# ShapeGTB: The role of local DNA shape in prioritization of functional variants in human promoters with machine learning

**Maja Malkowska** [1] , **Julian Zubek** [2] , **Dariusz Plewczynski** [Corresp., 2, 3] , **Lucjan S Wyrwicz** [Corresp. 1]

[1] Laboratory of Bioinformatics and Biostatistics, Maria Sklodowska-Curie Memorial Cancer Centre and Institute of Oncology, Warsaw, Poland

[2] Laboratory of Functional and Structural Genomics, Centre of New Technologies, University of Warsaw, Warsaw, Poland

[3] Faculty of Mathematics and Information Science, Warsaw University of Technology, Warsaw, Poland

Corresponding Authors: Dariusz Plewczynski, Lucjan S Wyrwicz
Email address: d.plewczynski@cent.uw.edu.pl, lucjan.wyrwicz@coi.pl

Motivation: The identification of functional sequence variations in regulatory DNA regions is one of the major challenges of modern genetics. Here, we report results of a combined multifactor analysis of properties characterizing functional sequence variants located in promoter regions of genes.

Results: We demonstrate that GC-content of the local sequence fragments and local DNA shape features play significant role in prioritization of functional variants and outscore features related to histone modifications, transcription factors binding sites, or evolutionary conservation descriptors. Those observations allowed us to build specialized machine learning classifier identifying functional SNPs within promoter regions – ShapeGTB. We compared our method with more general tools predicting pathogenicity of all non-coding variants. ShapeGTB outperformed them by a wide margin (AUC ROC 0.97 vs. 0.57-0.59). On the external validation set based on ClinVar database it displayed only slightly worse performance (AUC ROC 0.92 vs. 0.74-0.81). Such results suggest unique characteristics of mutations located within promoter regions and are a promising signal for the development of more accurate variant prioritization tools in the future.

Availability and implementation: The datasets and source code are publicly available at: https://github.com/zubekj/ShapeGTB.

# ShapeGTB: an analysis of the local DNA shape importance in the exploration of predictive features for accurate classification of functional variants in human promoters.

Maja Malkowska[1#], Julian Zubek[2#], Dariusz Plewczynski[2,3], Lucjan Wyrwicz[1]

[1]Laboratory of Bioinformatics and Biostatistics, Maria Sklodowska-Curie Memorial Cancer Centre and Institute of Oncology, W.K. Roentgena 5, 02-781 Warsaw, Poland
[2]Laboratory of Functional and Structural Genomics, Center of New Technologies, University of Warsaw, Banacha 2C, 02-097 Warsaw, Poland
[3]Faculty of Mathematics and Information Science, Warsaw University of Technology, Warsaw, Poland
[#]contributed equally

## ABSTRACT

<u>Motivation:</u> The identification of functional sequence variations in regulatory DNA regions is one of the major challenges of modern genetics. Here, we report results of a combined multifactor analysis of properties characterizing functional sequence variants located in promoter regions of genes.

<u>Results:</u> We demonstrate that GC-content of the local sequence fragments and local DNA shape features play significant role in prioritization of functional variants and outscore features related to histone modifications, transcription factors binding sites, or evolutionary conservation descriptors. Those observations allowed us to build specialized machine learning classifier identifying functional SNPs within promoter regions – ShapeGTB. We compared our method with more general tools predicting pathogenicity of all non-coding variants. ShapeGTB outperformed them by a wide margin (average precision 0.93 vs. 0.47-0.55). On the external validation set based on ClinVar database we observed that all methods decreased their performance (average precision 0.47 vs. 0.23-0.42). Such results suggest unique characteristics of mutations located within promoter regions and are a promising signal for the development of more accurate variant prioritization tools in the future.

<u>Availability and implementation:</u> The datasets and source code are publicly available at: https://github.com/zubekj/ShapeGTB.

<u>Contacts:</u> lwyrwicz@coi.pl or d.plewczynski@cent.uw.edu.pl

<u>Supplementary information:</u> Supplementary data are available at PeerJ online.

## INTRODUCTION

The concept of personalized medicine has made the functional annotation of genomic variations one of the major goals of human genetics. The research inquiries are done both at individual level of low-throughput methods and large-scale population studies. The results of genome-wide association studies (GWAS) of complex human traits have exposed enrichment for variations in the regulatory elements, such as promoters, enhancers, insulators or intergenic regions. Although about 90% of single nucleotide polymorphisms (SNPs) are located in non-coding regions of human genome, the knowledge about their role in pathology of diseases is limited. In this article, we propose a method for functional prioritization of variants in human promoters, which represent around 1% of all SNPs identified by the 1000 Genomes Project (Ignatieva et al. 2014).

In recent years, several computational methods have been developed to address the challenging

46 task of noncoding variants annotation. These methods differ in the adopted algorithms and utilized
47 data. The main three approaches used by currently available tools are: functional annotations,
48 sequence homology analysis and machine learning models integrating information from both
49 sources. Especially the third integrating machine learning approach is worth investigating. The last
50 decade has brought dramatic progress in application of machine learning algorithms in
51 computational biology. Their versatile predictions have been utilized to link noncoding variations
52 properties to their functional nature by i.e. genome-wide annotation of variants (GWAVA)
53 (Ritchie, et al., 2014), combined annotation-dependent depletion (CADD) (Kircher, et al., 2014),
54 deleterious annotation of genetic variants using neural networks (DANN) (Quang, et al., 2015),
55 FATHMM-MKL (Shihab, et al., 2015), deltaSVM (Lee, et al., 2015), DeepSEA (Zhou and
56 Troyanskaya, 2015).
57 Promoters are one of the key regulatory elements of transcription initiation. Several resources
58 indicate that promoter regions show distinct structural constrains when compared with non-
59 promoters (Kanhere and Bansal, 2005; Goni, et al., 2007; Morey, et al., 2011; Gan, et al., 2012).
60 The analysis by Freeman et al. shows that the sequence-dependent shape of DNA encodes histone
61 affinity and dominates molecular recognition in the problem of nucleosome positioning (Freeman,
62 et al., 2014). Since various DNA sequences can encode similar shapes (Gardiner, et al., 2004;
63 Greenbaum, et al., 2007), correlation between DNA shape descriptors and biological functions
64 becomes an interesting problem to investigate.
65 The development of DNAshape web server by Zhou et al. (Zhou, et al., 2013) allowed analyzing
66 DNA structural features on a genomic scale. The method computes four DNA shape features:
67 minor groove width (MGW), roll (Roll), propeller twist (ProT) and helix twist (HelT). Recent
68 studies have showed that combining DNA sequence with DNA local shape improves the prediction
69 accuracy of transcription binding sites in vitro (Rohs, et al., 2009; Dror, et al., 2014). Here, we
70 address the question of the usefulness of such data in predicting functional effects of sequence
71 variations in promoter regions of genes. We are convinced that the DNA shape features may
72 largely contribute to solving a demanding problem of regulatory variants interpretation and
73 assessment of their effects on disease pathology.
74 To test this hypothesis and demonstrate its applicability, we trained a machine learning classifier,
75 which uses local shape to predict functional prioritization of promoter sites. In this paper, we
76 compare structural predictor's performance with sequence-based methods, and analyze in detail the
77 statistical relevance of different types of features characterizing DNA molecule.
78 In the light of the unique promoter characteristics, inclusive GC distribution (Lenhard, et al., 2012;
79 Andersson, et al., 2014), transcription factor binding site composition (Rada-Iglesias, et al., 2011;
80 Shen, et al., 2012; Thurman, et al., 2012) and unique chromatin signatures (Heintzman, et al., 2007;
81 Hon, et al., 2009), we focused our analysis on the regions located upstream of the transcription start
82 site. To our best knowledge previously developed methods have not aimed the variant prioritization
83 in promoter regions by local DNA shape features but rather focused on non-coding sequence
84 variations without acknowledging genomic region.
85
86 **MATERIALS AND METHODS**
87 <u>Datasets</u>
88 To obtain the positive dataset we used single-nucleotide variants (SNVs) annotated as regulatory
89 mutations in The Human Gene Mutation Database (HGMD®) professional version (release
90 2016.2) within 5 kilobases (kb) upstream from the annotated transcription start sites (TSS) and
91 provided sequences (Stenson, et al., 2014). The total number of experimentally validated disease-

