

The impact of different negative training data on regulatory sequence predictions

Louisa-Marie Krützfeldt^{1,2}, Max Schubach^{1,2}, Martin Kircher^{1,2,*}

¹ Charité – Universitätsmedizin Berlin, 10117 Berlin, Germany

² Berlin Institute of Health (BIH), 10178 Berlin, Germany

* corresponding author

Corresponding author:

Dr. Martin Kircher

Charité – Universitätsmedizin Berlin

BIH JRG Computational Genome Biology

Charitéplatz 1

10117 Berlin

Germany

E-mail: martin.kircher@bihealth.de

Keywords:

regulatory sequence prediction

open chromatin regions

genomic background

sequence shuffles

convolutional neural networks (CNNs)

gapped k-mer SVMs (gkm-SVMs)

32 **Abstract**

33 Regulatory regions, like promoters and enhancers, cover an estimated 5-15% of the human
34 genome. Changes to these sequences are thought to underlie much of human phenotypic
35 variation and a substantial proportion of genetic causes of disease. However, our
36 understanding of their functional encoding in DNA is still very limited. Applying machine or
37 deep learning methods can shed light on this encoding and gapped k-mer support vector
38 machines (gkm-SVMs) or convolutional neural networks (CNNs) are commonly trained on
39 putative regulatory sequences.

40 Here, we investigate the impact of negative sequence selection on model performance. By
41 training gkm-SVM and CNN models on open chromatin data and corresponding negative
42 training dataset, both learners and two approaches for negative training data are compared.
43 Negative sets use either genomic background sequences or sequence shuffles of the positive
44 sequences. Model performance was evaluated on three different tasks: predicting elements
45 active in a cell-type, predicting cell-type specific elements, and predicting elements' relative
46 activity as measured from independent experimental data.

47 Our results indicate strong effects of the negative training data, with genomic backgrounds
48 showing overall best results. Specifically, models trained on highly shuffled sequences
49 perform worse on the complex tasks of tissue-specific activity and quantitative activity
50 prediction, and seem to learn features of artificial sequences rather than regulatory activity.
51 Further, we observe that insufficient matching of genomic background sequences results in
52 model biases. While CNNs achieved and exceeded the performance of gkm-SVMs for larger
53 training datasets, gkm-SVMs gave robust and best results for typical training dataset sizes
54 without the need of hyperparameter optimization.

55

56 Introduction

57 Regulatory sequences play an important role in the control of transcription initiation. Variants
58 in regulatory elements can lead to changes in gene expression patterns and are associated
59 with various diseases [1–3]. Deciphering the encryption of regulatory activity in genomic
60 sequences is an important goal and an improved understanding will inevitably contribute to a
61 better interpretation of personal genomes and phenotypes. While available approaches for
62 measuring changes in regulatory sequences activity in a native genomic context are still very
63 limited in their throughput [4], machine learning methods can be applied for regulatory activity
64 prediction directly from DNA sequence and reveal enriched sequences patterns and
65 arrangements [5].

66 There is a strong link between transcription factors (TFs) binding to regulatory elements and
67 general DNA accessibility, i.e. open chromatin. While the screening of individual TFs is tedious
68 and restricted by the availability of appropriate antibodies, chromatin accessibility can be
69 measured genome-wide and in multiple assays (e.g. DNase-seq, ATAC-seq or NOME-seq).
70 DNase I hypersensitive site sequencing (DNase-seq) provides a gold-standard for the
71 detection of chromatin accessibility [6] and is widely used by the ENCODE Consortium as a
72 sensitive and precise reference measure for mapping regulatory elements [7,8]. It allows the
73 detection of active regulatory elements, marked by DNase I hypersensitive sites (DHS), across
74 the whole genome [9,10].

75 Machine learning approaches identify regulatory elements among other coding or non-coding
76 DNA sequences based on structured patterns of their DNA sequences. Many of these patterns
77 can be matched to known transcription factor binding sites (TFBSs) [11,12] and their relative
78 orientation and positioning. TFs are known to have different binding affinities to DNA
79 sequences and to bind preferentially to a specific set of short nucleotide sequences named
80 binding motifs [11]. Further, TFs can have preferences for a three-dimensional structure of the
81 DNA [12]. While DNA structure can be predicted from the local sequence context, the same
82 DNA shape can be encoded by different nucleotide sequences. There are probably additional
83 patterns, but GC-related sequence features are commonly identified as predictors of
84 regulatory activity and can affect nucleosome occupancy due to differential DNA binding
85 affinity of histone molecules [13].

86 Gapped k-mer support vector machines (gkm-SVMs) [14–16] and convolutional neural
87 networks (CNNs) [17–19] have been recently applied in multiple studies to either predict
88 regulatory activity/function or to identify key elements of the activity-to-sequence encoding.
89 While DHS datasets serve as positive training data for these machine learning algorithms, the
90 ideal composition of the negative training dataset is still an unsolved question. There are two
91 commonly used approaches for the generation of negative training data, the selection of
92 sequences from genomic background [16] and k-mer shuffling of the positive sequences [20–
93 22].

94 In case of genomic background sequences, the negative training dataset is composed of
95 sequences from the genome that are not overlapping DHS regions. However, using non-DHS
96 regions does not guarantee selecting only inactive sequences, due to incomplete sampling of
97 the cell-type under consideration or activity in other cell types. Typically, when selecting
98 background sequences certain properties of the positive training set, e.g. sequence length
99 and repeat fraction, are preserved. Due to this matching of sequence features, this method
100 can be computationally expensive. An alternative approach, k-mer shuffling, is
101 computationally efficient and generates synthetic DNA sequences. A collection of negative
102 sequences according to this approach is composed of the shuffled DHS sequences while
103 preserving each original sequence' k-mer counts.

104 Our work investigates the choice of the negative training dataset and its impact on model
105 performance for predicting regulatory activity from DNA sequences. By applying gkm-SVM
106 and CNN models, both machine learning methods and approaches for negative training data
107 generation are compared. Models are trained on DHS regions from experiments in five
108 different cell lines and various matching negative sets. Performance of the resulting models is
109 evaluated on three different tasks. The first task is the binary classification of DNA sequences
110 into active and inactive for the specific cell line, i.e. classical hold-out performance for
111 individual DHS datasets. The second task tests the ability to learn tissue-specificity and
112 evaluates performance in identifying cell-type specific DHS sequences. In the third task,
113 models are applied to the prediction of enhancer activity and evaluated on an experimental
114 dataset of activity readouts from a reporter assay [23].

115 We show a large impact of the negative training dataset on model performance. Models
116 trained on highly shuffled sequences perform worse except for hold-out performance, while
117 models trained on genomic sequences excel on the more complex tasks of tissue-specific
118 activity prediction and quantitative activity prediction. We speculate that models trained on
119 sequence shuffles learn features of artificial sequence rather than regulatory activity. We also
120 note that insufficient matching of selected genomic background sequences may result in
121 model biases. While CNN performance was improved and exceeded gkm-SVMs for larger
122 training datasets, gkm-SVMs gave better results for small training dataset sizes.

123 **Materials and Methods**

124

125 *Training, validation and test data*

126 In general, positive and negative sequences (except for the independent liver enhancer
127 dataset, see [2.1.4.](#)) were split into three datasets for training, validation, and testing. The
128 validation (hyperparameter optimization) and test sets (performance evaluation) were
129 chromosome hold-out sets of chromosomes 21 and 8, respectively. Training was performed
130 on sequences located on the remaining autosomes and gonosomes.

131

132 *Positive training data: DNase I hypersensitive (DHS) data*

133 DNase-seq datasets were used as positive datasets for regulatory sequence prediction. Seven
134 DNase-seq datasets (narrow peak calls) from experiments in five different cell lines (A549,
135 HeLa-S3, HepG2, K562, MCF-7) were downloaded from ENCODE. Multiple technical
136 replicates were merged into one file per experiment, combining overlapping (minimum of 1 bp)
137 or adjacent sequences into a single spanning sequence. For cell lines A549 and MCF-7 two
138 pooled DHS datasets exist (S1 Table), we refer to those as experiments A and B. DHS regions
139 were defined 300 bp around the center of the narrow peaks and reference genome sequences
140 used (GRCh38 patch release 7, GRCh38.p7). Sequences located on alternative haplotypes,
141 on unlocalized genomic contigs, or containing non-ATCG bases were excluded. An overview
142 of the used DNase-seq datasets is presented in S1 Table.

143

144 *Negative training data: Genomic background data and k-mer shuffling*

145 To obtain genomic background sequences as negative training datasets, DNA sequences with
146 matching repeat and GC content (as in the DHS set) were randomly selected from the
147 genome. While matching repeat content is supposed to correct for potential alignment biases,
148 GC matching is performed to compensate for potential biases caused by better experimental
149 recovery of high GC sequences in DNA handling. Datasets were generated using the
150 genNullSeqs function of the R package gkmSVM [15]. For this purpose, genome sequences
151 (GRCh38.p7) were obtained from UCSC and stored in Biostrings
152 BSgenome.Hsapiens.UCSC hg38.masked

153 (<https://bioconductor.org/packages/release/data/annotation/html/BSgenome.Hsapiens.UCSC>
154 [hg38.masked.html](https://bioconductor.org/packages/release/data/annotation/html/BSgenome.Hsapiens.UCSC), accessed 02/26/2020). To make sure that matching sequences were
155 found for at least 80% of the samples in each dataset, the batch size and maximum number
156 of trials were increased (batchsize=10000, nMaxTrials=100). The tolerance for differences in
157 repeat ratio and relative sequence length were set to 0, but the tolerance for differences in GC
158 content was varied for different training datasets ($t_{GC}=\{0.02, 0.05, 0.1\}$).

