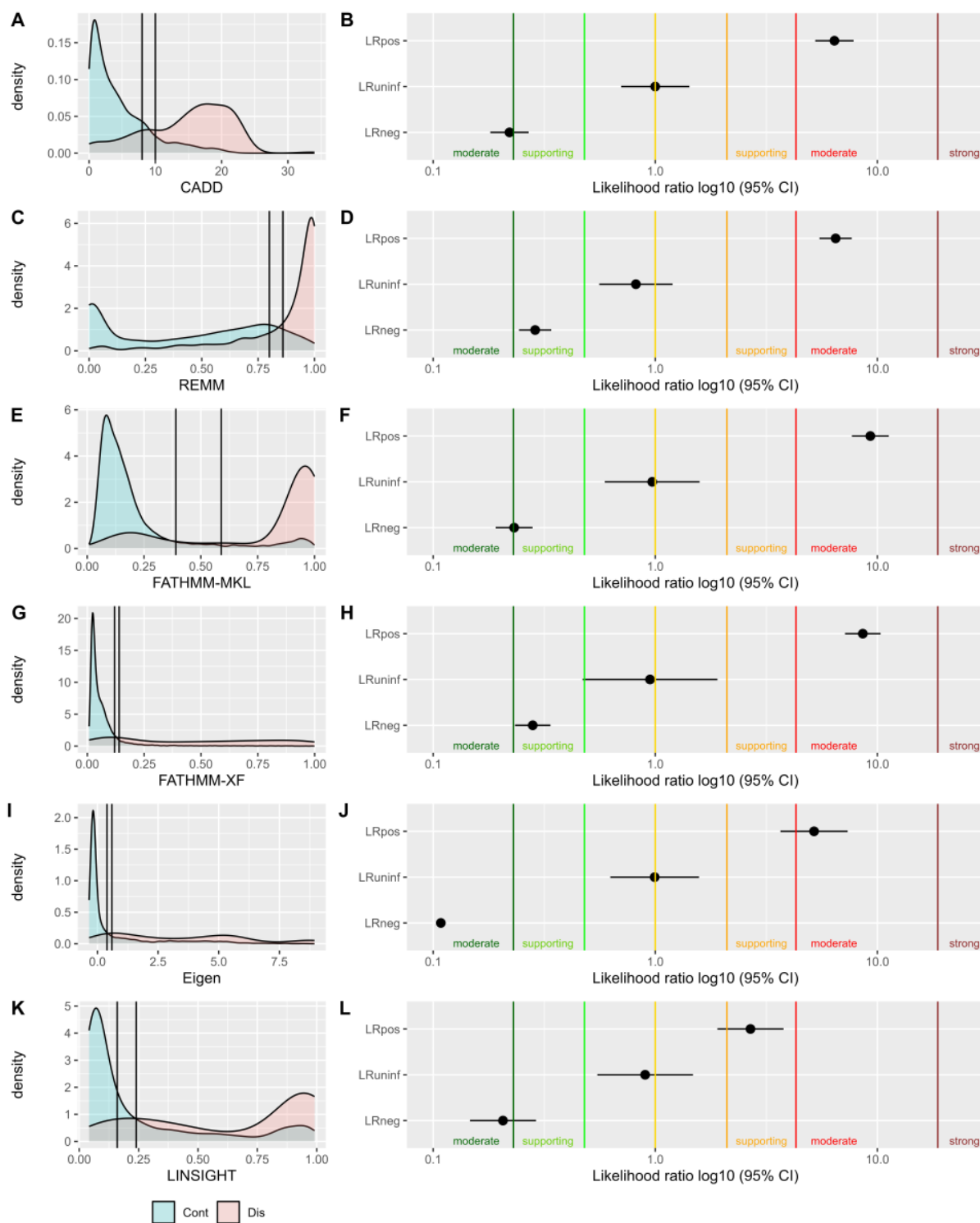
392

393 *Table 2. Bioinformatic tool performance and calibration using optimal thresholds.*

	CADD	REMM	FATHMM MKL	FATHMM XF	Eigen	LINSIGHT
<b># scored variants</b>	9949	9950	8866	8275	7975	1399
<b># control</b>	9504	9505	8452	7929	7632	1154
<b># disease</b>	445	445	413	340	343	227
<b># unscored</b>	1	0	1084	1680	1974	8568
<b>AUC</b>	0.88	0.87	0.91	0.90	0.9	0.82
<b>Youden's Index</b>	0.643	0.612	0.646	0.637	0.714	0.514
<b>Threshold</b>	8.93	0.83	0.50	0.13	0.49	0.21
<b>Lower threshold</b>	8.00	0.80	0.39	0.12	0.394	0.16
<b>Upper threshold</b>	10.00	0.86	0.59	0.14	0.594	0.24
<b>Sensitivity (% scored variants)</b>	80.72	74.94	78.84	74.40	91.08	86.19
<b>Specificity (% scored variants)</b>	87.40	88.28	91.51	91.33	82.43	67.58
<b>Accuracy (% scored variants)</b>	81.22	81.47	87.28	88.33	78.43	64.05
<b>Unscored (%)</b>	0.01	0.00	10.89	16.88	19.84	86.11
<b>Uninformative (%)</b>	6.74	7.09	3.56	2.06	4.23	1.14
<b>Undetermined (%)<sup>1</sup></b>	6.75	7.09	14.45	18.94	24.07	87.25
<b>LR negative<sup>2</sup></b>	0.22	0.29	0.23	0.28	0.11	0.21
<b>95% CI low (LR negative)</b>	0.18	0.24	0.19	0.23	0.08	0.15
<b>95% CI high (LR negative)</b>	0.27	0.34	0.28	0.34	0.15	0.29
<b>LR (uninformative)<sup>2</sup></b>	1.00	0.82	0.97	0.95	0.99	0.90
<b>95% CI low (uninformative)</b>	0.70	0.56	0.59	0.47	0.63	0.55
<b>95% CI high (uninformative)</b>	1.42	1.20	1.58	1.90	1.57	1.48
<b>LR positive<sup>2</sup></b>	6.41	6.49	9.30	8.60	5.19	2.68
<b>95% CI low (LR positive)</b>	5.25	5.49	7.68	7.15	3.66	1.90
<b>95% CI high (LR positive)</b>	7.81	7.67	11.25	10.34	7.35	3.78
<b>Clinical sensitivity (% all variants)</b>	75.28	70.56	70.34	55.51	66.52	40.67
<b>Clinical specificity (% all variants)</b>	81.49	81.98	78.12	74.30	62.69	7.52
<b>Clinical accuracy (% all variants)</b>	81.22	81.47	77.77	73.46	62.86	9.01