92   related variants in our dataset is equal to 1772. The control dataset contains SNVs from the 1000
93   Genomes Project (The 1000 Genomes Project Consortium, 2015) with a global minor allele
94   frequency ≥1%. The overlapping elements of both sets were removed. Only variants lying within 5
95   kb upstream of TSS were selected for further analysis (Rosenbloom, et al., 2015). The sequences of
96   neutral motifs (not associated with disease phenotype) were retrieved from Ensembl with BioMart
97   (Kinsella, et al., 2011). The total number of negative examples in our dataset is equal to 3806. We
98   ensured that positive and negative motif sets are matched in their basic properties (Kolmogorov–
99   Smirnov two sample test results for GC-content distributions are as follows D-statistic=0.02, p-
100  value=0.48, null hypothesis of identical distributions retained). Distributions of TSS distances in
101  the two sets differed, but we made sure that it does not affect obtained results (see Supplementary
102  Material 5).
103  Machine learning pipeline
104  We split the available data into training and test sets randomly keeping the ratio 8:2. Full training
105  set contained 1417 positives and 3045 negatives, full test set contained 355 positives and 761
106  negatives. Training set was used to build feature ranking, train classifiers and optimize their
107  parameters, while test set was left for final validation and for comparison with other prediction
108  methods. To validate our methods internally on the training set we used a cross-validation strategy
109  in which in each fold SNPs from a single chromosome formed test set and SNPs from other
110  chromosomes formed training set. This eliminated possibility of overfitting during parameter
111  tuning and feature selection procedures, and additionally demonstrated whether our method
112  generalizes across different chromosomes.
113  We applied Monte Carlo feature selection (MCFS) algorithm (Draminski, et al., 2008) to perform
114  feature importance ranking. It is a universal feature selection strategy combining random subspace
115  methods with decision trees. A random subset of the original features is drawn in each iteration of
116  the algorithm and an equivalent of random forest is induced using the selected variables. Feature
117  importance ranking is constructed based on all induced trees. Additionally, meaningful
118  interdependencies between features are discovered by calculating how often two features are used
119  together to predict the class value. MCFS aims at finding all features relevant for the classification
120  task, and it guarantees that with sufficient number of iterations all features can be tested. Following
121  general guidelines by the authors of the algorithm, we set the number of iterations to 1000 and the
122  subset of original features considered in each iteration to 0.25.
123  In the classification task gradient tree boosting was used (GTB) – a popular tree-based ensemble
124  algorithm (Friedman and Meulman, 2003). It is known to perform very well in many domains,
125  often outperforming methods such as random forest, support vector machines or neural networks
126  (Sheridan, et al., 2016; Ladds, et al., 2016; Babajide Mustapha and Saeed, 2016). The key idea
127  behind GTB is to build trees sequentially, training a tree at each step to explain the prediction error
128  made by the combination of existing trees. Usually the trees are regularized to prevent overfitting.
129  We used the state-of-art implementation provided by XGBoost library (Chen and Guestrin, 2016).
130  Through cross-validation performed on the training set we selected optimal parameter values
131  (number of trees – 300, maximal tree depth – 8, learning rate – 0.1).
132  Comparison with existing approaches
133  Presently, the field of prediction and prioritization of human noncoding regulatory variants still
134  lacks a large, independent and publicly available gold-standard dataset for training, testing and
135  validating existing *in silico* approaches. The comparison of our method to the current state-of-the-
136  art methods is hampered even further by different aims and objectives. To our best knowledge all
137  available tools were designed for genome-wide, regulatory variants prioritization and there are no

138 computational methods focused on promoter regions. Nonetheless, we compared performance of
139 our algorithm with other tools on our own hold-out test set and on independent high-quality data
140 from ClinVar database (Jan 5, 2017 release) after excluding variants present in our training data
141 (Landrum, et al., 2016). Our hold-out test set contained 355 positives from HGMD and 761
142 negative examples from 1000 Genomes Project. External validation set contained 32 positive
143 examples labeled as pathogenic in ClinVar database and 761 negative examples from 1000
144 Genomes Project (not present in our train set).
145
146
147 <u>Features groups</u>
148 We used the following feature groups to annotate each SNV in our pathogenic and control datasets
149 (more detailed description can be found in Supplementary material 1 and 4):
150 *1. DNA sequence* (52 variables): 9-nt sequence motifs centered on the mutated nucleotide. The
151 sequence was encoded using 4-bits binary coding. Additional 12 binary (4-nt by 3 mutations)
152 variables indicated what type of mutation occurred (e.g. $A \rightarrow C$, $G \rightarrow T$, etc.).
153 *2. Local DNA shape features* (88 variables): helix twist, minor groove width, propeller twist, roll
154 values in span of 9 nt. Differences (*_diff*) between reference and mutated scores were added as
155 additional features.
156 *3. GC-content* (8 variables): *GC-content* in span of 7- and 9-nt for reference and mutated sequences
157 separately. Differences between the reference and mutated scores were added as additional
158 features.
159 *4. Histone modifications* (38 variables): ChIP-seq data for histone 3 lysine 9 acetylation (H3K9ac)
160 and histone 3 lysine 4 trimethylation (H3K4me3) across 16 cell lines from ENCODE (Ram, et al.,
161 2011). For H3K9ac, H3K4me3 or either modification mean values over all cell lines and binary
162 variables indicating modifications in any cell line were added.
163 *5. Transcription Factor Binding Sites* (12 variables): TFBS ChIP-seq clusters (V3) from ENCODE
164 data retrieving binding sites of top 10 TFs with the highest binding site coverage. Mean value over
165 all TFs and 0-1 indicator of any TF occurrence were added in addition (ENCODE Project
166 Consortium, 2012)
167 *6. Transcription factor binding disruption* (1 variable):
168 P-value of disrupting putative strongest transcription factor binding site due to mutation was
169 calculated with Annotation of Regulatory Variants using Integrated Networks (ARVIN) algorithm
170 (Gao, et al., 2018) using Cis-BP database (Weirauch et al., 2014).
171 *7. Maximum transcription factor binding log-odds ratio score* (1 variable):
172 Maximum TF binding log-odds ratio score for reference and mutated sequences among scores
173 calculated with ARVIN algorithm (Gao, et al., 2018, Weirauch et al., 2014).
174 *8. DNase I hypersensitivity* (1 variable): ENCODE DNase clusters (V3) from 125 cell line types
175 (John, et al., 2011; Thurman, et al., 2012; Rosenbloom, et al., 2013).
176 *9. Evolutionary conservation* (10 variables):
177 <u>a) GERP ++:</u> Genomic Evolutionary Rate Profiling scores (Davydov, et al., 2010).
178 <u>b) PhastCons:</u> PhastCons conservation score by vtools (San Lucas, et al., 2012).
179 <u>c) Z-score:</u> recalculated Z-score values defined in our previous work (Wyrwicz, et al., 2007) on
180 whole genome human–mouse alignments (genome builds hg19 and mm9 (Chiaromonte, et al.,
181 2002; Kent, et al., 2003; Schwartz, et al., 2003) from UCSC Genome Browser (Kent, et al., 2002)
182 for the reference and mutated sequence and for window length 7 and 9. Differences of Z-scores for
183 the reference and mutated sequence were added.
184 *10. Dinucleotide content* (16 variables):

185 Observed vs. expected frequencies of 16 possible pairs of nucleotides appearing in the short
186 sequence motif.
187
188 **RESULTS**
189 <u>Feature importance</u>
190 From MCFS we obtained the ranking of all 227 features according to their relative importance in
191 the classification problem. Each feature group contained multiple individual features with different
192 ranks in the overall ranking. In the context of machine learning task, usefulness of a particular
193 group should be determined by the best performing features from this group.
194 Figure 1 presents detailed feature ranking including all features from each group. Generally,
195 features that contribute to the correct classification mostly belong to GC content group, shape
196 group and sequence group. Other feature groups were of lesser importance (the full ranking is
197 included as Supplementary material 2, feature names glossary as Supplementary material 4). The
198 most important feature was the difference in GC-content between the reference and the mutated
199 sequence fragment (rank 1). Features describing raw nucleotide sequence and dinucleotide content
200 appeared in the middle of the ranking. Among the shape features those describing the closest
201 neighborhood of the mutated nucleotide were the most important. This is not surprising because
202 differences in shape are expected to have local effects on DNA properties. Among the shape
203 features attributes concerning propeller twist were ranked as the most important, attributes
204 concerning helix twist and roll followed, and attributes concerning minor groove width occurred
205 lower in the ranking. What is notable, most of the features appearing among the top 20 concerned
206 differences in shape properties between SNP and wild type. Features derived from transcription
207 factors were less important than sequence-based features. Histone modifications, conservation
208 scores and DNase I hypersensitivity score were not identified as particularly informative features.
209 To investigate the role of individual features we calculated Welch's t-score capturing the
210 relationship between particular feature and class value. Decrease of GC-content between the
211 reference and the mutated sequence correlated negatively with functionality (t-score -8.2088 for
212 decrease for motif length 7, t-score -11.3710 for decrease for motif length 9), while increase of
213 propeller twist value correlated positively (t-score 9.7417 for increase immediately before the
214 modified nucleotide, t-score 5.5047 for increase immediately after the modified nucleotide).
215 The role of each feature in a classification task lies not only in its correlation with class value, but
216 also in how well it complements with other features. For example, Figure 2 presents joint
217 distributions of the two most important features in the two classes (difference of GC-content
218 between the reference and the mutated sequence, difference of propeller twist at $3^{rd}$ position
219 between mutated variant and wild type). For non-functional SNPs the features are uncorrelated, but
220 there is a visible negative correlation for functional SNPs. MCFS allows studying that kind of
221 dependencies through its interdependency discovery function. Full list of feature interdependencies
222 and their relative strength is included as Supplementary material 3. Figure 3 presents graph of the
223 strongest interdependencies among the top selected features (GCSCORE – GC composition, SEQ –
224 sequence feature, ROLL – roll, HELT – helix twist, PROT – propeller twist). Difference in GC-
225 content acts as a central hub and interacts strongly with all groups of shape features except minor
226 groove width. The simplified intuition is that functional SNPs should increase GC content of the
227 motif, and at the same time increase rotation of the DNA strand accordingly.
228
229 <u>Classifier performance</u>
230 Obtained feature ranking suggests that a large portion of information is contained in features
231 derived from the DNA sequence, and features describing evolutionary conservation and functional