159 To generate neutral DNA sequence for the negative training dataset, positive sequences were
160 shuffled while preserving the k-mer counts. Here, k-mer shuffling datasets were generated
161 using `fasta_ushuffle` (https://github.com/agordon/fasta_ushuffle, accessed 02/26/2020), a
162 wrapper for the `fasta` file format to `uShuffle` [24]. The parameter `k` which indicates the size of
163 the preserved k-mers was varied for different datasets ($k=[1,7]$). For each positive sequence,
164 200 shuffled sequences were generated and the sequence with minimal 8-mer overlap to the
165 respective positive sequence chosen.

166

167 *Tissue-specific test data*

168 Assessing the capability of models to predict tissue-specific regulatory activity, datasets with
169 tissue-specific DHS regions were used for testing. For each of the five cell lines, one positive
170 and one negative dataset was generated. For A549 and MCF-7, experiments B were chosen
171 based on best hold-out performance of the gkm-SVM model (shuffled, $k=2$). Positive datasets
172 contain non-overlapping DHS regions to the other four cell lines. The corresponding negative
173 datasets contain DHS regions of the other four cell lines not overlapping with DHS regions of
174 the cell line under consideration. A maximum 30% overlap of regions was tolerated. Tissue-
175 specific datasets were not used for training, but split up in validation and test (i.e. chromosome
176 hold-out sets of chromosomes 21 and 8, respectively; S2 Table) to exclude overlaps with
177 model training.

178

179 *Liver enhancer activity data*

180 Models were tested on an independent dataset of experimental activity readouts [23] to
181 evaluate the models' ability to quantitatively predict enhancer activity. The underlying
182 Massively Parallel Reporter Assay experiments were performed in HepG2 cells infected with
183 lentiviral reporter constructs bearing candidate enhancer sequences chosen on the basis of
184 ENCODE HepG2 chromatin immunoprecipitation sequencing (ChIP-seq) peaks for *EP300*
185 and H3K27ac marks. We used \log_2 RNA/DNA ratios reported for the wild-type integrase
186 experiments and excluded control/synthetic sequences. GRCh37 sequence coordinates were
187 converted to GRCh38.p7 and regulatory sequences where coordinate liftover changed the
188 fragment length were excluded (1 out of 2236). The original fragment size of 171 bp was
189 extended on both ends to a total of 300 bp.

190

191 *Merged datasets of different sizes*

192 A total of six DHS datasets of different sizes from a mixture of the five cell lines were created.
193 100k or 120k DHS regions from each cell line were randomly chosen and resulted in datasets
194 of 500k or 600k DHS regions, respectively. Derived from the 500k dataset, smaller datasets
195 (50k, 100k, 200k and 350k) were randomly sampled.

196

197 *Gapped k-mer support vector machine (gkm-SVM)*

198 Gkm-SVM models were trained with default parameters (word length $l=10$, informative
199 columns $k=6$) and a weighted gkm kernel, as these parameters were previously used for
200 regulatory sequence prediction [16]. To handle big training datasets, the R package LS-GKM
201 [15,25] was used.

202

203 *Convolutional neural network (CNN)*

204 Two different CNN architectures were used. The first architecture, named 4conv2pool4norm
205 (according to 4 convolutional layers, 2 max-pooling layers and 4 normalization layers), was
206 previously presented as DeepEnhancer for accurate prediction of enhancers based on DNA
207 sequence [26]. A smaller network named 2conv2norm (according to 2 convolutional layers
208 and 2 normalization layers), was derived from the 4conv2pool4norm network. Architecture and
209 layer properties of networks are described in S3 and S4 Tables.

210 Models were trained in the Python deep learning library Keras based on the tensorflow
211 interface [27]. The Adam optimizer [28] was used with default parameters as previously
212 suggested [29]. In addition to the default parameters for batch size (200) and learning rate
213 (0.001), a different parameter set was examined (batch size = 2000, learning rate = 0.0002).
214 For both architectures, the higher batch size and lower learning rate were chosen based on
215 accuracy and standard deviation on the validation set (chromosome 21 hold-out, regulatory
216 activity task). Models were trained over 20 epochs showing a convergence of the estimated
217 loss on the validation sets and no signs of overfitting (see S1 Figure and S2 Figure). Network
218 training was repeated 10 times using different seeds. For regulatory activity and tissue-specific
219 activity prediction, one out of the 10 models was chosen for further analysis based on median
220 model performance (chromosome 21 hold-out).

221

222 *Evaluation tasks and model evaluation*

223 Each model was evaluated on three tasks and different performance measures were chosen
224 depending on the task. Receiver Operating Characteristic (ROC) curve and area under ROC
225 curve (AUROC) values are commonly used and a good measure if test datasets are balanced
226 between classes [30] and if the confidence in class labels is similar. An alternative method for
227 imbalanced datasets are Precision-Recall (PR) curves. In contrast to AUROC, area under PR
228 curve (AUPRC) depends on the imbalance of the dataset [31]. A perfect model has an AUPRC
229 value of 1, a random model an AUPRC value equal to the proportion of positive samples in
230 the test set. The R packages PRROC [32,33] and pROC [34] were used to calculate the
231 respective values.

232 For task one (regulatory sequence prediction), AUROC, AUPRC and recall values were used
233 for model evaluation. First models were tested on validation sets to identify best parameters
234 for generating the negative training set based only on recall measures. Based on the test sets,
235 performance of models trained on genomic background or shuffled sequences were compared
236 for each classifier. We evaluated models on their respective hold-out and additionally the
237 models trained on shuffled data on hold-out using genomic background sequences as
238 negative sets. Pairwise comparisons of model performance were realized by Wilcoxon signed-
239 rank tests.

240 The second task considered the models' tissue-specificity. Again, negative training dataset
241 parameters were chosen according to validation dataset performance. Classifiers and types
242 of negative training sets were then compared based on the test datasets. To assess the model
243 performance on task 2 (tissue-specific prediction), PR and ROC curves and corresponding
244 AUPRC and AUROC values were used.

245 For the third task, models were tested on a regression problem and used to predict activity of
246 liver enhancer sequences for which experimental readouts were previously published [23].
247 Here, Spearman rank correlations were calculated between prediction scores and available
248 \log_2 activity ratios.

249

250 *Transcription factor (TF) binding motif analysis*

251 Training dataset sequences were searched for known TF binding profiles and for each dataset
252 the number of matched motifs per 300 bp calculated. A set of 460 non-redundant profiles

253 derived from human TFBSs was exported from the JASPAR CORE database [35]. Profile
254 matches were identified using FIMO [36] with default parameters and a maximum number of
255 motif occurrences retained in memory of 500,000.

256

257 *Frequency distribution of 8-mers*

258 All potential 8-mers consisting only of nucleotides A, C, G and T were extracted from all
259 autosomes (chromosomes 1-22) of the human reference genome sequence (GRCh37) with
260 their absolute count. Obtained 8-mer counts were Z-score transformed, i.e. mean-centered
261 and the standard deviation normalized to 1. Potential 8-mers were further extracted from test
262 and training sequences and the Z-score of their genomic frequency looked up. We also looked
263 up Z-scores for the top 100 scoring 8-mer sequences for each of 128 kernels in the first
264 convolutional layer of the CNN models. In all analyses, 8-mers not observed in the genomic
265 background were excluded from analysis.

266

267 *GC content distribution*

268 The GC content distribution was calculated for active DHS regions in HepG2, three
269 corresponding genomic background datasets with varied GC content tolerance and random
270 genomic sequences. One million random sequences of length 300 bp were selected from
271 GRCh38.p7 (excluding alternative haplotypes and unlocalized contigs) as a reference for the
272 composition of the human genome. For each sequence, GC content was calculated using the
273 R package 'seqinr' [37].

274 **Results**

275

276 *Training models for regulatory activity prediction*

277 To investigate the performance of machine learning methods for regulatory activity prediction
278 from DNA sequence and the impact of negative data set composition, multiple models were
279 compared. Two machine learning approaches, gkm-SVMs and CNNs with two different
280 architectures, were used. The CNN architectures were derived from DeepEnhancer [26] and
281 are referred to as 2conv2norm and 4conv2pool4norm (see Methods). Each model was trained
282 on a positive dataset of DHS regions in a specific cell line (active regulatory sequences) and
283 a corresponding set of negative sequences. Negative training datasets were generated using
284 two different approaches (genomic background, k-mer shuffles) and variation of parameters
285 led to ten different negative training sets per positive dataset. In the genomic background
286 approach three different GC content tolerances ($t_{GC}=\{0.02, 0.05, 0.1\}$) were tested. In the k-
287 mer shuffling approach, the size of the preserved k-mers varied from 1 to 7. The influence of
288 the negative training dataset on model performance was evaluated on chromosome hold-out
289 validation and test sets. First, model hyperparameters were selected on the validation sets,
290 then the models' capability to predict (tissue-specific) regulatory activity was assessed on the
291 test sets, as well as from a quantitative prediction of enhancer activity on an independent
292 experimental dataset.