394 <sup>1</sup>Undetermined refers to variants that are both unscored and/or uninformative

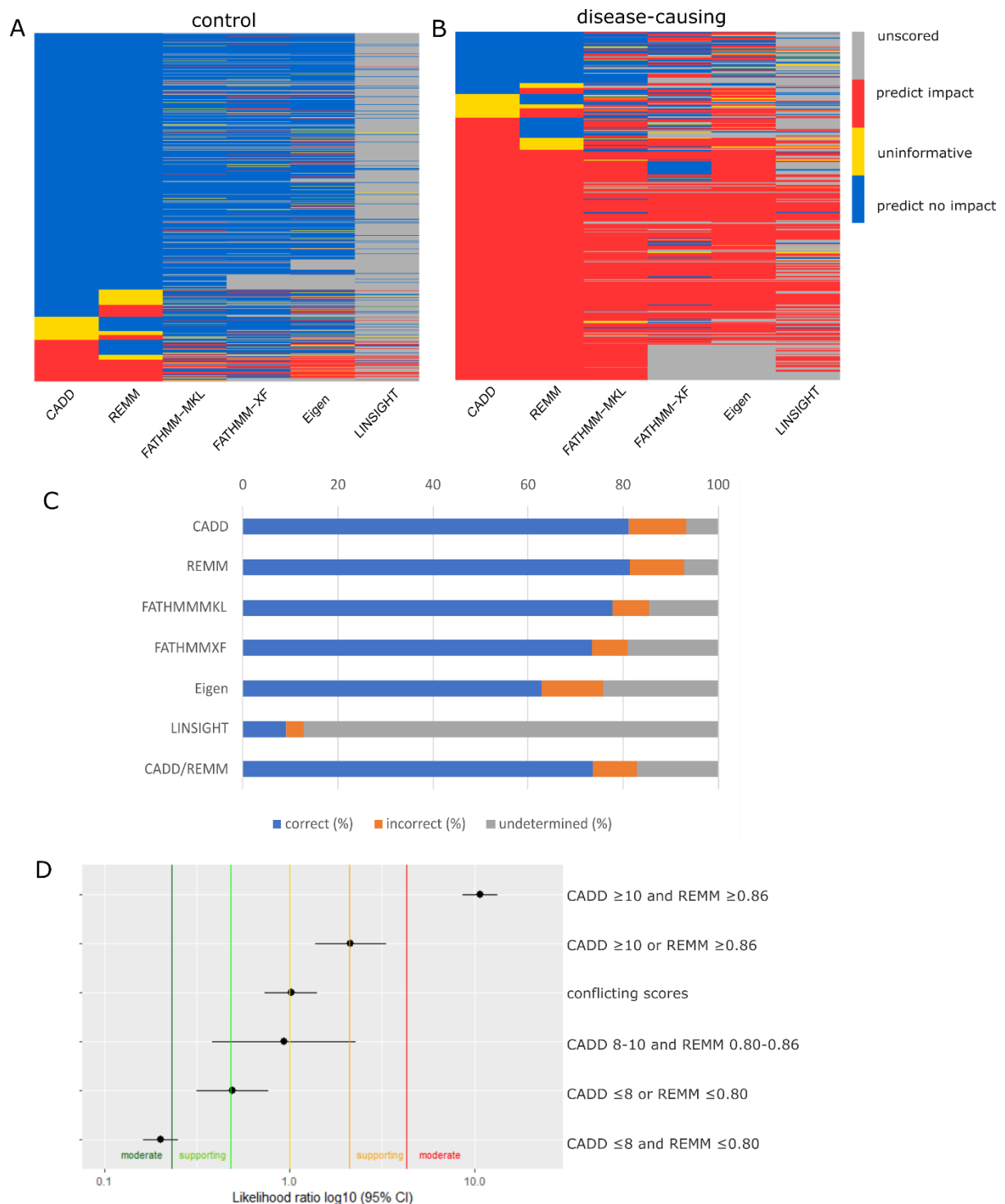
395 <sup>2</sup> LR negative refers to LR estimate for bioinformatic score range predicting no impact; LR positive  
396 refers to LR estimate for bioinformatic score range predicting impact; LR uninformative refers to LR  
397 estimate for the bioinformatic score range for variants in the middle "uninformative" category.  
398



399

400 *Figure 4. Bioinformatic tool calibration for prediction of cis-regulatory region variant pathogenicity.*  
 401 The panels on left show the distribution of scores for disease variants compared to the control  
 402 reference dataset variants, with designated optimal thresholds indicated by lines. The panels on the  
 403 right show results from bioinformatic tool calibration, with LRs for each of the three optimal categories  
 404 defined in Table 2, LRpos (LR positive) indicating the LR of the bioinformatic impact score category  
 405 predicting the variant as disease-causing, LRneg (LR negative) indicating the LR of negatively predicting  
 406 variant impact, or predicting as a control variant, and LRRuninf (LR uninformative) indicates the LR for  
 407 variants that score between categories. For all LRs, the 95% confidence interval (CI) is indicated by

408 horizontal lines. Colored vertical lines represent LR boundaries set for evidence strengths as per <sup>(27)</sup>,  
 409 with evidence strength categories indicated on the graphs. (A,B) CADD. (C,D) REMM. (E,F) FATHMM-  
 410 MKL. (G,H) FATHMM-XF. (I,J) Eigen. (K,L) LINSIGHT.  
 411



412  
 413 *Figure 5. Comparison of score prediction categories for the six tools assessed.*  
 414 Heatmap indicating score categories for (A) control reference dataset variants and (B) disease variants.  
 415 Blue, score category of variant predicts no impact (against pathogenicity); Red, score category of  
 416 variant predicts impact (towards pathogenicity); Yellow, score category of variant considered  
 417 uninformative; Grey, unscored (no score returned). (C) Comparison of overall performance of each tool  
 418 for the reported disease and control variants combined; correct (blue), referring to percentage of