232  properties play less significant role. To verify this hypothesis, we performed a cross-validation
233  experiment (with folds determined by chromosomes) on the train set by training gradient tree
234  boosting (GTB) classifier on different combinations of feature groups. Calculated values of
235  multiple performance measures are presented in Table 1.
236  Classifier based on all available features performed better than the classifier using only 25 best
237  ranked features. Among individual feature groups GC content produced classifier with the largest
238  AUC ROC (0.78). Combining GC content with shape features and sequence features allowed
239  achieving AUC ROC 0.98. No other combinations of features performed better. These results show
240  that shape features are more meaningful when combined with another feature. In further
241  experiments classifier trained on sequence, shape and GC content was used. We named this
242  classifier ShapeGTB.
243  We compared final ShapeGTB classifier with more general SNP prioritization methods, which did
244  not focus specifically on promoter regions: CADD, FATHMM-MKL and DeepSEA. Figure 4
245  present precision-recall curves calculated on the hold-out test set constructed from our data
246  (HGMD and 1000 Genomes Project) and for smaller experimental dataset (ClinVar and 1000
247  Genomes Project). Area under precision-recall curve can be interpreted as average precision (AP),
248  and is an aggregated measure of classifier performance. It is preferred over AUC ROC when
249  problem is characterized by large class imbalance. On the hold-out test set ShapeGTB
250  outperformed general-purpose methods by a large margin (AP 0.93 vs. 0.47-0.55). On the external
251  validation set ShapeGTB aggregated performance was comparable with FATHMM-MKL (AP 0.47
252  vs. AP 0.42). However, shapes of precision-recall curves for those methods were very different:
253  FATHMM-MKL displayed high precision only for small subset of examples, while ShapeGTB
254  precision was relatively stable even for large values of recall. Differences between results obtained
255  for the two datasets suggest that ClinVar-derived positives have different characteristics and pose a
256  greater challenge. We speculated that the gap between ShapeGTB and reference tools on the hold-
257  out test is due to inclusion of shape features and their interactions with GC content. To verify this,
258  we randomly permuted these features in our test set and evaluated performance of ShapeGTB again
259  on permuted data sets. AP of ShapeGTB with GC-derived features permuted was 0.80, with shape-
260  derived features permuted 0.44, and with both kinds of features permuted 0.35 (Figure 5). This
261  once more corroborates the hypothesis that shape features together with GC content provide
262  important information for distinguishing functional SNPs in our data set.
263
264  **DISCUSSION AND CONCLUSIONS**
265  Here, we report the influence of the combined multifactor analysis of DNA shape and other
266  descriptors in prediction of functional effect of promoter variants. Previously, Parker et al. has
267  demonstrated that the nucleotide alternations can significantly affect the DNA structure causing
268  changes in protein binding affinity and phenotype (Parker, et al., 2009). From our analysis, it is
269  clear that changes in the geometry of DNA molecule are important features for the task of
270  prioritization of functional regulatory variants within promoter regions. General conclusions that
271  can be drawn from our study are as follows: a) shape features work very locally, what is important
272  is what happens in the closest neighborhood of the mutated nucleotide, b) DNA chain rotations are
273  more important than minor groove width, c) differences of properties of the mutated variant and the
274  reference motif are the most meaningful. This picture is inherently complicated with the presence
275  of feature interdependencies – mostly between GC content and shape features. It is impossible to
276  make predictions based on DNA shape alone, it is meaningful only with respect to the sequence
277  content.

278  Interestingly, in our method the most informative indicator of variant functional impact is whether
279  the introduced nucleotide changes the GC-content. The GC composition has been previously linked
280  to DNA thermostability, bendability and potential for conformational transition between B- and Z-
281  forms, that relate to chromatin accessibility (Vinogradov, 2003). The instances of GC-rich
282  sequence motifs have been shown to play an important role in transcription regulation through their
283  connection with nucleosome occupancy and TF binding (Peckham, et al., 2007; Wang, et al.,
284  2012). In our opinion, high rank of GC-ratio derivatives is a result of promoter properties, which
285  distinguish it from other regulatory elements (Lenhard, et al., 2012; Andersson, et al., 2014). GC-
286  ratio may not be highly ranked if similar analysis would be performed on other regulatory
287  elements, which are not associated with promoter regions (e.g. splicing elements or insulators).
288  There is a vast amount of literature on complex networks of relations between nucleotide types and
289  various shape attributes (Yoon, et al., 1988; Florquin, et al., 2005; Rohs, et al., 2005; Samanta, et
290  al., 2009). For instance, the distribution of water around the minor groove shows specificity to the
291  DNA sequence as the availability of the hydrogen bond forming atoms changes. Variation in DNA
292  sequence may affect DNA flexibility by influencing the magnitude of propeller twist. Specific base
293  pairs combinations have different electrostatic potentials and prefer specific stacking geometry
294  (Samanta, et al., 2009). The results of Tillo and Hughes have highlighted that GC-ratio influences
295  nearly all aspects of DNA structure (Tillo and Hughes, 2009). The most pronounced dependency
296  has been observed between GC-ratio and propeller twist (Ponomarenko, et al., 1999). Deb et al.
297  previously reported the effect of an A/T base pair replacement by a G/C base pair on narrowing of
298  minor grows through negative propeller twisting (Deb, et al., 1987). This pair has also been rated
299  high in our feature interdependencies ranking. To sum up, it appears that only a specific
300  configuration of local structural feature values can meet the requirements of a functional genomic
301  element and that causative mutation substantially disrupt it consensus.
302  The data derived from ChIP-seq experiments and DNaseI hypersensitivity assays have relatively
303  low resolution generally ranging from 200 to 8 kbp (Park, 2009; Pique-Regi, et al., 2011;
304  ENCODE Project Consortium, 2012). Our analysis shows that histone modification and TFBS
305  ChIP-seq peaks along with TF disruption p-value and DNaseI hypersensitivity data, being used in
306  genome-wide setting, have no discriminative power for promoter region sequence variations. This
307  is especially true for TSS-balanced version of our data sets (Supplementary material 5). It is
308  important to stress that features based on histone modifications and TFBS have different meaning
309  than those derived directly from DNA sequence and shape. The former may represent statistical
310  relationships connected with high-level functioning of the organism, while the latter may
311  correspond to low-level binding mechanisms and biophysical properties of the DNA. Our method
312  is able to make successful predictions using only low-level features, which may inform the study of
313  low-level mechanisms behind functional SNP mutations.
314  There is a strong need in the field for entirely independent, high-quality collection of regulatory
315  elements variants categorized by type of non-coding sequence and functional status. Such
316  collection would allow constructing reliable tests sets to validate and compare available methods.
317  According to Li and Wang (2017) analysis, human genetic variants databases such as HGMD and
318  ClinVar contain contradictory entries and incorrectly categorized variants due to the lack of
319  primary review of evidence.
320  In our experiments, our method outperformed significantly the reference tools on our own dataset,
321  and exhibited better recall on external dataset. However, caution is required in drawing final
322  conclusions from the comparison. Our model targeted promoter regions specifically, while the
323  other tools were trained on larger subsets of non-coding regions. It is also possible that our

324  validation set, at least partially, overlapped with training sets used by other algorithms. We believe
325  that the main reason behind good performance of ShapeGTB is the inclusion of shape features.
326  Without them the expected performance is on par with the other methods (AP 0.44 on hold-out test
327  set).
328  In summary, we demonstrated that the local shape features of DNA surrounding single nucleotide
329  coupled with the GC-content and sequence composition are sufficient for single nucleotide variant
330  prioritization within promoter regions of human genes. Our results additionally confirmed the
331  interdependencies between alternations in the GC-content and local DNA shape features. Given
332  that the shape vectors implicitly reflect electrostatics, base stacking, hydration profiles (Przytycka
333  and Levens, 2015), including DNA shape into model results in functional reduction of the number
334  of features and therefore a great simplification of the method. We believe that local DNA shape
335  features carry a vast amount of information and their applicability should be investigated further. In
336  the future, we plan to extend our analysis on all types of regulatory elements in non-coding regions
337  of human genome.
338