293

294 *Model performance on chromosome hold-out sets*

295 To measure model performance, we calculated ratios of correctly predicted positive samples,
296 i.e. recall and the area under precision recall curve (AUPRC). For each classifier, we chose
297 one model trained on genomic background and one model using k-mer shuffles for further
298 experiments. To select these models, we compared their performance on a hold-out set of
299 active DHS regions on chromosome 21 (validation set). Since we did not observe relevant
300 effects for parameters of the genomic background set (S3 Figure and S4 Figure), we chose

301 the most stringent parameter ($t_{GC}=0.02$). In contrast, when comparing models trained on
302 shuffled sequences, model performance depended on the size of preserved k-mer k (S5
303 Figure and S6 Figure), with small k resulting in better performance and high k falling behind
304 the genomic background sets. We note that the value of k is anticorrelated to the number of
305 known transcription factor binding site (TFBS) motifs remaining in the negative training
306 sequences (S7 Figure) and suggests that models may identify positive samples based on
307 TFBS frequency. While models with $k=1$ show the best results, we chose $k=2$ as shuffled
308 sequences preserving dinucleotide composition are widely used [20].

309 Selected models were then compared across classifiers on a second chromosome hold-out
310 dataset (chromosome 8, test set). In accordance with previous studies, CNNs and gkm-SVM
311 classifiers are both able to predict active DHS regions from the hold-out sets with high recall
312 and AUPRC values (S8 Figure). We do not see a clear difference between the two CNNs
313 tested. However, models trained on highly shuffled data perform significantly better than
314 models trained on genomic background data; potentially the result of an improper evaluation
315 on varying compositions of the validation sets using different negative data.

316

317 **Fig 1: AUROC values for regulatory sequence prediction.** Models were trained on
318 sequences of DHS regions (positive) with corresponding sets of negative sequences. For each
319 classifier two different negative training sets are compared; sequences were either chosen
320 from genomic background ($t_{GC}=0.02$) or generated by shuffling positive sequences and
321 preserving k-mer counts ($k=2$). Models were tested on a chromosome 8 hold-out test set. The
322 top panels show the results for testing on hold-out sets using genomic background sequences
323 as negative sets, the bottom panels show the results for testing on hold-out sets using shuffled
324 sequences as negative sets. AUROC values were calculated to compare model performance.
325 Seven models were trained on data derived for specific cell lines, bars represent the mean
326 and error bars the standard deviations across models.

327

328 Fig 1 represents AUROC values for all selected models tested on hold-out sets including
329 genomic background sequences (top panels) or shuffled sequences as negative test sets
330 (bottom panels). Differences between CNN and gkm-SVM classifiers are marginal in this
331 comparison and models perform best on the composition that they were trained on. This is in
332 line with models relying on features from both negative and positive sequences. However,
333 models trained on shuffled sequences show a larger drop when tested on a test set using
334 natural sequences as negative class. For example, gkm-SVM models trained on shuffled
335 sequences drop from a mean AUROC of 0.96 to 0.64, while models trained on natural
336 sequences drop from a mean AUROC of 0.90 to 0.83. This suggests that model training may
337 focus more on the shuffled sequences in this case.

338 To explore further, how models were influenced by the negative sets, we analyzed 8-mers in
339 the different test data set classes as well as 8-mers prioritized in the first convolutional layer
340 of our CNN models. We compared these 8-mers based on genomic frequency across all
341 human autosomes. We observe that 8-mers in the genomic background negative sets are on
342 average more frequent than 8-mers from DHS sites (positive sets) and those are more
343 frequent than 8-mers from shuffled negative sequences (S10A Figure). While effects are more
344 subtle, similar effects propagate into 8-mers identified in the first convolutional layers (S10B
345 Figure), with models trained on genomic background sequences learning to identify more
346 common 8-mers (Wilcoxon rank tests, $p < 2.2e^{-16}$). Consequently, rare motifs in shuffled
347 negative sequences are learned by these models and may negatively impact model
348 performance.

349 For A549 and MCF-7 cell lines with two available DHS sets from ENCODE, two separate
350 models were trained and their performance on the test sets compared among all cell lines. We
351 see that performance generalizes well across diverse cell lines (e.g. breast, cervix, lung, liver
352 cancer and leukemia), suggesting that organismal rather than tissue-specific active regulatory
353 regions are predicted. As an example, Table 1 shows recall values for the gkm-SVM models

354 ranging from 0.79 to 0.88 for other cell types. Models trained on A549 training sets perform
355 best on A549 test sets (recall of 0.86 and 0.88, respectively) and MCF-7 models perform best
356 on MCF-7 datasets (recall of 0.90 and 0.91, respectively).

357

358 **Table 1: Recall of test set regulatory sequence prediction for different cell lines.** Gkm-
359 SVM models were trained on DHS datasets (positive) and corresponding sets of k-mer
360 shuffled ($k=2$) sequences (negative) for A549 or MCF-7 cells; cell lines with two training
361 datasets (A/B) each. Models were tested on seven different test sets derived from different
362 cell lines and recall values were calculated to compare model performance. Datasets are
363 named according to S1 Table.

364

		Model			
		A549 (A)	A549 (B)	MCF-7 (A)	MCF-7 (B)
Test set	A549 (A)	0.896	0.863	0.873	0.859
	A549 (B)	0.882	0.880	0.855	0.846
	HeLa-S3	0.877	0.852	0.863	0.848
	HepG2	0.838	0.822	0.813	0.799
	K562	0.843	0.802	0.809	0.793
	MCF-7 (A)	0.872	0.844	0.905	0.893
	MCF-7 (B)	0.870	0.853	0.906	0.900

365

366 *Prediction of tissue-specific regulatory sequences*

367 As seen in the previous experiments, models trained on data derived from one cell line may
368 generalize in predicting active DHS regions in other cell lines. While some regulatory
369 sequences are active in multiple cell types, others are specifically active in only one cell type.
370 To further assess the models' capability to predict tissue-specific regulatory activity, we used
371 datasets containing tissue-specific DHS sequences for further testing. We selected DHS
372 sequences only active in the training cell line (positive samples) and DHS regions not active
373 in this cell line but active in at least one of the other cell lines (negative samples).

374 Again, we first tested parameter choice on a validation set (chromosome 21 hold-out). We
375 notice that performance is considerably reduced compared to the first task and see big
376 differences regarding model performance across different training cell lines (S5 and S6
377 Tables). Since HeLa-S3 models performed best, we focused on this cell line. While models
378 trained using genomic background showed similar performance independent of the GC
379 content tolerance (S11 Figure), performance was dependent on k for shuffled sequences.
380 Model performance tends to increase with higher size of preserved k-mers in shuffled
381 sequences (S12 Figure). For the genomic background set, we chose again the most stringent
382 parameter ($t_{GC}=0.02$) and for shuffled sequences $k=7$ based on precision recall. This high
383 value of k preserves a number of TFBS motifs (46 ± 2 motifs per 300 bp) similar to the positive
384 set (47 ± 2 motifs per 300 bp, S7 Figure), suggesting that presence of tissue-specific factors as
385 well as relative positioning may be most critical for model performance.

386

387

388 **Fig 2: HeLa-S3 model performance for tissue-specific regulatory sequence prediction.**
389 Models were trained on sequences of DHS regions active in HeLa-S3 cells (positive) and
390 negative sequence sets of either matched genomic background sequences ($t_{GC}=0.1$) or k-mer
391 shuffled ($k=7$) sequences. Models were tested on DHS sequences only active in HeLa-S3
392 (positive) and DHS sequences active only in one or multiple other cell lines (A549, HepG2,
393 K562, MCF-7) (negative). Dashed lines represent random model performance. Panels (A) and
394 (B) show ROC and PR curves for 2conv2norm models, (C) and (D) show ROC and PR curves
395 for 4conv2pool4norm models, (E) and (F) show ROC and PR curves for gkm-SVM models.
396 Corresponding AUROC and AUPRC values are provided.

397 We present HeLa-S3 models for the final evaluation on the hold-out test set (chromosome 8).
398 Fig 2 shows ROC and PR curves for 2conv2norm (Fig. 2A/2B), 4conv2pool4norm (Fig. 2C/2D)
399 and gkm-SVM (Fig. 2E/2F) models. Predicting tissue-specific regulatory activity, the
400 performance of models is low, but models trained on genomic background data generally
401 perform better than models trained on shuffled sequences (e.g. AUROC differences of
402 6.7/6.8% for the two different CNN architectures). We do not measure a clear performance
403 difference between the two different CNN architectures, but observe that the gkm-SVM model
404 performed a bit better (AUROC +2%) on this task.

405

406 *Quantitative enhancer activity prediction*

407 Lastly, we evaluated the models capability of predicting quantitative enhancer activity for an
408 independent experimental dataset. For this purpose, we used enhancer activity readouts from
409 published data [23] and calculated Spearman correlation of predicted scores with known
410 activity readouts.