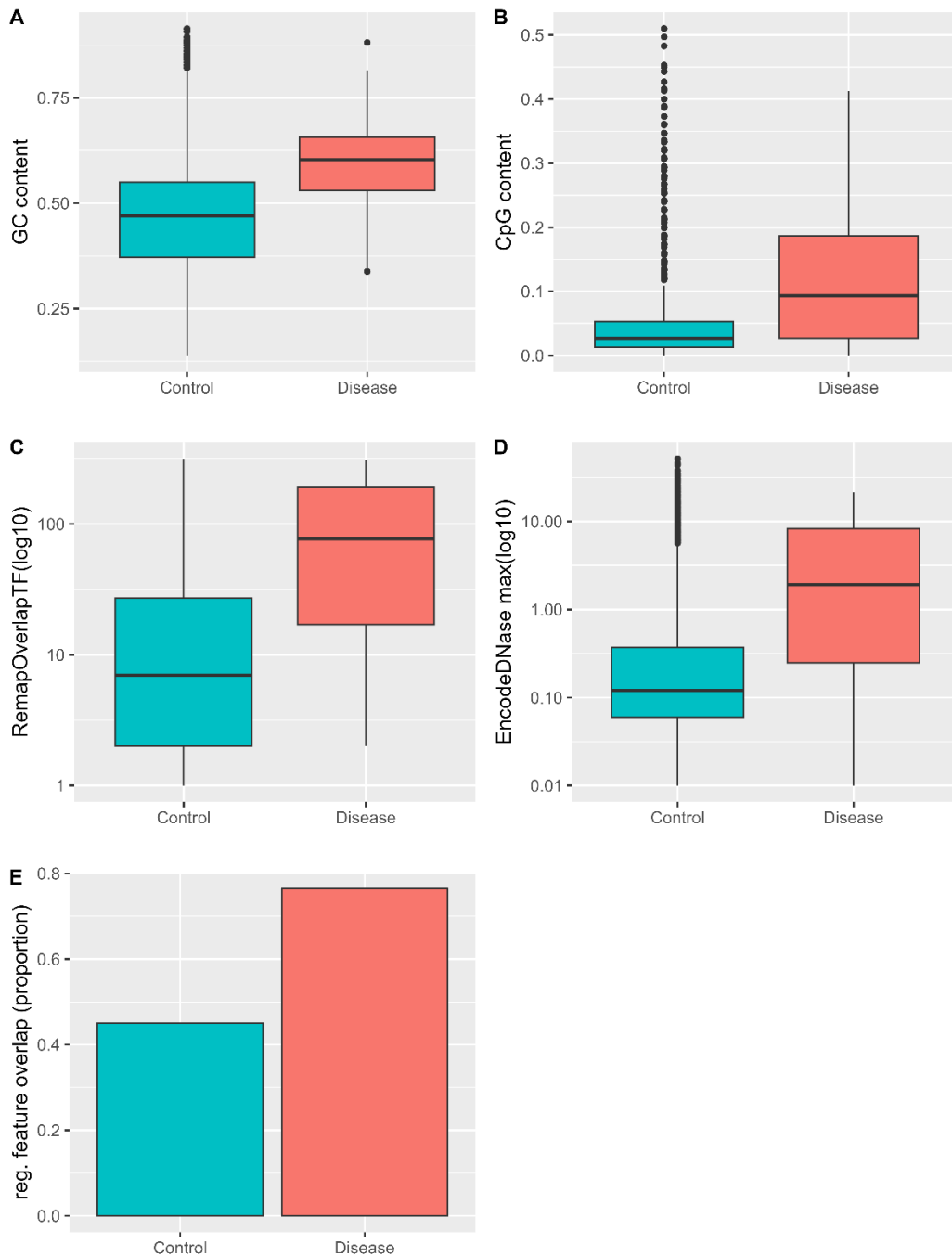
419 variants with a score category that aligns with reference dataset group, incorrect (orange), referring  
420 to percentage of variants with a score category in contradiction with reference dataset group, and  
421 undetermined (grey), referring to percentage of variants that were unscored or had scores that did not  
422 reach sufficient strength to provide classification evidence (uninformative). (D) LRs obtained when  
423 considering both CADD and REMM optimal categories combined, the specific score limits of each  
424 combined category indicated. The LR of each combinatorial category are indicated (black dot) with the  
425 95% confidence interval (CI) indicated by horizontal lines. Colored vertical lines represent LR  
426 boundaries set for evidence strengths as per <sup>(27)</sup>, with evidence strength categories indicated on the  
427 graph.

428

429 Concordance of score categories between tools was assessed to consider potential value in combining  
430 outputs of different tools for improved performance (Figure 5). This highlighted the absence of scores  
431 returned by LINSIGHT relative to the other tools, for control variants especially (Figure 5A), but also  
432 for reported disease variants (Figure 5B). CADD and REMM showed the highest concordance, while  
433 FATHMM-MKL, FATHMM-XF and Eigen showed a generally similar pattern to CADD and REMM for  
434 scored variants. Since CADD and REMM also showed the highest performance when considering  
435 accuracy for the entire dataset (representing clinical diagnostic application), we investigated if  
436 combining CADD and REMM score categories would improve prediction over use of either tool alone.  
437 As expected, combining categories increased the proportion of variants with an undetermined call,  
438 due to variants with a conflict in category assignment by the two tools (Figure 5C). As shown in Figure  
439 5D, the LR towards pathogenicity for a variant with high CADD and high REMM score category was  
440 increased (LR 10.73) compared to that for either tool alone (CADD LR 6.41, REMM LR 6.49), but  
441 remained within the moderate evidence strength range. Similarly, the LR against pathogenicity was  
442 shifted more clearly into moderate evidence for a variant with low CADD and low REMM category  
443 (CADD/REMM LR 0.20 vs CADD LR 0.22 and REMM LR 0.29, which each alone had met only supporting  
444 level of evidence).

#### 445 **Using genomic features to improve disease variant impact prediction**

446 We next assessed if specific genomic features of core-promoter regions differed between reference  
447 dataset disease or control variants, to determine if these features might be useful to improve  
448 prediction accuracy (Figure 6). Reported disease variants showed increased GC percentage, CpG  
449 percentage and TFB overlap, and higher max DNase hypersensitivity scores. Further, a considerably  
450 higher proportion of reported disease variants (75.7%) overlapped with an Ensembl regulatory feature  
451 (combining annotated regulatory elements from the Ensembl Regulatory Build <sup>(51)</sup>) compared to  
452 control variants (43.9%). The enrichment of promoter-related features in reported disease versus  
453 control variants highlighted features underlying the bioinformatic tool prediction performance, and  
454 showed the value of considering promoter region overlap for pathogenicity prediction.