339  **Acknowledgements**

342

343  **REFERENCES**

344  Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y, Zhao X,
345  Schmidl C, Suzuki T, Ntini E, Arner E, Valen E, Li K, Schwarzfischer L, Glatz D, Raithel J, Lilje
346  B, Rapin N, Bagger FO, Jorgensen M, Andersen PR, Bertin N, Rackham O, Burroughs AM, Baillie
347  JK, Ishizu Y, Shimizu Y, Furuhata E, Maeda S, Negishi Y, Mungall CJ, Meehan TF, Lassmann T,
348  Itoh M, Kawaji H, Kondo N, Kawai J, Lennartsson A, Daub CO, Heutink P, Hume DA, Jensen TH,
349  Suzuki H, Hayashizaki Y, Muller F, Forrest ARR, Carninci P, Rehli M, and Sandelin A. 2014. An
350  atlas of active enhancers across human cell types and tissues. Nature 507:455-461.
351  10.1038/nature12787
352  Babajide Mustapha I, and Saeed F. 2016. Bioactive Molecule Prediction Using Extreme Gradient
353  Boosting. Molecules 21. 10.3390/molecules21080983
354  Chiaromonte F, Yap VB, and Miller W. 2002. Scoring pairwise genomic sequence alignments. Pac
355  Symp Biocomput:115-126.
356  Chiu TP, Comoglio F, Zhou T, Yang L, Paro R, and Rohs R. 2016. DNAshapeR: an
357  R/Bioconductor package for DNA shape prediction and feature encoding. Bioinformatics 32:1211-
358  1213. 10.1093/bioinformatics/btv735
359  Cock PJ, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T, Kauff F,
360  Wilczynski B, and de Hoon MJ. 2009. Biopython: freely available Python tools for computational
361  molecular biology and bioinformatics. Bioinformatics 25:1422-1423.
362  10.1093/bioinformatics/btp163
363  Consortium EP. 2012. An integrated encyclopedia of DNA elements in the human genome. Nature
364  489:57-74. 10.1038/nature11247
365  Davydov EV, Goode DL, Sirota M, Cooper GM, Sidow A, and Batzoglou S. 2010. Identifying a
366  high fraction of the human genome to be under selective constraint using GERP++. PLoS Comput
367  Biol 6:e1001025. 10.1371/journal.pcbi.1001025
368  Deb S, Tsui S, Koff A, DeLucia AL, Parsons R, and Tegtmeyer P. 1987. The T-antigen-binding
369  domain of the simian virus 40 core origin of replication. J Virol 61:2143-2149.

370  Draminski M, Rada-Iglesias A, Enroth S, Wadelius C, Koronacki J, and Komorowski J. 2008.
371  Monte Carlo feature selection for supervised classification. Bioinformatics 24:110-117.
372  10.1093/bioinformatics/btm486
373  Dror I, Zhou T, Mandel-Gutfreund Y, and Rohs R. 2014. Covariation between homeodomain
374  transcription factors and the shape of their DNA binding sites. Nucleic Acids Res 42:430-441.
375  10.1093/nar/gkt862
376  Florquin K, Saeys Y, Degroeve S, Rouze P, and Van de Peer Y. 2005. Large-scale structural
377  analysis of the core promoter in mammalian and plant genomes. Nucleic Acids Res 33:4255-4264.
378  10.1093/nar/gki737
379  Freeman GS, Lequieu JP, Hinckley DM, Whitmer JK, and de Pablo JJ. 2014. DNA shape
380  dominates sequence affinity in nucleosome formation. Phys Rev Lett 113:168101.
381  10.1103/PhysRevLett.113.168101
382  Friedman JH, and Meulman JJ. 2003. Multiple additive regression trees with application in
383  epidemiology. Stat Med 22:1365-1381. 10.1002/sim.1501
384  Gan Y, Guan J, and Zhou S. 2012. A comparison study on feature selection of DNA structural
385  properties for promoter prediction. BMC Bioinformatics 13:4. 10.1186/1471-2105-13-4
386  Gao L, Uzun Y, Gao P, He B, Ma X, Wang J, Han S, and Tan K. 2018. Identifying noncoding risk
387  variants using disease-relevant gene regulatory networks. Nat Commun 9:702. 10.1038/s41467-
388  018-03133-y
389  Gardiner EJ, Hunter CA, Lu XJ, and Willett P. 2004. A structural similarity analysis of double-
390  helical DNA. J Mol Biol 343:879-889. 10.1016/j.jmb.2004.08.092
391  Genomes Project C, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO,
392  Marchini JL, McCarthy S, McVean GA, and Abecasis GR. 2015. A global reference for human
393  genetic variation. Nature 526:68-74. 10.1038/nature15393
394  Gerstein MB, Kundaje A, Hariharan M, Landt SG, Yan KK, Cheng C, Mu XJ, Khurana E,
395  Rozowsky J, Alexander R, Min R, Alves P, Abyzov A, Addleman N, Bhardwaj N, Boyle AP,
396  Cayting P, Charos A, Chen DZ, Cheng Y, Clarke D, Eastman C, Euskirchen G, Frietze S, Fu Y,
397  Gertz J, Grubert F, Harmanci A, Jain P, Kasowski M, Lacroute P, Leng JJ, Lian J, Monahan H,
398  O'Geen H, Ouyang Z, Partridge EC, Patacsil D, Pauli F, Raha D, Ramirez L, Reddy TE, Reed B,
399  Shi M, Slifer T, Wang J, Wu L, Yang X, Yip KY, Zilberman-Schapira G, Batzoglou S, Sidow A,
400  Farnham PJ, Myers RM, Weissman SM, and Snyder M. 2012. Architecture of the human
401  regulatory network derived from ENCODE data. Nature 489:91-100. 10.1038/nature11245
402  Goni JR, Perez A, Torrents D, and Orozco M. 2007. Determining promoter location based on DNA
403  structure first-principles calculations. Genome Biol 8:R263. 10.1186/gb-2007-8-12-r263
404  Greenbaum JA, Pang B, and Tullius TD. 2007. Construction of a genome-scale structural map at
405  single-nucleotide resolution. Genome Res 17:947-953. 10.1101/gr.6073107
406  Guenther MG, Levine SS, Boyer LA, Jaenisch R, and Young RA. 2007. A chromatin landmark and
407  transcription initiation at most promoters in human cells. Cell 130:77-88.
408  10.1016/j.cell.2007.05.042
409  Heintzman ND, Stuart RK, Hon G, Fu Y, Ching CW, Hawkins RD, Barrera LO, Van Calcar S, Qu
410  C, Ching KA, Wang W, Weng Z, Green RD, Crawford GE, and Ren B. 2007. Distinct and
411  predictive chromatin signatures of transcriptional promoters and enhancers in the human genome.
412  Nat Genet 39:311-318. 10.1038/ng1966
413  Hon GC, Hawkins RD, and Ren B. 2009. Predictive chromatin signatures in the mammalian
414  genome. Hum Mol Genet 18:R195-201. 10.1093/hmg/ddp409
415  Ignatieva EV, Levitsky VG, Yudin NS, Moshkin MP, and Kolchanov NA. 2014. Genetic basis of