411 Since enhancer activity was measured in HepG2 cells, we first applied our models trained on
412 HepG2 DHS data. In contrast to earlier results, model performance differs across models
413 trained using different GC content matching of the genomic background datasets. Models
414 trained on sequences that varied most from positive sequences regarding their GC content,
415 performed best (S13 Figure). Therefore, this less stringent matching parameter was
416 considered here. Next, the shuffling parameter k was evaluated on enhancer activity prediction
417 for HepG2 models. Here, the extremes, i.e. models trained on highly shuffled sequences ($k=1$)
418 or models with low shuffling ($k=7$) performed worse for the different model types (S14 Figure).
419 Best performance is achieved for $k=\{3,4\}$ for gkm-SVM, while for the CNN architectures
420 $k=\{5,3\}$ perform best. Based on these results, the parameter $k=3$ was chosen.

421

422 **Fig 3: HepG2 and K562 model performance for enhancer activity prediction.** Models were
423 trained either on DHS sequences active in HepG2 or K562 cells (positive) and negative
424 sequences, where sets are either composed of genomic background ($t_{GC}=0.1$) or shuffled
425 ($k=3$) sequences. Models were tested on enhancer sequence activity readouts previously
426 published for HepG2 cells [23]. Spearman rank correlation of predicted scores and \log_2
427 RNA/DNA ratios was used to evaluate model performance. For 2conv2norm and
428 4conv2pool4norm bars represent the median of multiple model training runs ($n=10$) while error
429 bars represent 1st and 3rd quartiles. The dashed black line (Spearman's $\rho=0.276$) represents
430 a reference value which was previously achieved [23].

431 Our HepG2 models did not achieve the performance of a Spearman's ρ of 0.28 reported before
432 [23] (see Fig 3). Therefore, other cell-type models were also tested and A549, HeLa-S3 and
433 K562 models achieved or exceeded the reference performance (Fig 3 incl. HepG2 and K562,
434 further cell-types see S15 Figure). Compared to others, the HepG2 training set is smaller
435 (123k compared to 281k HeLa-S3, 222k K562 and 192k A549, S1 Table). To investigate
436 whether the size of the training dataset influences model performance, new models were
437 trained on datasets of varying size (50k to 600k), by sampling sequences from all cell lines

438 (see Methods). We note that sampling across cell lines dilutes a tissue-specific signal and we
439 expect that correlation with experimental readouts might be reduced.

440 Again, we evaluated the correlation of prediction scores and activity readouts. Results are
441 presented in Fig 4. Model performance of gkm-SVM classifier seems very stable across
442 training set sizes and repeated training runs, but due to runtime we did not test more than
443 350,000 positive training examples. Using genomic background sequences clearly
444 outperformed shuffled sequences. For CNNs, the more complex architecture
445 (4conv2pool4norm) outperformed 2conv2norm on both negative sets. To achieve or exceed
446 the gkm-SVM performance, 4conv2pool4norm required larger training datasets (6-7x more
447 data). Looking across 10 trained CNN models per data set, we see considerable variance in
448 model performance, suggesting high stochasticity in training, likely originating from non-
449 optimal parameters (e.g. batch size, learning rate, convergence). Gkm-SVM (0.29) and
450 4conv2pool4norm models (0.30) both exceeded the reference Spearman's ρ value (0.28, Fig
451 4), despite effects of pooling training datasets across cell lines.
452

453 **Fig 4: Model performance in enhancer activity prediction for different training set sizes.**
454 Models were trained on datasets of different sizes composed of DHS sequences (positive)
455 created by sampling of multiple DHS sets of different cell types, and corresponding negative
456 sequence sets, composed of genomic background ($t_{GC}=0.1$) (on the left) or shuffled ($k=3$)
457 sequences (on the right). Classifiers are represented with different colors. Due to long training
458 durations, gkm-SVM models were trained up to a maximum size of 350k positive samples.
459 Models were tested on enhancer sequences active in HepG2 cells from which activity readouts
460 were previously published [23]. Spearman rank correlation of predicted scores and log2
461 RNA/DNA ratios was used to evaluate model performance. Dots represent median values of
462 repeated model training ($n=10$) while ribbons represent 1st and 3rd quartiles. The dashed black
463 line (Spearman's $\rho=0.276$) represents a reference value achieved previously [23].

464 Discussion

465 We found that CNN models and gkm-SVM models are equally suited for active DHS
466 prediction. While similar in performance, CNN models showed larger variance across training
467 runs and the smaller 2conv2norm network architecture reduced performance on genomic
468 background sets. These and results of k-mer shuffled negative sets suggest that models
469 primarily learn representation differences of short motifs. We note that we selected all shuffles
470 to minimize the 8-mer overlap with the positive sequence template, i.e. sequences that mutate
471 the overall motif positioning. We could also show that k-mer size is correlated to the number
472 of known TFBS motifs found in the negative training sequences and that shuffled sequences
473 have a higher proportion of rare genomic 8-mers than DHS sequences and genomic
474 background sequences. We suggest that learning rare motifs is the reason that model
475 performance for active DHS prediction seems highest when using highly shuffled sequences
476 ($k=\{1..3\}$) as negative training data, but drops considerably when applying models to validation
477 sets using genomic background negative sets. Independent of that effect, genomic
478 background sequences also outperformed shuffles for k higher than 4 for active DHS
479 prediction.

480 Since shuffled sequences are artificial and lack biological constraints, models based on this
481 kind of negative set may learn differential sequence motif representations that correspond to
482 genuine TFBS motifs (both active or inactive in the specific cell-type) and differential motif
483 representation due to other biological constraints (e.g. underrepresentation of CpG
484 dinucleotides). While density of binding sites was previously shown to be predictive of
485 regulatory activity [38,23], quantitative and tissue-specific predictions require the models to
486 learn motifs directly related to sequence activity (e.g. active TFBS in a certain cell-type).
487 Consequently, for the two tasks of tissue-specific activity and quantitative activity prediction,
488 genomic background sequences perform always better than sequence shuffles. In line with

489 these observations, models trained on longer preserved k-mers perform better for these tasks,
490 while still falling behind models using the genomic background. We conclude that with
491 background genomic sequences as negative training data, model training tends to ignore
492 patterns present in natural DNA sequences and is able to focus on more subtle differences in
493 binding site representation.

494 These patterns are consistent across gkm-SVM and CNN models. On the "complex" tasks,
495 gkm-SVM models outperformed the CNN models in our setup. While we do not see a clear
496 difference between CNN architectures for tissue-specific DHS regions, in the quantitative
497 enhancer activity predictions, the more complex 4conv2pool4norm architecture performs
498 considerably better. For biologically meaningful results, appropriate training datasets are
499 always required and we showed on this last task that training set sizes for CNNs need to be
500 much larger to reach gkm-SVM model performance. The amount of training data is also just
501 one parameter that influences CNN model performance and there are many other network
502 and training hyperparameters that can be tuned.

503 The quantitative predictions also revealed an issue with the commonly used software package
504 for drawing background sequences from the genome. While in the first two tests, the GC
505 matching parameter did not seem to make a difference, a larger deviation in GC matching
506 provided a performance increase in quantitative enhancer activity prediction. Concurrently, the
507 HepG2 enhancer activity readouts show a positive correlation of GC content with enhancer
508 activity (Spearman ρ of 0.24 with MaxGC feature in the previous publication, [23]). We
509 therefore looked more rigorously at the GC matching and noticed that even for the most
510 stringent setting, high GC-content DHS regions are not sufficiently matched with genomic
511 background sequences (S16 Figure). This causes the models to learn sequence GC content
512 as predictive of regulatory activity rather than specific sequence patterns. We need to highlight
513 a necessary balance in sequence matching attempts though. While trying to compensate for
514 experimental biases in open chromatin data, we might need to acknowledge a real GC signal
515 due to an enrichment of open chromatin in GC-rich active open chromatin regions, like CpG
516 island promoters [39].

517 **Conclusions**

518 Regulatory sequences are essential for all cellular processes as well as cell-type specific
519 expression in multicellular organisms. A better understanding of the encoding of regulatory
520 activity in DNA sequences is critical and will help to decipher the complex mechanisms of gene
521 expression. Supervised machine learning methods like gkm-SVMs and CNNs can identify
522 associated patterns in DNA sequences [5], however to build the respective models, positive
523 sets of active regulatory sequences and negative sets of inactive sequences are required.
524 While proxies for active regions (e.g. DHS open chromatin sites) are widely available for many
525 cell-types and organisms, negative sets are typically computationally derived from genomic
526 background sequences or shuffles of the positive sequences.

527 To assess whether one approach is preferable over the other, we contrasted both in several
528 experiments. Our results indicate an important influence of negative training data on model
529 performance. Multiple results show that genomic sequences are the better choice for more
530 biologically meaningful results and, when using shuffled sequences, the model performance
531 highly depends on the size of the preserved k-mers.