455

456 *Figure 6. Promoter-related components of CADD score are enriched in disease versus control variants.*  
 457 Comparison of promoter-related CADD score component annotations for control variants (n=9,505  
 458 scored) and reported disease variants (n=445 scored). A) GC percent averages. B) CpG percent averages  
 459 (Percent GC in a window of +/- 75bp; default: 0.42). C) REmapOverlapTF average per CADD bin (Remap  
 460 number of different transcription factors binding; default: - 0.5). D) Encode DNase Hypersensitivity  
 461 max score. A) to D) Plots show median (box center line), 25<sup>th</sup> and 75<sup>th</sup> percentiles (upper and lower box  
 462 boundaries respectively), the inter-quartile range (whisker line) with outlier points plotted individually  
 463 (dots). E) The proportion of the variants in the test group overlapping with an Ensembl Regulatory  
 464 Feature (all regulatory features combined) for control variants compared to reported disease variants.

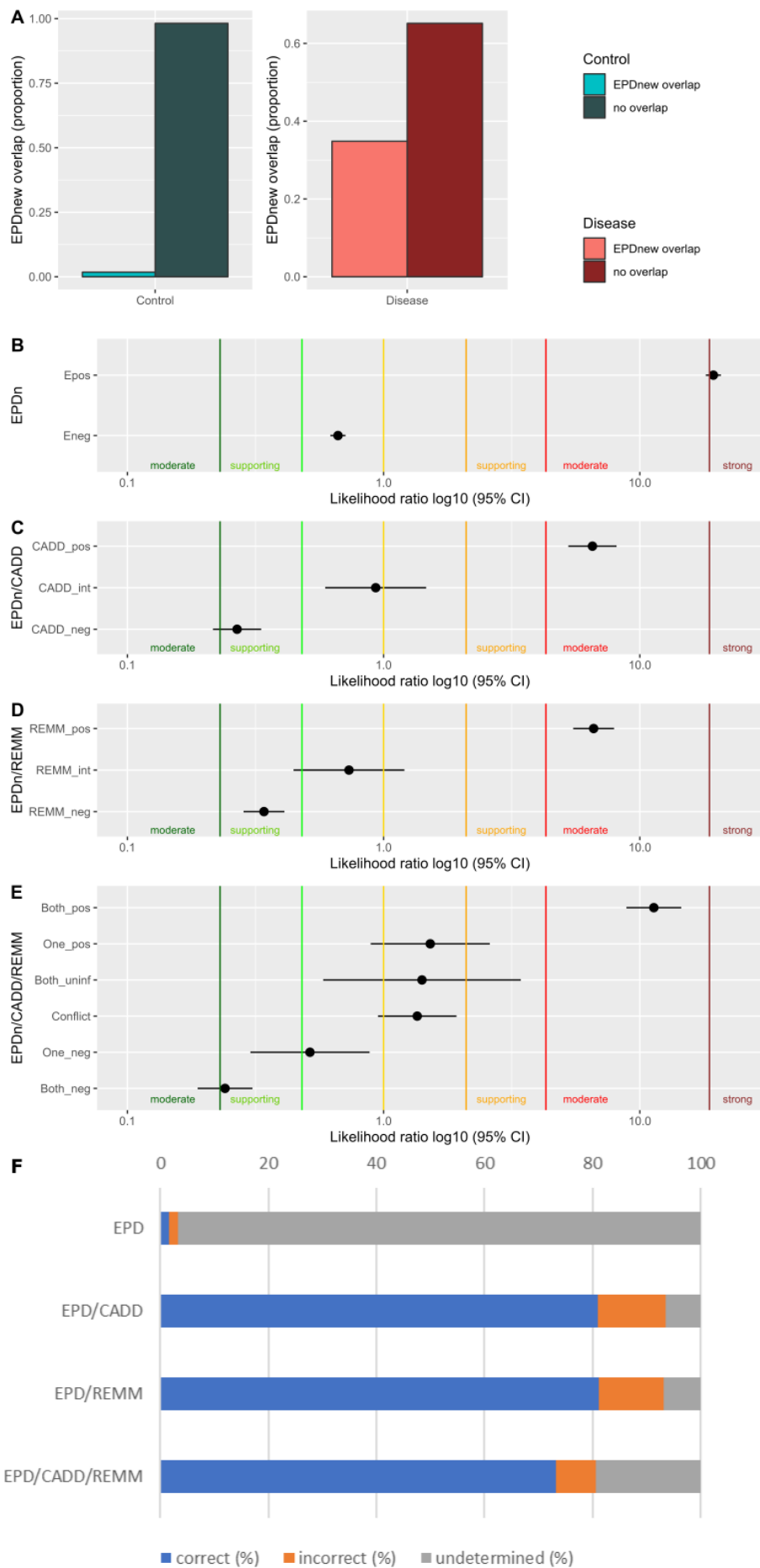


465 **Annotation of variant location within a promoter region improves pathogenicity prediction.**

466 As promoter region features were enriched in disease variants compared to controls (see Figure 6), we  
467 compared the proportion of control to disease variants located within a promoter region, as defined  
468 using the EPDnew promoter prediction database annotation <sup>(41)</sup>. A significantly greater proportion of  
469 reported disease variants (34.8%) were located within EPDnew-defined promoter regions compared  
470 to control reference dataset variants (1.8%) (Figure 7). Using EPDnew-region as the definition of  
471 promoter location, we calculated that variant location within a promoter region provides moderate  
472 evidence towards pathogenicity (LR 19.36, 95% CI 18.08-20.72); location outside of a promoter region  
473 did not reach LR thresholds required to provide evidence against pathogenicity (LR 0.66, 95% CI 0.62-  
474 0.71).

475 Based on these findings, we reassessed the evidence strength based on CADD and REMM for the  
476 subset of variants outside of the promoter region by recalculating the likelihood ratio for variants  
477 outside of EPDnew regions (Figure 7B, 7C). Reassessment showed that the impact prediction tool score  
478 thresholds calibrated based on the complete reference dataset remain appropriate for providing  
479 evidence towards and against pathogenicity for variants outside of the promoter region. For variants  
480 outside of an EPDnew promoter region, LRs for the CADD categories were:  $\leq 8$  LR 0.27, 8-10 LR 0.93  
481 (including 7% control, 4% disease variants),  $\geq 10$  LR 6.53 (Figure 7C). LRs for the REMM categories were:  
482  $\leq 0.8$  LR 0.34, 0.8-0.86 LR 0.73 (including 4% control, 4% disease variants),  $\geq 0.86$  LR 6.60 (Figure 7D).  
483 Considering CADD and REMM together for variants outside of the promoter region (Figure 7E), overall  
484 findings were similar to those for combined CADD and REMM without considering promoter region  
485 location (Figure 5). The LRs were further increased if both scores were in the high category or low  
486 category compared to using a single scores information, but there was no change in the evidence  
487 strength applicable (Table S11). The proportion of incorrect calls decreased to 7.37%, but at the  
488 expense of proportion of correct predictions (Figure 7F). Additional details are in Table S12.