416  olfactory cognition: extremely high level of DNA sequence polymorphism in promoter regions of
417  the human olfactory receptor genes revealed using the 1000 Genomes Project dataset. Front
418  Psychol 5:247. 10.3389/fpsyg.2014.00247
419  John S, Sabo PJ, Thurman RE, Sung MH, Biddie SC, Johnson TA, Hager GL, and
420  Stamatoyannopoulos JA. 2011. Chromatin accessibility pre-determines glucocorticoid receptor
421  binding patterns. Nat Genet 43:264-268. 10.1038/ng.759
422  Kanhere A, and Bansal M. 2005. Structural properties of promoters: similarities and differences
423  between prokaryotes and eukaryotes. Nucleic Acids Res 33:3165-3175. 10.1093/nar/gki627
424  Kent WJ, Baertsch R, Hinrichs A, Miller W, and Haussler D. 2003. Evolution's cauldron:
425  duplication, deletion, and rearrangement in the mouse and human genomes. Proc Natl Acad Sci U
426  S A 100:11484-11489. 10.1073/pnas.1932072100
427  Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, and Haussler D. 2002. The
428  human genome browser at UCSC. Genome Res 12:996-1006. 10.1101/gr.229102
429  Kinsella RJ, Kahari A, Haider S, Zamora J, Proctor G, Spudich G, Almeida-King J, Staines D,
430  Derwent P, Kerhornou A, Kersey P, and Flicek P. 2011. Ensembl BioMarts: a hub for data retrieval
431  across taxonomic space. Database (Oxford) 2011:bar030. 10.1093/database/bar030
432  Kircher M, Witten DM, Jain P, O'Roak BJ, Cooper GM, and Shendure J. 2014. A general
433  framework for estimating the relative pathogenicity of human genetic variants. Nat Genet 46:310-
434  315. 10.1038/ng.2892
435  Ladds MA, Thompson AP, Slip DJ, Hocking DP, and Harcourt RG. 2016. Seeing It All:
436  Evaluating Supervised Machine Learning Methods for the Classification of Diverse Otariid
437  Behaviours. PLoS One 11:e0166898. 10.1371/journal.pone.0166898
438  Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, Gu B, Hart J, Hoffman D,
439  Hoover J, Jang W, Katz K, Ovetsky M, Riley G, Sethi A, Tully R, Villamarin-Salomon R,
440  Rubinstein W, and Maglott DR. 2016. ClinVar: public archive of interpretations of clinically
441  relevant variants. Nucleic Acids Res 44:D862-868. 10.1093/nar/gkv1222
442  Landt SG, Marinov GK, Kundaje A, Kheradpour P, Pauli F, Batzoglou S, Bernstein BE, Bickel P,
443  Brown JB, Cayting P, Chen Y, DeSalvo G, Epstein C, Fisher-Aylor KI, Euskirchen G, Gerstein M,
444  Gertz J, Hartemink AJ, Hoffman MM, Iyer VR, Jung YL, Karmakar S, Kellis M, Kharchenko PV,
445  Li Q, Liu T, Liu XS, Ma L, Milosavljevic A, Myers RM, Park PJ, Pazin MJ, Perry MD, Raha D,
446  Reddy TE, Rozowsky J, Shoresh N, Sidow A, Slattery M, Stamatoyannopoulos JA, Tolstorukov
447  MY, White KP, Xi S, Farnham PJ, Lieb JD, Wold BJ, and Snyder M. 2012. ChIP-seq guidelines
448  and practices of the ENCODE and modENCODE consortia. Genome Res 22:1813-1831.
449  10.1101/gr.136184.111
450  Lee D, Gorkin DU, Baker M, Strober BJ, Asoni AL, McCallion AS, and Beer MA. 2015. A
451  method to predict the impact of regulatory variants from DNA sequence. Nat Genet 47:955-961.
452  10.1038/ng.3331
453  Lenhard B, Sandelin A, and Carninci P. 2012. Metazoan promoters: emerging characteristics and
454  insights into transcriptional regulation. Nat Rev Genet 13:233-245. 10.1038/nrg3163
455  Li Q, and Wang K. 2017. InterVar: Clinical Interpretation of Genetic Variants by the 2015 ACMG-
456  AMP Guidelines. Am J Hum Genet 100:267-280. 10.1016/j.ajhg.2017.01.004
457  Lu Q, Hu Y, Sun J, Cheng Y, Cheung KH, and Zhao H. 2015. A statistical framework to predict
458  functional non-coding regions in the human genome through integrated analysis of annotation data.
459  Sci Rep 5:10576. 10.1038/srep10576
460  Morey C, Mookherjee S, Rajasekaran G, and Bansal M. 2011. DNA free energy-based promoter
461  prediction and comparative analysis of Arabidopsis and rice genomes. Plant Physiol 156:1300-

462  1315. 10.1104/pp.110.167809
463  Nishida H, Suzuki T, Kondo S, Miura H, Fujimura Y, and Hayashizaki Y. 2006. Histone H3
464  acetylated at lysine 9 in promoter is associated with low nucleosome density in the vicinity of
465  transcription start site in human cell. Chromosome Res 14:203-211. 10.1007/s10577-006-1036-7
466  Nishizaki SS, and Boyle AP. 2017. Mining the Unknown: Assigning Function to Noncoding Single
467  Nucleotide Polymorphisms. Trends Genet 33:34-45. 10.1016/j.tig.2016.10.008
468  Park PJ. 2009. ChIP-seq: advantages and challenges of a maturing technology. Nat Rev Genet
469  10:669-680. 10.1038/nrg2641
470  Parker SC, Hansen L, Abaan HO, Tullius TD, and Margulies EH. 2009. Local DNA topography
471  correlates with functional noncoding regions of the human genome. Science 324:389-392.
472  10.1126/science.1169050
473  Peckham HE, Thurman RE, Fu Y, Stamatoyannopoulos JA, Noble WS, Struhl K, and Weng Z.
474  2007. Nucleosome positioning signals in genomic DNA. Genome Res 17:1170-1177.
475  10.1101/gr.6101007
476  Pique-Regi R, Degner JF, Pai AA, Gaffney DJ, Gilad Y, and Pritchard JK. 2011. Accurate
477  inference of transcription factor binding from DNA sequence and chromatin accessibility data.
478  Genome Res 21:447-455. 10.1101/gr.112623.110
479  Ponomarenko JV, Ponomarenko MP, Frolov AS, Vorobyev DG, Overton GC, and Kolchanov NA.
480  1999. Conformational and physicochemical DNA features specific for transcription factor binding
481  sites. Bioinformatics 15:654-668.
482  Przytycka TM, and Levens D. 2015. Shapely DNA attracts the right partner. Proc Natl Acad Sci U
483  S A 112:4516-4517. 10.1073/pnas.1503951112
484  Quang D, Chen Y, and Xie X. 2015. DANN: a deep learning approach for annotating the
485  pathogenicity of genetic variants. Bioinformatics 31:761-763. 10.1093/bioinformatics/btu703
486  Rada-Iglesias A, Bajpai R, Swigut T, Brugmann SA, Flynn RA, and Wysocka J. 2011. A unique
487  chromatin signature uncovers early developmental enhancers in humans. Nature 470:279-283.
488  10.1038/nature09692
489  Ram O, Goren A, Amit I, Shoresh N, Yosef N, Ernst J, Kellis M, Gymrek M, Issner R, Coyne M,
490  Durham T, Zhang X, Donaghey J, Epstein CB, Regev A, and Bernstein BE. 2011. Combinatorial
491  patterning of chromatin regulators uncovered by genome-wide location analysis in human cells.
492  Cell 147:1628-1639. 10.1016/j.cell.2011.09.057
493  Ritchie GR, Dunham I, Zeggini E, and Flicek P. 2014. Functional annotation of noncoding
494  sequence variants. Nat Methods 11:294-296. 10.1038/nmeth.2832
495  Rohs R, Sklenar H, and Shakked Z. 2005. Structural and energetic origins of sequence-specific
496  DNA bending: Monte Carlo simulations of papillomavirus E2-DNA binding sites. Structure
497  13:1499-1509. 10.1016/j.str.2005.07.005
498  Rohs R, West SM, Sosinsky A, Liu P, Mann RS, and Honig B. 2009. The role of DNA shape in
499  protein-DNA recognition. Nature 461:1248-1253. 10.1038/nature08473
500  Rosenbloom KR, Armstrong J, Barber GP, Casper J, Clawson H, Diekhans M, Dreszer TR, Fujita
501  PA, Guruvadoo L, Haeussler M, Harte RA, Heitner S, Hickey G, Hinrichs AS, Hubley R,
502  Karolchik D, Learned K, Lee BT, Li CH, Miga KH, Nguyen N, Paten B, Raney BJ, Smit AF, Speir
503  ML, Zweig AS, Haussler D, Kuhn RM, and Kent WJ. 2015. The UCSC Genome Browser
504  database: 2015 update. Nucleic Acids Res 43:D670-681. 10.1093/nar/gku1177
505  Rosenbloom KR, Sloan CA, Malladi VS, Dreszer TR, Learned K, Kirkup VM, Wong MC,
506  Maddren M, Fang R, Heitner SG, Lee BT, Barber GP, Harte RA, Diekhans M, Long JC, Wilder
507  SP, Zweig AS, Karolchik D, Kuhn RM, Haussler D, and Kent WJ. 2013. ENCODE data in the