532 While k-mer shuffling is computationally efficient and generates synthetic DNA sequences,
533 selection of genomic background sequences involves matching of certain properties of the
534 positive training set (e.g. length, GC content, repeat fraction) which makes it computationally
535 more expensive. With the genomic background method applied here [15], we notice that GC
536 matching should be improved to closely reproduce the continuous GC density distribution of
537 the positive set rather than a mean and standard deviation. Further, for both types of negative
538 sets, it is only assumed that sequences are regulatory inactive. For the shuffles this

539 assumption is based on the artificial nature of sequences, for the background it is based on
540 the excluded overlap with active sequences. While this might generally argue for semi-
541 supervised learning approaches, comprehensive positive sets may somewhat alleviate the
542 issue for genomic background sets.

543 Comparing two different machine learning approaches, we show that gkm-SVMs give very
544 robust and good results, while CNNs performance could be improved by larger training
545 datasets. This is inline with gkm-SVMs being the simpler machine learning approach (despite
546 being slower in their current implementation) and we see this as a cautionary reminder to keep
547 models simple, especially if training data is limited. Apart from the negative training data
548 analyzed here, network architecture and training parameters of CNNs should be explored and
549 optimized in future work. The parameter space of CNNs is immense and remains largely
550 underexplored. Further, multi-task CNN implementations show improved performance [18,40],
551 potentially also due to the effective increase in training data. However, to focus our analysis
552 on the effects of the negative set and to keep comparisons to gkm-SVMs possible, we did not
553 include these here.

554 To conclude, this study provided relevant insights about how regulatory activity is encoded in
555 DNA sequence, like highlighting the importance of short sequence motifs, and yielded
556 important insights for training machine learning models. We show that negative training data
557 is of high importance for model performance and that the best results are obtained when using
558 sufficiently large and well-matched genomic background datasets. Comparing different
559 learners, we see that gkm-SVMs are very robust and provide good overall performance. While
560 CNNs have the potential to outperform these simpler models, they require careful attention to
561 the selection of adequate architectures and hyperparameter optimization. While not a focus of
562 this work, models may be further interpreted with respect to their sequence features learned
563 [41,42], in order to shed more light upon the sequence encoding of gene regulation.

564

565 **Acknowledgements**

566 We thank current and previous members of the Kircher group for helpful discussions and
567 suggestions. Specifically, we would also like to acknowledge input from Giorgio Valentini and
568 his lab at Università degli Studi di Milano, as well as Dirk Walther at the University of Potsdam.
569 This work was supported by the Berlin Institute of Health and Charité – Universitätsmedizin
570 Berlin. The funder had no involvement in study design; in the collection, analysis and
571 interpretation of data; in the writing of the report; and in the decision to submit the article for
572 publication.

573

574 **CRedit author statement**

575 **Louisa-Marie Krütfeldt:** Data curation, Formal analysis, Investigation, Methodology,
576 Visualization, Writing-Original draft preparation. **Max Schubach:** Conceptualization, Formal
577 analysis, Investigation, Supervision, Writing- Reviewing and Editing, Project administration.
578 **Martin Kircher:** Conceptualization, Investigation, Supervision, Writing- Reviewing and
579 Editing, Project administration, Funding acquisition.

580

581

582

References

583

- 584 1. Gupta RM, Hadaya J, Trehan A, Zekavat SM, Roselli C, Klarin D, et al. A Genetic
585 Variant Associated with Five Vascular Diseases Is a Distal Regulator of Endothelin-1
586 Gene Expression. *Cell*. 2017;170: 522–533.
- 587 2. Jostins L, Ripke S, Weersma RK, Duerr RH, McGovern DP, Hui KY, et al. Host–
588 microbe interactions have shaped the genetic architecture of inflammatory bowel
589 disease. *Nature*. 2012;491: 119–124.
- 590 3. Vinagre J, Almeida A, Pópulo H, Batista R, Lyra J, Pinto V, et al. Frequency of *TERT*
591 promoter mutations in human cancers. *Nat Commun*. 2013;4: 2185.
- 592 4. Gasperini M, Tome JM, Shendure J. Towards a comprehensive catalogue of validated
593 and target-linked human enhancers. *Nat Rev Genet*. 2020; 1–19.
- 594 5. Ching T, Himmelstein DS, Beaulieu-Jones BK, Kalinin AA, Do BT, Way GP, et al.
595 Opportunities and obstacles for deep learning in biology and medicine. *J R Soc*
596 *Interface*. 2018;15.
- 597 6. Elkon R, Agami R. Characterization of noncoding regulatory DNA in the human
598 genome. *Nat Biotechnol*. 2017;35: 732–746.
- 599 7. ENCODE Project Consortium. A user’s guide to the encyclopedia of DNA elements
600 (ENCODE). *PLoS Biol*. 2011;9: e1001046.
- 601 8. The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the
602 human genome. *Nature*. 2012;489: 57–74.
- 603 9. Liu Y, Fu L, Kaufmann K, Chen D, Chen M. A practical guide for DNase-seq data
604 analysis: from data management to common applications. *Brief Bioinform*. 2018;
605 bby057.
- 606 10. Song L, Crawford GE. DNase-seq: a high-resolution technique for mapping active gene
607 regulatory elements across the genome from mammalian cells. *Cold Spring Harb*
608 *Protoc*. 2010;2010: pdb.prot5384.
- 609 11. Boeva V. Analysis of Genomic Sequence Motifs for Deciphering Transcription Factor
610 Binding and Transcriptional Regulation in Eukaryotic Cells. *Front Genet*. 2016;7: 24.
- 611 12. Samee MdAH, Bruneau BG, Pollard KS. A De Novo Shape Motif Discovery Algorithm
612 Reveals Preferences of Transcription Factors for DNA Shape Beyond Sequence Motifs.
613 *Cell Syst*. 2019;8: 27–42.
- 614 13. Tillo D, Hughes TR. G+C content dominates intrinsic nucleosome occupancy. *BMC*
615 *Bioinformatics*. 2009;10: 442.
- 616 14. Beer MA. Predicting enhancer activity and variant impact using gkm-SVM. *Hum Mutat*.
617 2017;38: 1251–1258.
- 618 15. Ghandi M, Mohammad-Noori M, Ghareghani N, Lee D, Garraway L, Beer MA.
619 gkmSVM: an R package for gapped-kmer SVM. *Bioinformatics*. 2016;32: 2205–2207.
- 620 16. Lee D, Gorkin DU, Baker M, Strober BJ, Asoni AL, McCallion AS, et al. A method to
621 predict the impact of regulatory variants from DNA sequence. *Nat Genet*. 2015;47:
622 955–961. doi:10.1038/ng.3331
- 623 17. Wang M, Tai C, E W, Wei L. DeFine: deep convolutional neural networks accurately
624 quantify intensities of transcription factor-DNA binding and facilitate evaluation of
625 functional non-coding variants. *Nucleic Acids Res*. [cited 9 Apr 2018].
626 doi:10.1093/nar/gky215
- 627 18. Zhou J, Troyanskaya OG. Predicting effects of noncoding variants with deep learning-
628 based sequence model. *Nat Methods*. 2015;12: 931–934. doi:10.1038/nmeth.3547
- 629 19. Zou J, Huss M, Abid A, Mohammadi P, Torkamani A, Telenti A. A primer on deep
630 learning in genomics. *Nat Genet*. 2019;51: 12–18.
- 631 20. Alipanahi B, DeLong A, Weirauch MT, Frey BJ. Predicting the sequence specificities of
632 DNA- and RNA-binding proteins by deep learning. *Nat Biotechnol*. 2015;33: 831–838.
- 633 21. Gesell T, Washietl S. Dinucleotide controlled null models for comparative RNA gene
634 prediction. *BMC Bioinformatics*. 2008;9: 248.
- 635 22. Reid J, Wernisch L. STEME: A robust, accurate motif finder for large data sets. *PLOS*