489 Overall these analyses inform a process for variant annotation and bioinformatic categorization, where  
490 combining information from EPDnew and impact prediction scores into increasingly defined categories  
491 can be applied in cis-regulatory region variant classification. By first determining location in an EPDnew  
492 promoter region, followed by annotation of CADD or REMM score for variants outside of the promoter  
493 region, the process can provide evidence reaching at least supporting strength for classification of  
494 variants located within a cis-regulatory region. This combined two-step process increased the number  
495 of variants with a bioinformatic score category applicable, without compromising accuracy. While  
496 decreased sensitivity and evidence strengths based on LR estimates do not justify combined use of  
497 CADD and REMM, this might nevertheless be considered a more cautious approach in the clinical  
498 setting due to improved specificity.



500 *Figure 7. Location within a core-promoter region provides evidence towards variant pathogenicity*  
501 A) Proportion of variants located within EPDnew-defined promoter regions. Significantly more  
502 reported disease variants (n= 155/445 34.83%) than control variants (171/9,505, 1.78%) were located  
503 in a promoter region ( $\chi$ -squared = 1,550.6, df = 1, p-value <2.2e-16). B) LRs estimated for variants  
504 overlapping EPDnew promoter regions (Epos), compared to variants outside of EPDnew region (Eneg).  
505 C) LRs estimated from CADD score categories outside of EPDnew regions (EPDnew negative regions).  
506 D) LRs estimated from REMM score categories outside of EPDnew regions (EPDnew negative regions).  
507 E) LR estimates of combined CADD and REMM score categories outside of EPDnew regions. (C-E) are  
508 calculated from EPDnew location-negative variants. F) Breakdown of process accuracy for each score  
509 combination showing proportion of correctly called control and disease variants combined (blue),  
510 incorrectly called control and disease variants combined (orange), and variants with undetermined  
511 bioinformatic category, reflecting both variants unscored and for which evidence criteria thresholds  
512 were not met (grey).

513

## 514 Discussion

515 This multi-step study was undertaken to provide an evidence base for selecting and applying  
516 bioinformatic approaches for use in classification of 5' cis-regulatory region variants, in the context of  
517 Mendelian disease.

518 Analysis of existing public data highlighted the need to establish tool thresholds according to variant  
519 location and type (inferring likely molecular consequence). Further, our observation that population  
520 control frequency is negatively correlated with conservation scores informed selection of control  
521 reference dataset variants with substantially lower allele frequency. This provided a reference dataset  
522 that was not inherently enriched for lower conservation and thereby lower overall tool scores. An  
523 additional advantage to this approach to control reference dataset collection is that the bioinformatic  
524 calibration process better reflects application in the clinical variant curation setting, where variants are  
525 prioritized for more detailed curation generally after exclusion of common variants that meet  
526 ACMG/AMP population frequency codes. We also demonstrate the need for careful compilation of  
527 reported disease and control variants using various filtering strategies, particularly to ensure that  
528 reference dataset variants are exclusively located in 5' regulatory regions.

529 Our analyses showed that all six impact prediction tools assessed, when appropriately calibrated using  
530 refined reference sets, could potentially be used to inform regulatory region variant classification  
531 based on the thresholds optimized for this variant type. However, it is critical to consider the  
532 proportion of variants that will be scored by a given tool, to measure accuracy in the clinical context. The  
533 extremely high proportion of unscored variants for LINSIGHT (86.11%) would render this tool unusable  
534 in the diagnostic laboratory setting. Although FATHMM-MKL and FATHMM-XF showed the highest  
535 accuracy based on correct predictions for scored variants, they were unable to return scores for indels,  
536 which comprised 11.29% of the full variant reference dataset. When considering all reference dataset  
537 variants, REMM and CADD achieved similar clinical accuracy (CADD 81.22%, REMM 81.47%) and  
538 provided similar strengths of evidence towards and against pathogenicity. Combining CADD and REMM  
539 increased the strength of evidence both towards and against pathogenicity and resulted in fewer  
540 incorrectly assigned evidence categories, but compromised accuracy (fewer variants with correctly  
541 applied evidence). Combining scores therefore represents a cautious approach focused on minimizing  
542 false prediction of impact.

543 To facilitate the application of bioinformatic annotations for interpretation of 5' cis-regulatory region  
544 variants, we summarize in Figure 8 a staged process by which to consider the potential impact of such  
545 variants. The application of thresholds as derived from our reference datasets is considered  
546 appropriate for the interpretation of non-coding variants within the 5kb upstream and 5kb  
547 untranslated/UTR of the clinically relevant transcript. After confirmation of variant location as

548 exclusively non-coding, and without potential splicing impact, variants would be annotated for location  
549 within a promoter region using EPDnew, followed by scoring using CADD and/or REMM for variants  
550 outside an EPDnew-defined promoter region. These annotations combined may be used individually  
551 or combined to determine whether a computational code (PP3, BP4) may be applied for the variant.  
552 Although LRs derived from this study indicate that specific categories may be used to provide moderate  
553 evidence towards or against pathogenicity, we suggest a conservative approach would be to apply this  
554 bioinformatic evidence at supporting level in the first instance. The justification for such a conservative  
555 approach is that the variants identified as disease-causing to date may be biased towards those that  
556 were prioritized for functional and clinical follow-up precisely because they lay in recognizable  
557 promoter elements. Replication studies, using independent reference dataset variants, would be  
558 helpful to assess if such caution is justified.

559 To reiterate the need to calibrate thresholds considering variant location and type, we refer to a recent  
560 study calibrating tools for missense variant impact prediction, which reported that CADD scores  $\geq$   
561 25.3 provide supporting evidence towards pathogenicity (PP3) and CADD scores  $\leq 22.7$  provide  
562 supporting evidence against pathogenicity (BP4) <sup>(28)</sup>. Our findings show clearly that these thresholds  
563 are inappropriate for regulatory region variants; use of CADD  $< 22.7$  would incorrectly assign a benign  
564 prediction code for the majority of reported disease variants located in a genuine cis-regulatory region  
565 (93.5% in our final disease reference dataset). Calibration using our compiled cis-regulatory region  
566 reference datasets determined that CADD score  $\geq 10$  provides moderate evidence towards  
567 pathogenicity (LR 6.41), and score  $\leq 8$  or provides moderate evidence against pathogenicity (LR 0.22).  
568 Only 6.88% of variants would be considered “no code applicable/undetermined”. Additionally, a tiered  
569 approach combining EPDnew promoter region location with CADD and/or REMM scores enabled  
570 increased evidence strength (LR  $> 10$  rather than  $> 6$ ) for a proportion of variants, and fewer variants  
571 with incorrectly designated evidence (7% all predictions combined, rather than  $> 9\%$  via all other  
572 approaches). However, this approach comes with a compromise in terms of fewer variants assigned a  
573 bioinformatic category (84% with all predictions combined rather than  $> 90\%$  with a single tool).