508  UCSC Genome Browser: year 5 update. Nucleic Acids Res 41:D56-63. 10.1093/nar/gks1172
509  Samanta S, Mukherjee S, Chakrabarti J, and Bhattacharyya D. 2009. Structural properties of
510  polymeric DNA from molecular dynamics simulations. J Chem Phys 130:115103.
511  10.1063/1.3078797
512  San Lucas FA, Wang G, Scheet P, and Peng B. 2012. Integrated annotation and analysis of genetic
513  variants from next-generation sequencing studies with variant tools. Bioinformatics 28:421-422.
514  10.1093/bioinformatics/btr667
515  Schwartz S, Kent WJ, Smit A, Zhang Z, Baertsch R, Hardison RC, Haussler D, and Miller W.
516  2003. Human-mouse alignments with BLASTZ. Genome Res 13:103-107. 10.1101/gr.809403
517  Shen Y, Yue F, McCleary DF, Ye Z, Edsall L, Kuan S, Wagner U, Dixon J, Lee L, Lobanenkov
518  VV, and Ren B. 2012. A map of the cis-regulatory sequences in the mouse genome. Nature
519  488:116-120. 10.1038/nature11243
520  Sheridan RP, Wang WM, Liaw A, Ma J, and Gifford EM. 2016. Extreme Gradient Boosting as a
521  Method for Quantitative Structure-Activity Relationships. J Chem Inf Model 56:2353-2360.
522  10.1021/acs.jcim.6b00591
523  Shihab HA, Rogers MF, Gough J, Mort M, Cooper DN, Day IN, Gaunt TR, and Campbell C. 2015.
524  An integrative approach to predicting the functional effects of non-coding and coding sequence
525  variation. Bioinformatics 31:1536-1543. 10.1093/bioinformatics/btv009
526  Sivolob AV, and Khrapunov SN. 1995. Translational positioning of nucleosomes on DNA: the role
527  of sequence-dependent isotropic DNA bending stiffness. J Mol Biol 247:918-931.
528  10.1006/jmbi.1994.0190
529  Stenson PD, Mort M, Ball EV, Shaw K, Phillips A, and Cooper DN. 2014. The Human Gene
530  Mutation Database: building a comprehensive mutation repository for clinical and molecular
531  genetics, diagnostic testing and personalized genomic medicine. Hum Genet 133:1-9.
532  10.1007/s00439-013-1358-4
533  Thurman RE, Rynes E, Humbert R, Vierstra J, Maurano MT, Haugen E, Sheffield NC, Stergachis
534  AB, Wang H, Vernot B, Garg K, John S, Sandstrom R, Bates D, Boatman L, Canfield TK, Diegel
535  M, Dunn D, Ebersol AK, Frum T, Giste E, Johnson AK, Johnson EM, Kutyavin T, Lajoie B, Lee
536  BK, Lee K, London D, Lotakis D, Neph S, Neri F, Nguyen ED, Qu H, Reynolds AP, Roach V, Safi
537  A, Sanchez ME, Sanyal A, Shafer A, Simon JM, Song L, Vong S, Weaver M, Yan Y, Zhang Z,
538  Zhang Z, Lenhard B, Tewari M, Dorschner MO, Hansen RS, Navas PA, Stamatoyannopoulos G,
539  Iyer VR, Lieb JD, Sunyaev SR, Akey JM, Sabo PJ, Kaul R, Furey TS, Dekker J, Crawford GE, and
540  Stamatoyannopoulos JA. 2012. The accessible chromatin landscape of the human genome. Nature
541  489:75-82. 10.1038/nature11232
542  Tillo D, and Hughes TR. 2009. G+C content dominates intrinsic nucleosome occupancy. BMC
543  Bioinformatics 10:442. 10.1186/1471-2105-10-442
544  Vinogradov AE. 2003. DNA helix: the importance of being GC-rich. Nucleic Acids Res 31:1838-
545  1844.
546  Wang J, Zhuang J, Iyer S, Lin X, Whitfield TW, Greven MC, Pierce BG, Dong X, Kundaje A,
547  Cheng Y, Rando OJ, Birney E, Myers RM, Noble WS, Snyder M, and Weng Z. 2012. Sequence
548  features and chromatin structure around the genomic regions bound by 119 human transcription
549  factors. Genome Res 22:1798-1812. 10.1101/gr.139105.112
550  Weirauch MT, Yang A, Albu M, Cote AG, Montenegro-Montero A, Drewe P, Najafabadi HS,
551  Lambert SA, Mann I, Cook K, Zheng H, Goity A, van Bakel H, Lozano JC, Galli M, Lewsey MG,
552  Huang E, Mukherjee T, Chen X, Reece-Hoyes JS, Govindarajan S, Shaulsky G, Walhout AJM,
553  Bouget FY, Ratsch G, Larrondo LF, Ecker JR, and Hughes TR. 2014. Determination and inference

554  of eukaryotic transcription factor sequence specificity. Cell 158:1431-1443.
555  10.1016/j.cell.2014.08.009
556  Wyrwicz LS, Gaj P, Hoffmann M, Rychlewski L, and Ostrowski J. 2007. A common cis-element
557  in promoters of protein synthesis and cell cycle genes. Acta Biochim Pol 54:89-98.
558  Yoon C, Prive GG, Goodsell DS, and Dickerson RE. 1988. Structure of an alternating-B DNA
559  helix and its relationship to A-tract DNA. Proc Natl Acad Sci U S A 85:6332-6336.
560  Zhou J, and Troyanskaya OG. 2015. Predicting effects of noncoding variants with deep learning-
561  based sequence model. Nat Methods 12:931-934. 10.1038/nmeth.3547
562  Zhou T, Yang L, Lu Y, Dror I, Dantas Machado AC, Ghane T, Di Felice R, and Rohs R. 2013.
563  DNAshape: a method for the high-throughput prediction of DNA structural features on a genomic
564  scale. Nucleic Acids Res 41:W56-62. 10.1093/nar/gkt437

**Table 1**(on next page)

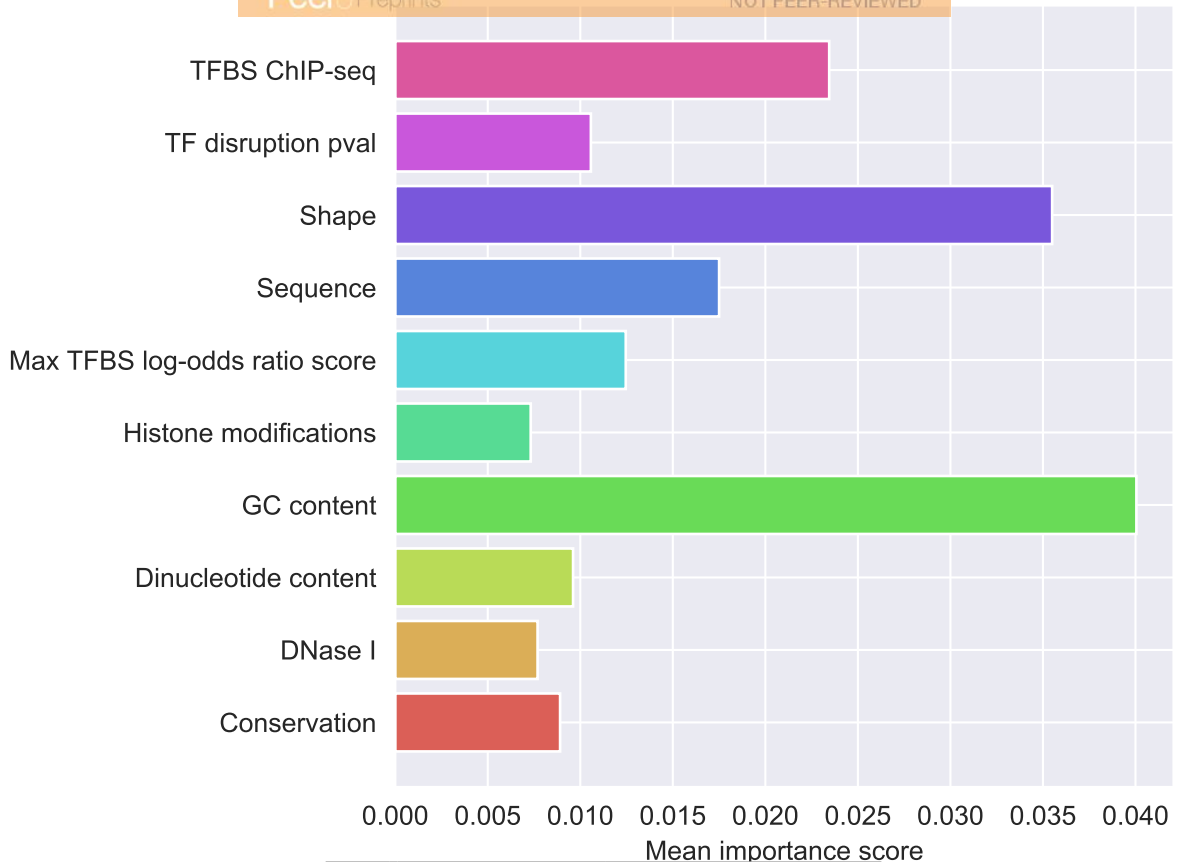Cross-validation classification results for different feature groups on TSS-balanced data set.