- 636 ONE. 2014;9: e90735.
- 637 23. Inoue F, Kircher M, Martin B, Cooper GM, Witten DM, McManus MT, et al. A systematic
638 comparison reveals substantial differences in chromosomal versus episomal encoding
639 of enhancer activity. *Genome Res.* 2017;27: 38–52. doi:10.1101/gr.212092.116
- 640 24. Jiang M, Anderson J, Gillespie J, Mayne M. uShuffle: A useful tool for shuffling
641 biological sequences while preserving the k-let counts. *BMC Bioinformatics.* 2008;9:
642 192.
- 643 25. Lee D. LS-GKM: a new gkm-SVM for large-scale datasets. *Bioinformatics.* 2016;32:
644 2196–2198.
- 645 26. Min X, Zeng W, Chen S, Chen N, Chen T, Jiang R. Predicting enhancers with deep
646 convolutional neural networks. *BMC Bioinformatics.* 2017;18: 478. doi:10.1186/s12859-
647 017-1878-3
- 648 27. Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, et al. TensorFlow: large-
649 scale machine learning on heterogeneous distributed systems. *arXiv.* 2016;
650 1603.04467.
- 651 28. Kingma DP, Ba J. Adam: a method for stochastic optimization. *arXiv.* 2014; 1412.6980.
- 652 29. Reddi SJ, Kale S, Kumar S. On the Convergence of Adam and Beyond. *Int Conf Learn*
653 *Represent.* 2018 [cited 26 Apr 2019]. Available:
654 <https://openreview.net/forum?id=ryQu7f-RZ>
- 655 30. Davis J, Goadrich M. The relationship between Precision-Recall and ROC curves. *Proc*
656 *23rd Int Conf Mach Learn - ICML 06.* 2006; 233–240.
- 657 31. Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot
658 when evaluating binary classifiers on imbalanced datasets. *PLOS ONE.* 2015;10:
659 e0118432.
- 660 32. Grau J, Grosse I, Keilwagen J. PRROC: computing and visualizing precision-recall and
661 receiver operating characteristic curves in R. *Bioinformatics.* 2015;31: 2595–2597.
- 662 33. Keilwagen J, Grosse I, Grau J. Area under precision-recall curves for weighted and
663 unweighted data. *PLOS ONE.* 2014;9: e92209.
- 664 34. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez J-C, et al. pROC: an open-
665 source package for R and S+ to analyze and compare ROC curves. *BMC*
666 *Bioinformatics.* 2011;12: 77.
- 667 35. Khan A, Fornes O, Stigliani A, Gheorghe M, Castro-Mondragon JA, van der Lee R, et
668 al. JASPAR 2018: update of the open-access database of transcription factor binding
669 profiles and its web framework. *Nucleic Acids Res.* 2018;46: D260–D266.
- 670 36. Grant CE, Bailey TL, Noble WS. FIMO: scanning for occurrences of a given motif.
671 *Bioinformatics.* 2011;27: 1017–1018.
- 672 37. Charif D, Lobry J. SeqinR 1.0-2: A Contributed Package to the R Project for Statistical
673 Computing Devoted to Biological Sequences Retrieval and Analysis. *Struct Approaches*
674 *Seq Evol Mol Netw Popul Biol Med Phys Biomed Engineering Springer Verl.* 2007;
675 207–232.
- 676 38. Smith RP, Taher L, Patwardhan RP, Kim MJ, Inoue F, Shendure J, et al. Massively
677 parallel decoding of mammalian regulatory sequences supports a flexible
678 organizational model. *Nat Genet.* 2013;45: 1021–1028.
- 679 39. Fenouil R, Cauchy P, Koch F, Descostes N, Cabeza JZ, Innocenti C, et al. CpG islands
680 and GC content dictate nucleosome depletion in a transcription-independent manner at
681 mammalian promoters. *Genome Res.* 2012;22: 2399–2408. doi:10.1101/gr.138776.112
- 682 40. Kelley DR, Snoek J, Rinn JL. Basset: learning the regulatory code of the accessible
683 genome with deep convolutional neural networks. *Genome Res.* 2016;26: 990–999.
684 doi:10.1101/gr.200535.115
- 685 41. Movva R, Greenside P, Marinov GK, Nair S, Shrikumar A, Kundaje A. Deciphering
686 regulatory DNA sequences and noncoding genetic variants using neural network
687 models of massively parallel reporter assays. *PLOS ONE.* 2019;14: e0218073.
688 doi:10.1101/393926
- 689 42. Shrikumar A, Prakash E, Kundaje A. GkmExplain: fast and accurate interpretation of
690 nonlinear gapped k-mer SVMs. *Bioinformatics.* 2019;35: i173–i182.

691 Supporting information

692

693 **S1 Figure: Estimated loss on the training and validation sets over training epochs for**
694 **2conv2norm models.** Each model was trained on a HeLa-S3 DHS (positive) training dataset
695 and a 2-mer shuffled (negative) training dataset using the 2conv2norm classifier. Training was
696 repeated 10 times and results are represented in different shades of blue while the mean
697 values are represented in orange. Estimated loss in the training set and the validation set are
698 displayed on the left and right, respectively.

699

700 **S2 Figure: Estimated loss on the training and validation sets over training epochs for**
701 **4conv2pool4norm models.** Each model was trained on a HeLa-S3 DHS (positive) training
702 dataset and a 2-mer shuffled (negative) training dataset using the 4conv2pool4norm classifier.
703 Training was repeated 10 times and results are represented in different shades of blue while
704 the mean values are represented in orange. Estimated loss in the training set and the
705 validation set are displayed on the left and right, respectively.

706

707 **S3 Figure: Recall values for regulatory sequence prediction on validation sets of**
708 **models trained on genomic background sequences.** Each model was trained on a DHS
709 (positive) training dataset and a genomic background (negative) training dataset and tested
710 on a chromosome 21 hold-out validation set. Recall was calculated as a measure of model
711 performance. For each classifier three different negative training sets are compared where the
712 tolerances of differences in GC content composition (t_{GC}) is varied. Each model was trained
713 on data derived from one cell line. Bars represent the mean of multiple cell lines and technical
714 replicates ($n=7$ for gkm-SVM, $n=70$ for CNNs: 10 replicates per cell line) while error bars
715 represent the standard deviation.

716

717 **S4 Figure: AUPRC values for regulatory sequence prediction on validation sets of**
718 **models trained on genomic background sequences.** Each model was trained on a DHS
719 (positive) training dataset and a genomic background (negative) training dataset and tested
720 on a chromosome 21 hold-out validation set. Area under precision recall curve (AUPRC) was
721 calculated as a measure of model performance. For each classifier three different negative
722 training sets are compared where the tolerances of differences in GC content composition
723 (t_{GC}) is varied. Each model was trained on data derived from one cell line. Bars represent the
724 mean of multiple cell lines and technical replicates ($n=7$ for gkm-SVM, $n=70$ for CNNs: 10
725 replicates per cell line) while error bars represent the standard deviation.

726

727 **S5 Figure: Recall values for regulatory sequence prediction on validation sets of**
728 **models trained on shuffled sequences.** Each model was trained on a DHS (positive)
729 training dataset and a k-mer shuffled (negative) training dataset and tested on a chromosome
730 21 hold-out validation set. Recall was calculated as a measure of model performance. For
731 each classifier seven different negative training sets are compared where the size of preserved
732 k-mers during shuffling is varied. Each model was trained on data derived from one cell line.
733 Bars represent the mean of multiple cell lines and technical replicates ($n=7$ for gkm-SVM, $n=70$
734 for CNNs: 10 replicates per cell line) while error bars represent the standard deviation.

735

736 **S6 Figure: AUPRC values for regulatory sequence prediction on validation sets of**
737 **models trained on shuffled sequences.** Each model was trained on a DHS (positive)
738 training dataset and a k-mer shuffled (negative) training dataset and tested on a chromosome
739 21 hold-out validation set. Area under precision recall curve (AUPRC) was calculated as a
740 measure of model performance. For each classifier seven different negative training sets are
741 compared where the size of preserved k-mers during shuffling is varied. Each model was
742 trained on data derived from one cell line. Bars represent the mean of multiple cell lines and
743 technical replicates ($n=7$ for gkm-SVM, $n=70$ for CNNs: 10 replicates per cell line) while error

744 bars represent the standard deviation.

745

746 **S7 Figure: Number of transcription factor binding motifs in training sequences.**

747 Known human transcription factor binding site (TFBS) motifs were matched in training
748 sequences of different datasets from different cell lines (n=7). Bars represent the mean value,
749 error bars the standard deviation.

750

751 **S8 Figure: Recall values for regulatory sequence prediction.** Models were trained on
752 sequences of DHS regions (positive) with corresponding sets of negative sequences and
753 tested on a chromosome 8 hold-out test set. For each classifier two different negative training
754 sets are compared; sequences were either chosen from genomic background ($t_{GC}=0.02$) or
755 generated by shuffling positive sequences and preserving k-mer counts ($k=2$). Recall was
756 calculated to compare model performance. Seven models were trained on data derived for
757 specific cell lines, bars represent the mean and error bars the standard deviations across
758 models. Pairwise comparisons were performed with Wilcoxon signed-rank tests and asterisks
759 represent significance levels (* $p<0.05$, ** $p<0.01$, *** $p<0.001$).

760

761 **S9 Figure: AUPRC values for regulatory sequence prediction on test sets.** Each model
762 was trained on a DHS (positive) training dataset and a set of neutral sequences (negative)
763 and tested on a chromosome 8 hold-out test set. Recall was calculated as a measure of model
764 performance. For each classifier two different negative training sets are compared. Sequences
765 were either chosen from genomic background ($t_{GC}=0.02$) or generated by shuffling positive
766 sequences and preserving k-mer counts ($k=2$). Each model was trained on data derived from
767 one cell line. Bars represent the mean of multiple cell lines (n=7) while error bars represent
768 standard deviations. Pairwise comparisons were performed with Wilcoxon signed-rank test
769 and asterisks represent significance levels (* $p<0.05$, ** $p<0.01$, *** $p<0.001$).

770

771 **S10 Figure: Genomic frequency of 8-mers in different classes of the test sets and the**
772 **first convolutional layer of the CNN models.** Exemplary for all cell-types, the figure shows
773 results for HeLa-S3. Genomic frequency of 8-mers was extracted across all human autosomes
774 and Z-Score transformed (i.e. mean-centered and standard deviation normalized to one).
775 Eight-mers absent from the genome were discarded in the plots. Panel (A) shows the genomic
776 frequency of 8-mers in the test sets split out as DHS sites (black, positive class), genomic
777 background sequences (red, negative class) and different k-mer shuffles (blue, alternative
778 negative class). Smaller k-mer shuffles contain more rare genomic 8-mers. Panel (B) shows
779 the distribution of the genomic 8-mer frequency for the top 100 sequences for each of 128
780 kernels in the first convolutional layer for 2conv2norm (left) and 4conv2pool4norm (right)
781 architectures.