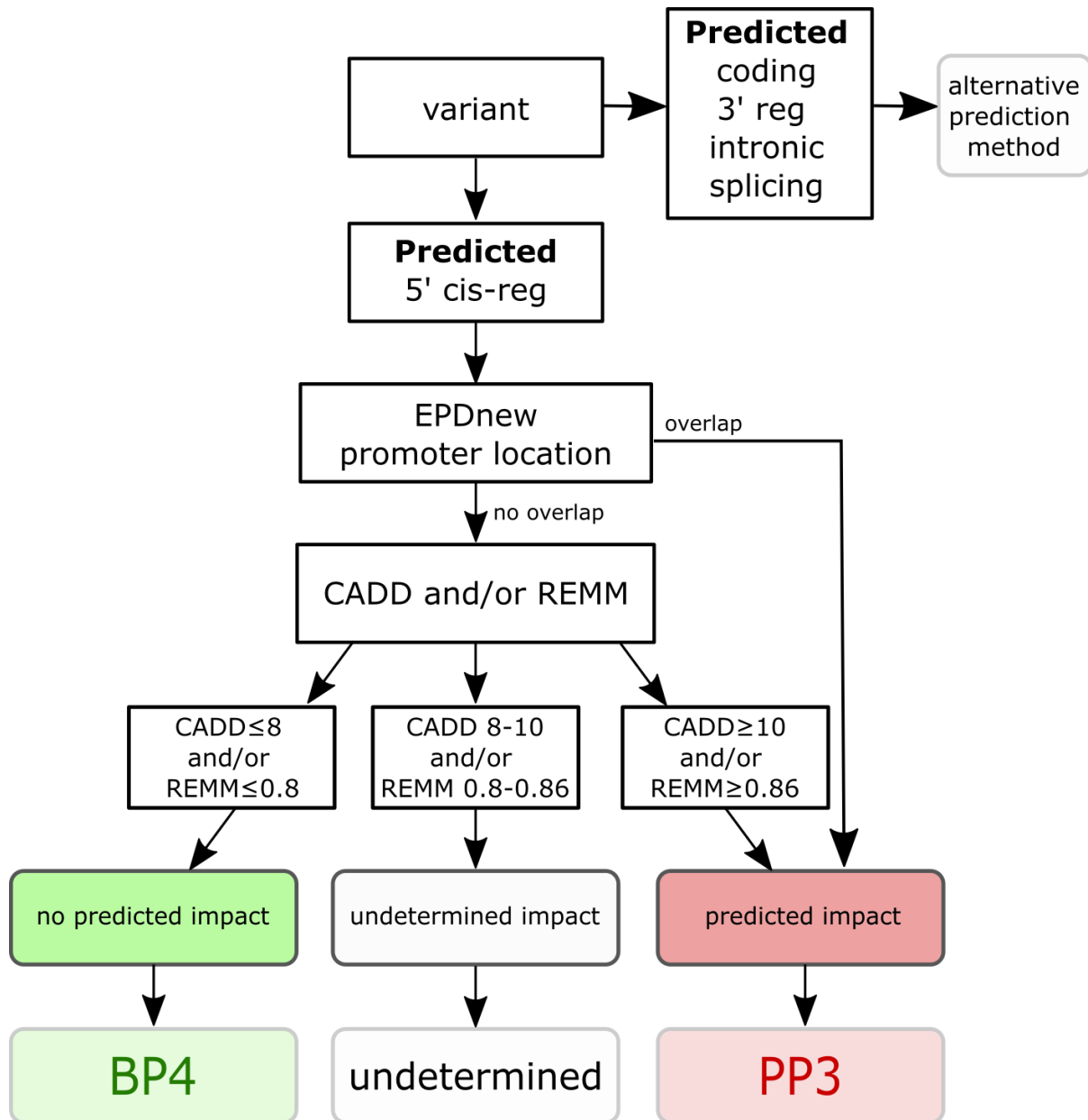
574 We stress that our study design has not provided tool calibration for regulatory region variants that  
575 also overlap in location with a coding region, for which bioinformatic score thresholds are likely to be  
576 different. Recognizing this limitation, our calibration study using carefully refined reference datasets  
577 represents an important advance for use of bioinformatic prediction evidence in the clinical  
578 classification of variants located exclusively within 5' cis-regulatory regions of Mendelian disease genes.

579

580

581

582



583

584 *Figure 8. Recommended process for assigning computational evidence of predicted impact for cis-*  
 585 *regulatory region variant classification.*

586 The presented calibration metrics are specific for use in cis-regulatory region variants. We suggest that  
 587 before applying the metrics the following be verified: the variant is definitely within the cis-regulatory  
 588 region (5kb 5' to the transcription start site to the translation start site of the clinically relevant  
 589 transcript); the variant is not also within a protein-coding region since thresholds defined here have  
 590 not been validated for variants that overlap with any coding region; the variant is not predicted to  
 591 impact splicing, a reasonably well predicted mechanism of variant impact likely to take molecular  
 592 precedence over variant impact on gene regulation. When the variant is verified as a predicted cis-  
 593 regulatory region variant, EPDnew overlap (location in a promoter region), CADD and REMM scores  
 594 can be used to determine if the variant has predicted impact, no predicted impact, or if the impact  
 595 remains undetermined (no evidence provided). Based on the categories as defined above,  
 596 computational evidence can then be used to assign at least supporting evidence for computational  
 597 ACMG/AMP code PP3 (predicted impact/towards pathogenicity) or BP4 (predicted no impact/against  
 598 pathogenicity).

599

600 **Appendices**

601 None included.

602 **Declaration of interests**

603 The authors declare no competing interests.

604 **Acknowledgments**

605 We would like to thank Daffodil Canson and Jonathan Beesley for their advice throughout results  
606 generation and manuscript preparation. RV, MM, ALD and ABS were supported by NHMRC Funding  
607 (APP177524). The work of A.L.D. was also supported in part by National Institutes of Health grant R01  
608 CA264971.