1

2

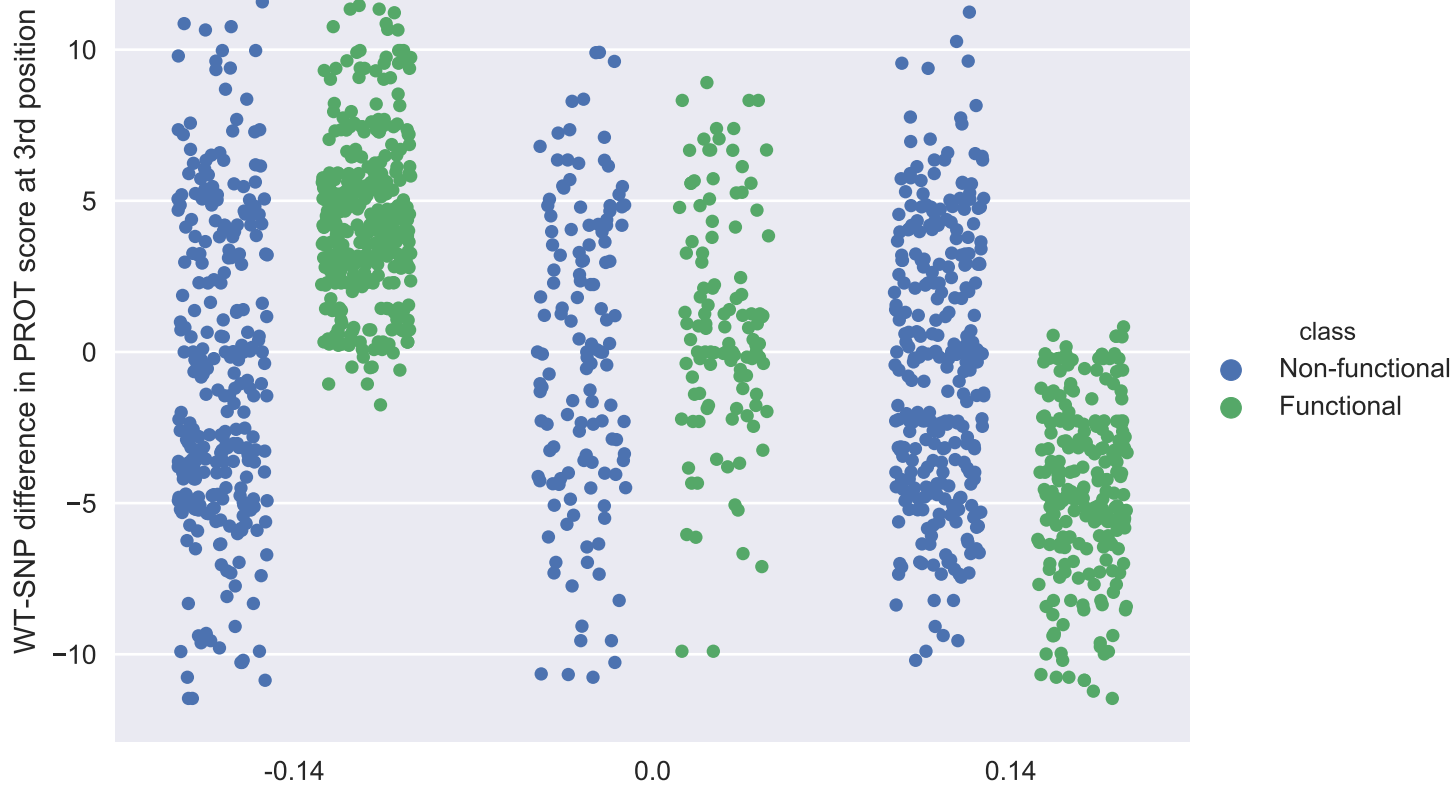| | AUC | AUC_std | Accuracy | Accuracy_std | F1 | F1_std | Precision | Precision_std | Recall | Recall_std | size |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **All** | 0.9764 | 0.0133 | 0.9258 | 0.0247 | 0.8803 | 0.0456 | 0.8840 | 0.0643 | 0.8792 | 0.0480 | 227.0 |
| **Best 25** | 0.9243 | 0.0345 | 0.8449 | 0.0418 | 0.7551 | 0.0785 | 0.7456 | 0.1079 | 0.7713 | 0.0710 | 25.0 |
| **Sequence** | 0.5555 | 0.0473 | 0.6162 | 0.0584 | 0.3170 | 0.0416 | 0.3766 | 0.0878 | 0.2834 | 0.0453 | 52.0 |
| **GC content** | 0.7765 | 0.0525 | 0.7051 | 0.0626 | 0.4934 | 0.0634 | 0.5560 | 0.1054 | 0.4546 | 0.0713 | 8.0 |
| **Shape** | 0.5571 | 0.0566 | 0.6251 | 0.0690 | 0.2546 | 0.0597 | 0.3574 | 0.0994 | 0.2039 | 0.0551 | 88.0 |
| **Conservation** | 0.5440 | 0.0416 | 0.6569 | 0.0522 | 0.2693 | 0.0764 | 0.4313 | 0.1547 | 0.2003 | 0.0545 | 10.0 |
| **TFBS ChIP-seq** | 0.5255 | 0.0482 | 0.6674 | 0.0755 | 0.2416 | 0.0707 | 0.4722 | 0.1589 | 0.1683 | 0.0550 | 12.0 |
| **Histone modifications** | 0.5664 | 0.0641 | 0.6270 | 0.0690 | 0.3342 | 0.0702 | 0.3987 | 0.1069 | 0.2994 | 0.0844 | 38.0 |
| **DNase I** | 0.5846 | 0.0622 | 0.6662 | 0.0817 | 0.1474 | 0.0674 | 0.4088 | 0.1921 | 0.0914 | 0.0431 | 1.0 |
| **Dinucleotide content** | 0.5205 | 0.0615 | 0.6211 | 0.0614 | 0.2354 | 0.0798 | 0.3407 | 0.1323 | 0.1858 | 0.0647 | 16.0 |
| **Max TFBS log-odds ratio score + TF disruption pval** | 0.5141 | 0.0613 | 0.6773 | 0.0824 | 0.0364 | 0.0381 | 0.3812 | 0.3618 | 0.0193 | 0.0205 | 2.0 |
| **Sequence + GC content** | 0.7689 | 0.0404 | 0.6997 | 0.0465 | 0.5029 | 0.0578 | 0.5426 | 0.1159 | 0.4816 | 0.0477 | 60.0 |
| **Shape + GC content** | 0.9175 | 0.0313 | 0.8395 | 0.0333 | 0.7399 | 0.0627 | 0.7557 | 0.1052 | 0.7332 | 0.0583 | 96.0 |
| **Sequence + GC content + Shape** | 0.9787 | 0.0140 | 0.9446 | 0.0208 | 0.9124 | 0.0381 | 0.8894 | 0.0616 | 0.9400 | 0.0437 | 148.0 |
| **Sequence + GC content + Shape + TF disruption pval** | 0.9787 | 0.0132 | 0.9471 | 0.0231 | 0.9161 | 0.0400 | 0.8899 | 0.0624 | 0.9468 | 0.0401 | 149.0 |
| **Sequence + GC content + Shape + TF disruption pval + Max TFBS log-odds ratio score** | 0.9782 | 0.0139 | 0.9442 | 0.0189 | 0.9118 | 0.0318 | 0.8933 | 0.0595 | 0.9346 | 0.0374 | 150.0 |
| **Sequence + GC content + TFBS ChIP-seq** | 0.7902 | 0.0332 | 0.7206 | 0.0410 | 0.5252 | 0.0614 | 0.5698 | 0.0934 | 0.4933 | 0.0616 | 72.0 |
| **Sequence + GC content + Histone modifications** | 0.7981 | 0.0426 | 0.7249 | 0.0464 | 0.5359 | 0.0656 | 0.5882 | 0.1170 | 0.5054 | 0.0664 | 98.0 |

3

**Figure 1**(on next page)

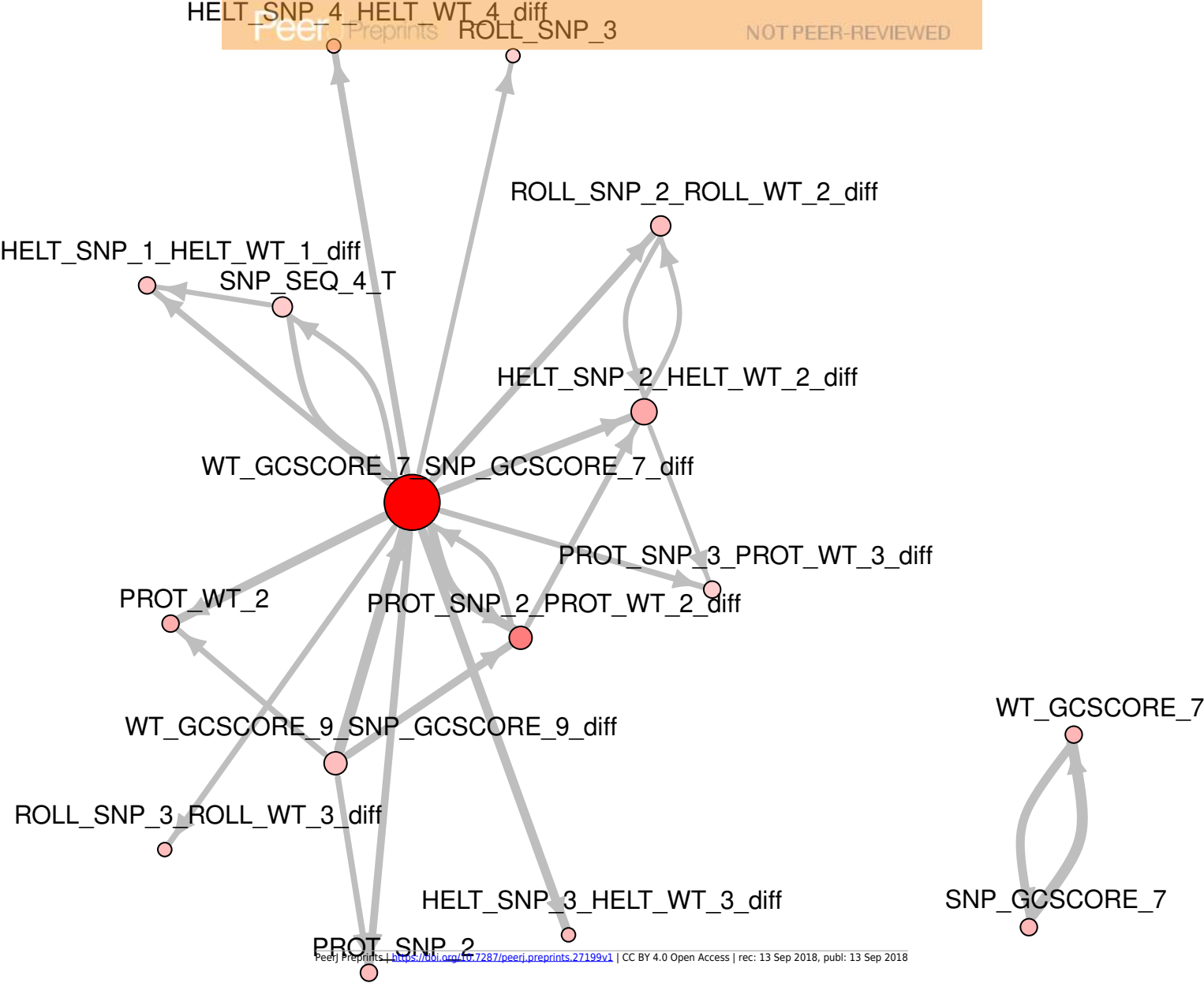Mean importance of 5 best scoring features in each feature group.