782

783 **S11 Figure: HeLa-S3 model performance for tissue-specific regulatory sequence**
784 **prediction on validation sets of models trained on genomic background sequences.**

785 Models were trained on DHS sequences (positive) active in HeLa-S3 cells and neutral
786 sequences from genomic background (negative) with varied GC content tolerance (t_{GC}).
787 Models were tested on DHS sequences specifically active in HeLa-S3 (positive) and DHS
788 sequences active only in one or multiple other cell lines (A549, HepG2, K562, MCF-7)
789 (negative). (A) and (B) show ROC and PR curves for 2conv2norm models, (C) and (D) show
790 ROC and PR curves for 4conv2pool4norm models, (E) and (F) show ROC and PR curves for
791 gkm-SVM models. Corresponding AUROC and AUPRC values are included.

792

793 **S12 Figure: HeLa-S3 model performance for tissue-specific regulatory sequence**
794 **prediction on validation sets of models trained on shuffled sequences.**

795 Models were trained on DHS sequences (positive) active in HeLa-S3 cells and neutral sequences from
796 genomic background (negative) with varied size of preserved k-mers. Models were tested on
797 DHS sequences specifically active in HeLa-S3 (positive) and DHS sequences active only in
798 one or multiple other cell lines (A549, HepG2, K562, MCF-7) (negative). (A) and (B) show

799 ROC and PR curves for 2conv2norm models, (C) and (D) show ROC and PR curves for
800 4conv2pool4norm models, (E) and (F) show ROC and PR curves for gkm-SVM models.
801 Corresponding AUROC and AUPRC values are included.

802

803 **S13 Figure: HepG2 model performance for enhancer activity prediction of models**
804 **trained on genomic background sequences.** Models were trained on HepG2 DHS
805 sequences (positive) and genomic background sequences (negative), where different
806 genomic background sets result from a variation of the GC content tolerance (t_{GC}). Models
807 were tested on enhancer activity readouts in HepG2 cells [23]. Spearman rank correlation of
808 predicted scores and log₂ RNA/DNA ratios was used to evaluate model performance.

809

810 **S14 Figure: HepG2 model performance for enhancer activity prediction of models**
811 **trained on shuffled sequences.** Models were trained on HepG2 DHS sequences (positive)
812 and genomic background sequences (negative), where different genomic background sets
813 result from a variation of the size of preserved k-mers. Models were tested on enhancer activity
814 readouts in HepG2 cells [23]. Spearman rank correlation of predicted scores and log₂
815 RNA/DNA ratios was used to evaluate model performance.

816

817 **S15 Figure: Model performance for enhancer activity prediction of A549, HeLa-S3 and**
818 **MCF-7 models.** Models were trained either on DHS sequences active in A549, HeLa-S3 or
819 MCF-7 cells (positive) and neutral sequences (negative), where different negative sets are
820 composed of genomic background ($t_{GC}=0.1$) or shuffled ($k=3$) sequences. Models were tested
821 on activity readouts of enhancer sequences in HepG2 cells [23]. Spearman rank correlation
822 of predicted scores and log₂ RNA/DNA ratios was used to evaluate model performance. For
823 2conv2norm and 4conv2pool4norm bars represent the median of multiple replicates (n=10)
824 while error bars represent 1st and 3rd quartiles. The dashed black line represents a reference
825 value (Spearman's $\rho=0.276$) achieved previously [23].

826

827 **S16 Figure: Distribution of GC content in sequences of HepG2 training datasets.** The
828 distribution of the sequences' GC contents in a dataset of active DHS regions in HepG2, three
829 corresponding genomic background datasets with varied GC content tolerance (t_{GC}) and a set
830 of random 300 bp sequences from the genome is shown.

831

832 **S1 Table: Overview of DNase-seq datasets.** The number of DHS sequences is given after
833 merging replicates and exclusion of alternative haplotypes, unlocalized genomic contigs and
834 sequences containing non-ATCG bases. The datasets were split up into training, validation
835 (chromosome 21) and test (chromosome 8) sets. The number of samples in these sets are
836 given in the respective columns. Experiment and Replicate IDs are referring to ENCODE
837 accessions[8].

838

839 **S2 Table: Overview of tissue-specific validation and test sets.** Tissue-specific positive
840 samples are DHS sequences of one cell line not overlapping with DHS sequences of the other
841 cell lines. In contrast, negative samples are DHS sequences of other cell lines not overlapping
842 with the first cell line. For A549, one dataset was chosen (B, named according to S1 Table).
843 For MCF-7 one dataset was chosen (B, named according to S1 Table). The number of DHS
844 sequences is given after exclusion of alternative haplotypes, unlocalized genomic contigs and
845 sequences containing non-ATCG bases. The validation and test sets contain sequences
846 located on chromosome 21 and 8, respectively.

847

848 **S3 Table: Layer properties of 4conv2pool4norm network.** The column named 'Size'
849 provides the convolutional kernel size, the max-pooling window size, the relative dropout size
850 and the dense layer size depending on information given in column 'Layer type'.

851

852 **S4 Table: Layer properties of 2conv2norm network.** The column named 'Size' provides the
853 convolutional kernel size, the max-pooling window size, the relative dropout size and the

854 dense layer size depending on information given in column 'Layer type'.

855

856 **S5 Table: AUROC values for tissue-specific regulatory sequence prediction on**
857 **validation sets.** Models were trained on DHS sequences (positive) with corresponding sets
858 of negative sequences and tested on a set of tissue-specific chromosome 21 test set. For
859 each classifier two different negative training sets are compared; sequences were either
860 chosen from genomic background ($t_{GC}=0.1$) or generated by shuffling positive sequences and
861 preserving k-mer counts ($k=7$). AUROC value was calculated to compare model performance
862

863 **S6 Table: AUPRC values for tissue-specific regulatory sequence prediction on**
864 **validation sets.** Models were trained on DHS sequences (positive) with corresponding sets
865 of negative sequences and tested on a set of tissue-specific chromosome 21 test set. For
866 each classifier two different negative training sets are compared; sequences were either
867 chosen from genomic background ($t_{GC}=0.1$) or generated by shuffling positive sequences and
868 preserving k-mer counts ($k=7$). AUPRC value was calculated to compare model performance.
869