609 **Author contributions**

610 RMV, Conceptualization, Formal analysis, Methodology, Investigation, Visualization, Writing.

611 MM, Data curation, Formal analysis, Visualization, Writing.

612 ALD, Methodology, Writing, Software, Resources.

613 ABS, Conceptualization, Funding acquisition, Methodology, Writing, Supervision.

614

615 **Web resources**

616 Web-based resources and URLs are provided in Table S9.

617

618 **Data and code availability**

619 All information to replicate the findings of this study are available in the supplemental material. The  
620 code and datasets generated during this study are available through github, at cisregulatoryV,  
621 <https://github.com/ReeVee2006/cisregulatoryV>.

## 622 References

- 623 1. Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, et al. Standards and guidelines for  
624 the interpretation of sequence variants: a joint consensus recommendation of the American College  
625 of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet Med*.  
626 2015;17(5):405-24.
- 627 2. Ellingford JM, Ahn JW, Bagnall RD, Baralle D, Barton S, Campbell C, et al. Recommendations  
628 for clinical interpretation of variants found in non-coding regions of the genome. *Genome Med*.  
629 2022;14(1):73.
- 630 3. Smith M, Flodman PL. Expanded Insights Into Mechanisms of Gene Expression and Disease  
631 Related Disruptions. *Front Mol Biosci*. 2018;5:101.
- 632 4. Carninci P, Sandelin A, Lenhard B, Katayama S, Shimokawa K, Ponjavic J, et al. Genome-wide  
633 analysis of mammalian promoter architecture and evolution. *Nat Genet*. 2006;38(6):626-35.
- 634 5. Burke TW, Kadonaga JT. The downstream core promoter element, DPE, is conserved from  
635 *Drosophila* to humans and is recognized by TAFII60 of *Drosophila*. *Genes Dev*. 1997;11(22):3020-31.
- 636 6. Vo Ngoc L, Cassidy CJ, Huang CY, Duttke SH, Kadonaga JT. The human initiator is a distinct and  
637 abundant element that is precisely positioned in focused core promoters. *Genes Dev*. 2017;31(1):6-  
638 11.
- 639 7. Saxonov S, Berg P, Brutlag DL. A genome-wide analysis of CpG dinucleotides in the human  
640 genome distinguishes two distinct classes of promoters. *Proc Natl Acad Sci U S A*. 2006;103(5):1412-  
641 7.
- 642 8. Zambelli F, Pavesi G. Genome wide features, distribution and correlations of NF-Y binding  
643 sites. *Biochim Biophys Acta Gene Regul Mech*. 2017;1860(5):581-9.
- 644 9. Boyle AP, Davis S, Shulha HP, Meltzer P, Margulies EH, Weng Z, et al. High-resolution mapping  
645 and characterization of open chromatin across the genome. *Cell*. 2008;132(2):311-22.
- 646 10. Vavouri T, Lehner B. Human genes with CpG island promoters have a distinct transcription-  
647 associated chromatin organization. *Genome Biol*. 2012;13(11):R110.
- 648 11. Zoghbi HY, Beaudet AL. Epigenetics and Human Disease. *Cold Spring Harb Perspect Biol*.  
649 2016;8(2):a019497.
- 650 12. Soto LF, Li Z, Santoso CS, Berenson A, Ho I, Shen VX, et al. Compendium of human  
651 transcription factor effector domains. *Mol Cell*. 2022;82(3):514-26.
- 652 13. Phornphutkul C, Anikster Y, Huizing M, Braun P, Brodie C, Chou JY, et al. The promoter of a  
653 lysosomal membrane transporter gene, CTNS, binds Sp-1, shares sequences with the promoter of an  
654 adjacent gene, CARKL, and causes cystinosis if mutated in a critical region. *Am J Hum Genet*.  
655 2001;69(4):712-21.
- 656 14. Teresi RE, Zbuk KM, Pezzolesi MG, Waite KA, Eng C. Cowden syndrome-affected patients with  
657 PTEN promoter mutations demonstrate abnormal protein translation. *Am J Hum Genet*.  
658 2007;81(4):756-67.
- 659 15. Savinkova L, Drachkova I, Arshinova T, Ponomarenko P, Ponomarenko M, Kolchanov N. An  
660 experimental verification of the predicted effects of promoter TATA-box polymorphisms associated  
661 with human diseases on interactions between the TATA boxes and TATA-binding protein. *PLoS One*.  
662 2013;8(2):e54626.
- 663 16. Lin Y, Lin S, Baxter MD, Lin L, Kennedy SM, Zhang Z, et al. Novel APC promoter and exon 1B  
664 deletion and allelic silencing in three mutation-negative classic familial adenomatous polyposis  
665 families. *Genome Med*. 2015;7(1):42.
- 666 17. Hesson LB, Packham D, Kwok CT, Nunez AC, Ng B, Schmidt C, et al. Lynch syndrome  
667 associated with two MLH1 promoter variants and allelic imbalance of MLH1 expression. *Hum Mutat*.  
668 2015;36(6):622-30.
- 669 18. Pabinger S, Dander A, Fischer M, Snajder R, Sperk M, Efremova M, et al. A survey of tools for  
670 variant analysis of next-generation genome sequencing data. *Brief Bioinform*. 2014;15(2):256-78.
- 671 19. Drubay D, Gautheret D, Michiels S. A benchmark study of scoring methods for non-coding  
672 mutations. *Bioinformatics*. 2018;34(10):1635-41.