**Figure 2**(on next page)

Joint distributions of the two most important features in the two classes. WT-SNP difference corresponds to difference of scores between reference (wild type) and mutated (SNP) variants.

**Figure 3**(on next page)

The strongest feature interdependencies.

HELT_SNP_4_HELT_WT_4_diff

ROLL_SNP_3

ROLL_SNP_2_ROLL_WT_2_diff

HELT_SNP_1_HELT_WT_1_diff

SNP_SEQ_4_T

HELT_SNP_2_HELT_WT_2_diff

WT_GCSCORE_7_SNP_GCSCORE_7_diff

PROT_SNP_3_PROT_WT_3_diff

PROT_WT_2

PROT_SNP_2_PROT_WT_2_diff

WT_GCSCORE_9_SNP_GCSCORE_9_diff

WT_GCSCORE_7

ROLL_SNP_3_ROLL_WT_3_diff
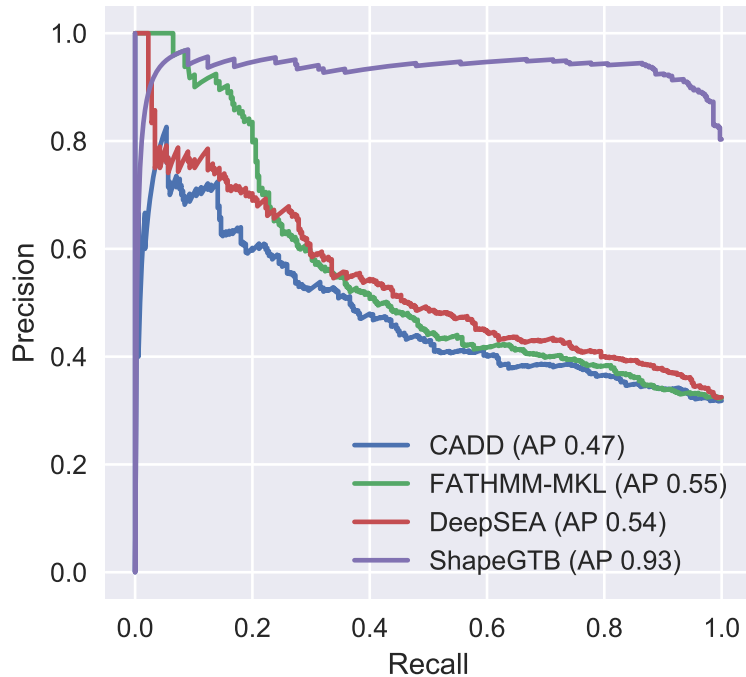
SNP_GCSCORE_7
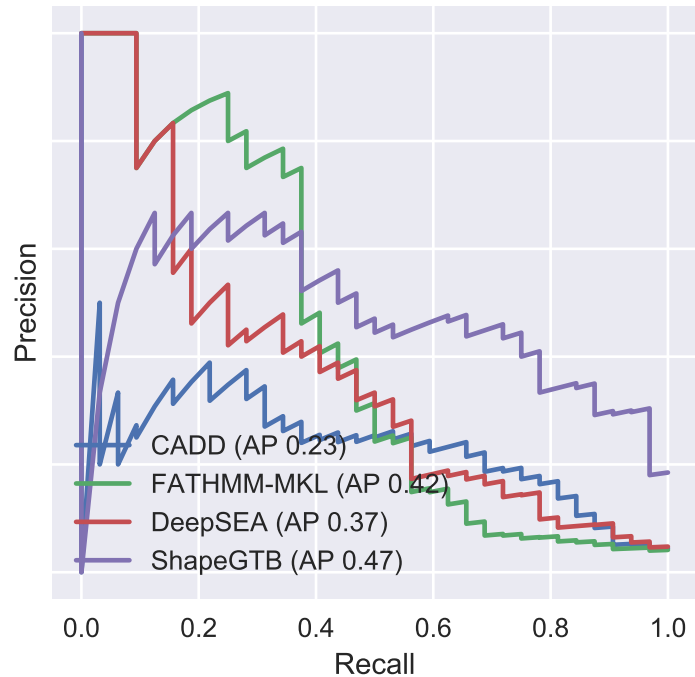
HELT_SNP_3_HELT_WT_3_diff

PROT_SNP_2

**Figure 4**(on next page)

Precision-recall curves for different classifiers.

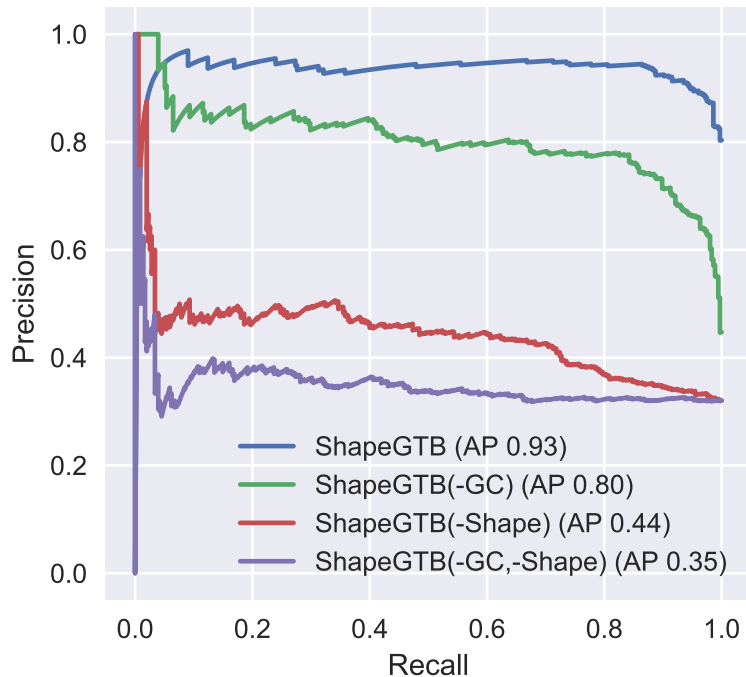a) Hold-out test set (HGMD+1000G)    b) External validation set (ClinVar+1000G)

CADD (AP 0.47)
FATHMM-MKL (AP 0.55)
DeepSEA (AP 0.54)
ShapeGTB (AP 0.93)

CADD (AP 0.23)
FATHMM-MKL (AP 0.42)
DeepSEA (AP 0.37)
ShapeGTB (AP 0.47)

# Figure 5(on next page)

Precision-recall curves for variants of ShapeGTB in which feature vectors from specific feature groups were permuted (effectively reducing their usefulness).

-GC corresponds to classifier with GC-derived features permuted, -Shape corresponds to classifier.

a) Hold-out test set (HGMD+1000G)

ShapeGTB (AP 0.93)
ShapeGTB(-GC) (AP 0.80)
ShapeGTB(-Shape) (AP 0.44)
ShapeGTB(-GC,-Shape) (AP 0.35)

b) External validation set (ClinVar+1000G)

ShapeGTB (AP 0.47)
ShapeGTB(-GC) (AP 0.31)
ShapeGTB(-Shape) (AP 0.07)
ShapeGTB(-GC,-Shape) (AP 0.03)