- 673 20. Biggs H, Parthasarathy P, Gavryushkina A, Gardner PP. ncVarDB: a manually curated database  
674 for pathogenic non-coding variants and benign controls. Database (Oxford). 2020;2020.
- 675 21. Gloss BS, Dinger ME. Realizing the significance of noncoding functionality in clinical  
676 genomics. *Exp Mol Med*. 2018;50(8):1-8.
- 677 22. Rojano E, Seoane P, Ranea JAG, Perkins JR. Regulatory variants: from detection to predicting  
678 impact. *Brief Bioinform*. 2019;20(5):1639-54.
- 679 23. Kuksa PP, Greenfest-Allen E, Cifello J, Ionita M, Wang H, Nicaretta H, et al. Scalable  
680 approaches for functional analyses of whole-genome sequencing non-coding variants. *Hum Mol*  
681 *Genet*. 2022;31(R1):R62-R72.
- 682 24. Wang Z, Zhao G, Li B, Fang Z, Chen Q, Wang X, et al. Performance comparison of  
683 computational methods for the prediction of the function and pathogenicity of non-coding variants.  
684 *Genomics Proteomics Bioinformatics*. 2022.
- 685 25. Tabarini N, Biagi E, Uva P, Iovino E, Pippucci T, Seri M, et al. Exploration of Tools for the  
686 Interpretation of Human Non-Coding Variants. *Int J Mol Sci*. 2022;23(21).
- 687 26. Wilcox EH, Sarmady M, Wulf B, Wright MW, Rehm HL, Biesecker LG, et al. Evaluating the  
688 impact of in silico predictors on clinical variant classification. *Genet Med*. 2022;24(4):924-30.
- 689 27. Tavgigian SV, Greenblatt MS, Harrison SM, Nussbaum RL, Prabhu SA, Boucher KM, et al.  
690 Modeling the ACMG/AMP variant classification guidelines as a Bayesian classification framework.  
691 *Genet Med*. 2018;20(9):1054-60.
- 692 28. Pejaver V, Byrne AB, Feng BJ, Pagel KA, Mooney SD, Karchin R, et al. Calibration of  
693 computational tools for missense variant pathogenicity classification and ClinGen recommendations  
694 for PP3/BP4 criteria. *Am J Hum Genet*. 2022;109(12):2163-77.
- 695 29. Davydov EV, Goode DL, Sirota M, Cooper GM, Sidow A, Batzoglou S. Identifying a high  
696 fraction of the human genome to be under selective constraint using GERP++. *PLoS Comput Biol*.  
697 2010;6(12):e1001025.
- 698 30. Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A. Detection of nonneutral substitution rates  
699 on mammalian phylogenies. *Genome Res*. 2010;20(1):110-21.
- 700 31. Rentzsch P, Schubach M, Shendure J, Kircher M. CADD-Splice-improving genome-wide variant  
701 effect prediction using deep learning-derived splice scores. *Genome Med*. 2021;13(1):31.
- 702 32. Feng BJ. PERCH: A Unified Framework for Disease Gene Prioritization. *Hum Mutat*.  
703 2017;38(3):243-51.
- 704 33. Landrum MJ, Lee JM, Benson M, Brown GR, Chao C, Chitipiralla S, et al. ClinVar: improving  
705 access to variant interpretations and supporting evidence. *Nucleic Acids Res*. 2018;46(D1):D1062-D7.
- 706 34. Caron B, Luo Y, Rausell A. NCBoost classifies pathogenic non-coding variants in Mendelian  
707 diseases through supervised learning on purifying selection signals in humans. *Genome Biol*.  
708 2019;20(1):32.
- 709 35. Smedley D, Schubach M, Jacobsen JOB, Kohler S, Zemojtel T, Spielmann M, et al. A Whole-  
710 Genome Analysis Framework for Effective Identification of Pathogenic Regulatory Variants in  
711 Mendelian Disease. *Am J Hum Genet*. 2016;99(3):595-606.
- 712 36. Davidson AL, Kondrashova O, Leonard C, Wood S, Tadini E, Hollway GE, et al. Analysis of  
713 hereditary cancer gene variant classifications from ClinVar indicates a need for regular reassessment  
714 of clinical assertions. *Hum Mutat*. 2022;43(12):2054-62.
- 715 37. Smedley D, Haider S, Ballester B, Holland R, London D, Thorisson G, et al. BioMart--biological  
716 queries made easy. *BMC Genomics*. 2009;10:22.
- 717 38. Niu YN, Roberts EG, Denisko D, Hoffman MM. Assessing and assuring interoperability of a  
718 genomics file format. *Bioinformatics*. 2022;38(13):3327-36.
- 719 39. Chen S, Francioli LC, Goodrich JK, Collins RL, Kanai M, Wang Q, et al. A genome-wide  
720 mutational constraint map quantified from variation in 76,156 human genomes. *bioRxiv*.  
721 2022:2022.03.20.485034.



- 722 40. Alsheikh AJ, Wollenhaupt S, King EA, Reeb J, Ghosh S, Stolzenburg LR, et al. The landscape of  
723 GWAS validation; systematic review identifying 309 validated non-coding variants across 130 human  
724 diseases. *BMC Med Genomics*. 2022;15(1):74.
- 725 41. Dreos R, Ambrosini G, Groux R, Cavin Perier R, Bucher P. The eukaryotic promoter database  
726 in its 30th year: focus on non-vertebrate organisms. *Nucleic Acids Res*. 2017;45(D1):D51-D5.
- 727 42. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, et al. The Ensembl Variant Effect  
728 Predictor. *Genome Biol*. 2016;17(1):122.
- 729 43. Jaganathan K, Kyriazopoulou Panagiotopoulou S, McRae JF, Darbandi SF, Knowles D, Li YI, et  
730 al. Predicting Splicing from Primary Sequence with Deep Learning. *Cell*. 2019;176(3):535-48 e24.
- 731 44. Walker LC, Hoya M, Wiggins GAR, Lindy A, Vincent LM, Parsons MT, et al. Using the  
732 ACMG/AMP framework to capture evidence related to predicted and observed impact on splicing:  
733 Recommendations from the ClinGen SVI Splicing Subgroup. *Am J Hum Genet*. 2023;110(7):1046-67.
- 734 45. Parsons MT, Tadini E, Li H, Hahnen E, Wappenschmidt B, Feliubadalo L, et al. Large scale  
735 multifactorial likelihood quantitative analysis of BRCA1 and BRCA2 variants: An ENIGMA resource to  
736 support clinical variant classification. *Hum Mutat*. 2019;40(9):1557-78.
- 737 46. Davidson AL, Leonard C, Koufariotis LT, Parsons MT, Hollway GE, Pearson JV, et al.  
738 Considerations for using population frequency data in germline variant interpretation: Cancer  
739 syndrome genes as a model. *Hum Mutat*. 2021;42(5):530-6.
- 740 47. Shihab HA, Rogers MF, Gough J, Mort M, Cooper DN, Day IN, et al. An integrative approach to  
741 predicting the functional effects of non-coding and coding sequence variation. *Bioinformatics*.  
742 2015;31(10):1536-43.
- 743 48. Rogers MF, Shihab HA, Mort M, Cooper DN, Gaunt TR, Campbell C. FATHMM-XF: accurate  
744 prediction of pathogenic point mutations via extended features. *Bioinformatics*. 2018;34(3):511-3.
- 745 49. Ionita-Laza I, McCallum K, Xu B, Buxbaum JD. A spectral approach integrating functional  
746 genomic annotations for coding and noncoding variants. *Nat Genet*. 2016;48(2):214-20.
- 747 50. Huang YF, Gulko B, Siepel A. Fast, scalable prediction of deleterious noncoding variants from  
748 functional and population genomic data. *Nat Genet*. 2017;49(4):618-24.
- 749 51. Zerbino DR, Wilder SP, Johnson N, Juettemann T, Flicek PR. The ensembl regulatory build.  
750 *Genome Biol*. 2015;16(1):56.

751