DHS prediction on test set

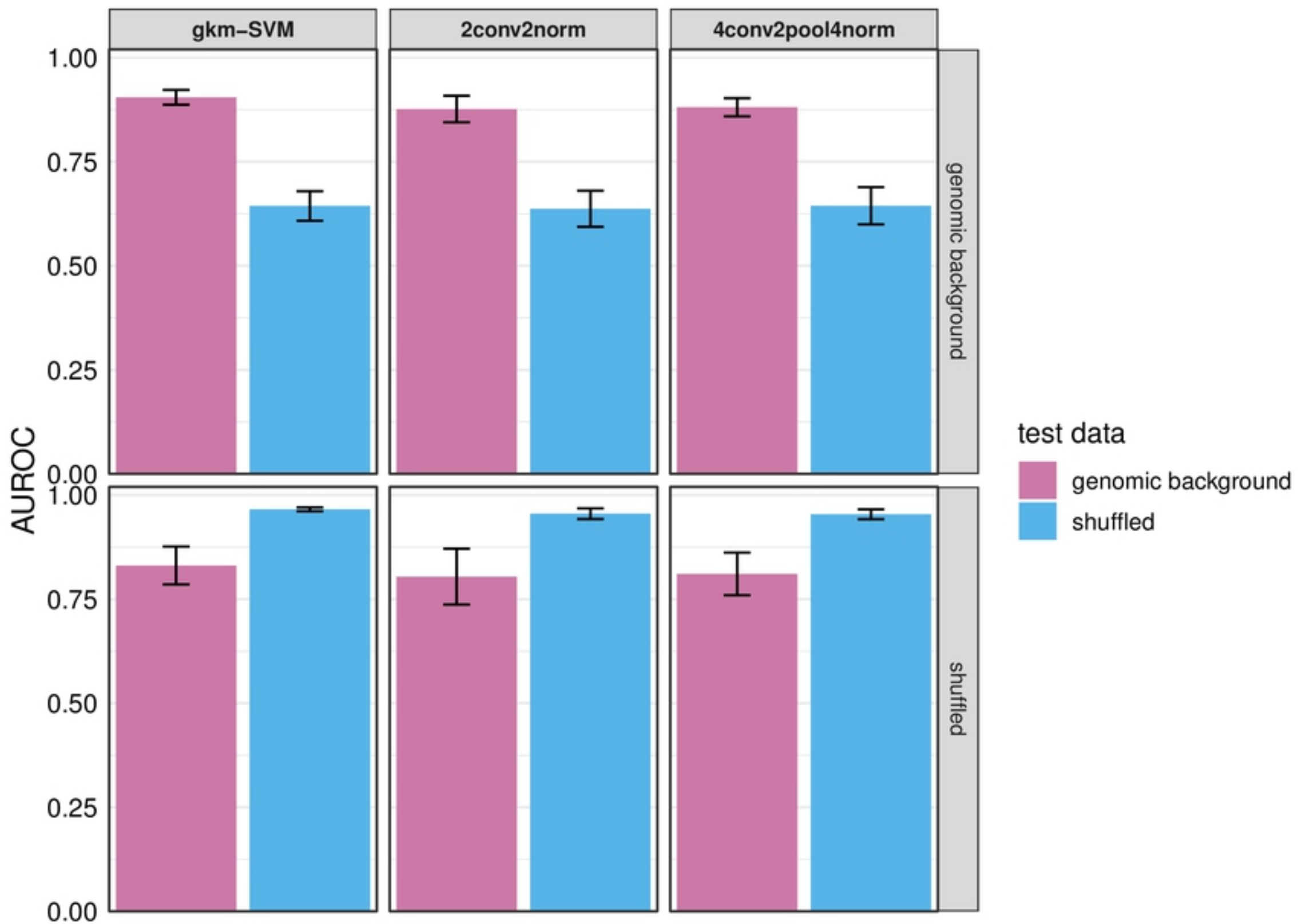


Fig 1

Tissue-specific DHS prediction on test set

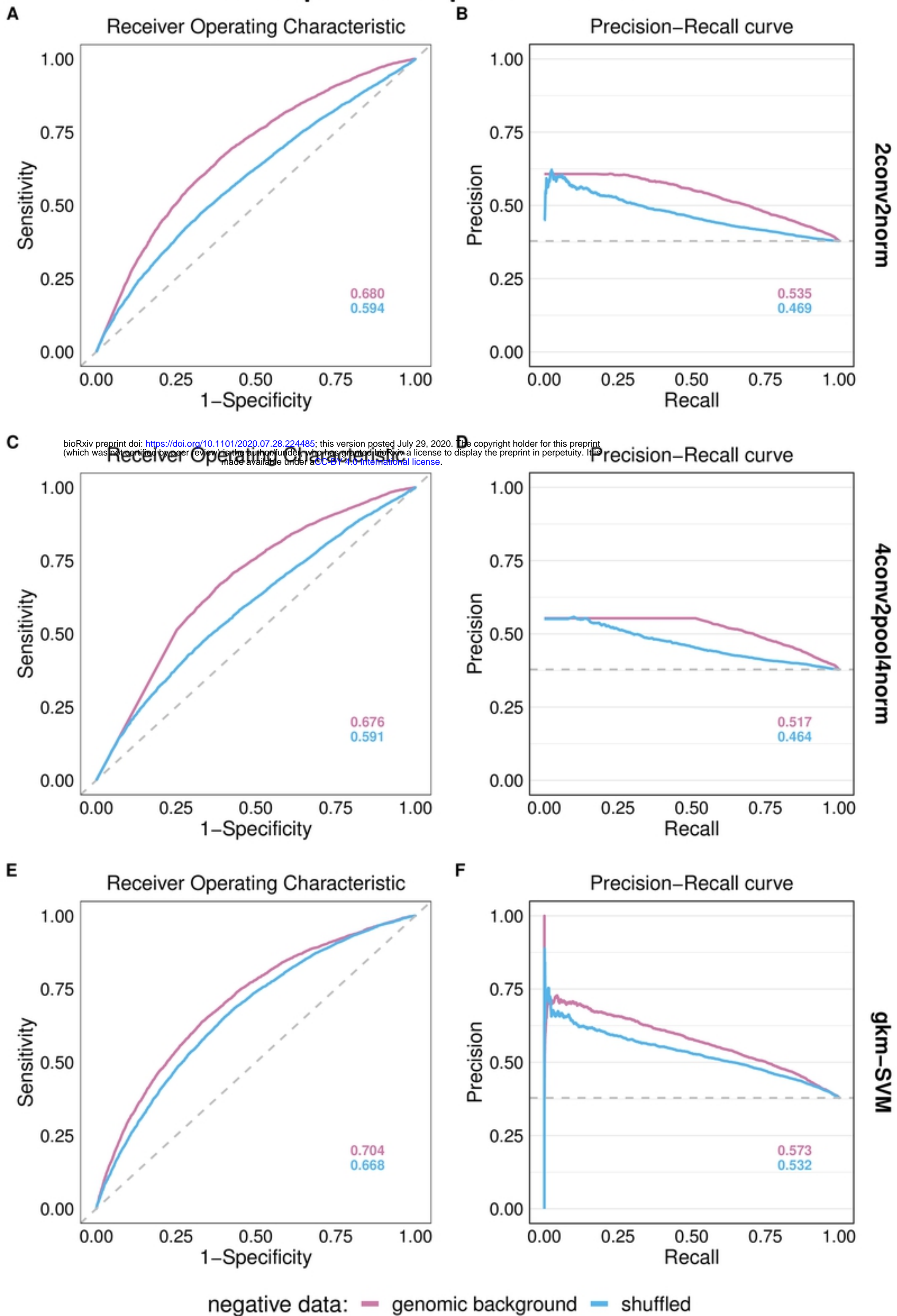


Fig 2

Liver enhancer activity prediction

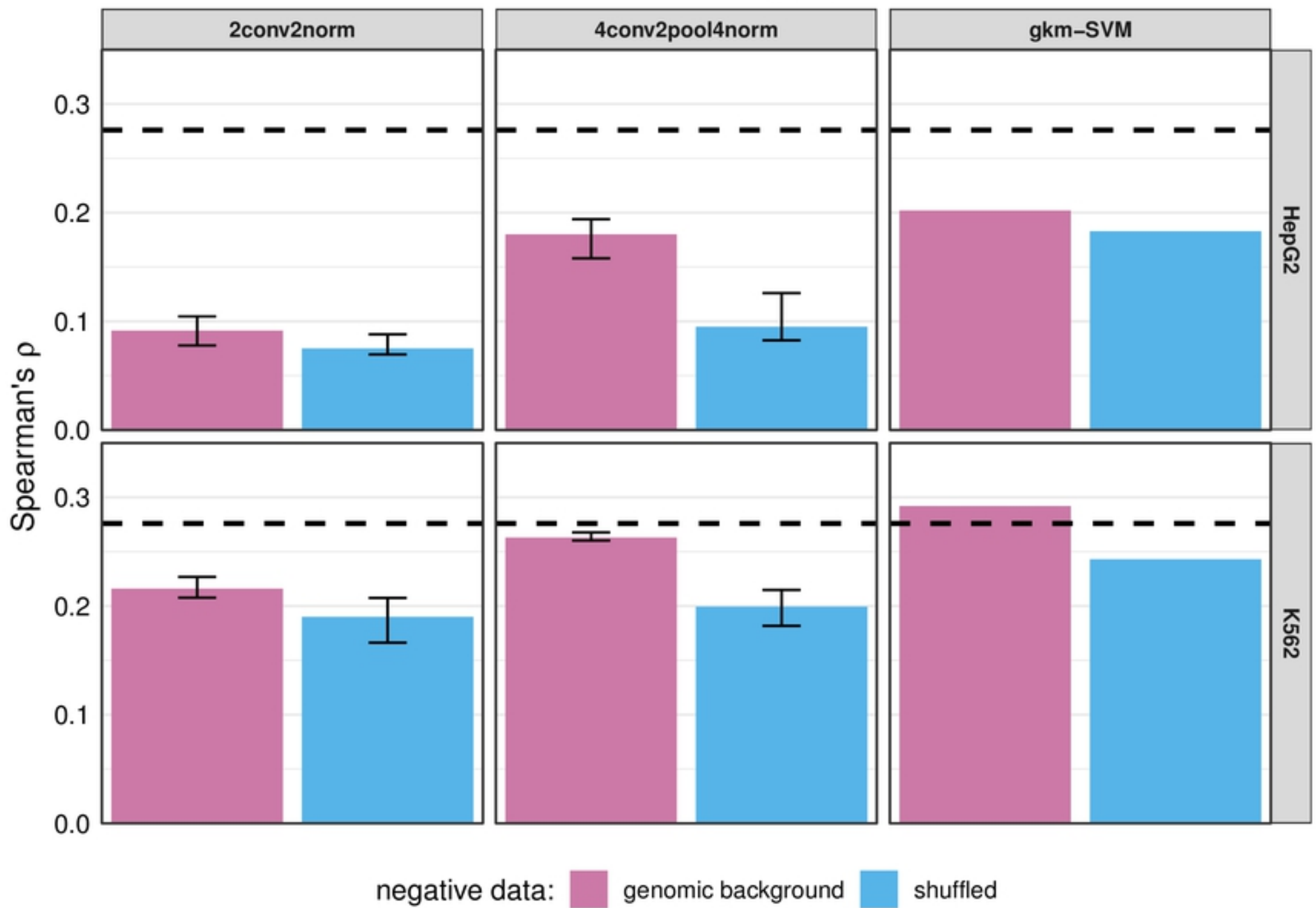


Fig 3

Liver enhancer activity prediction

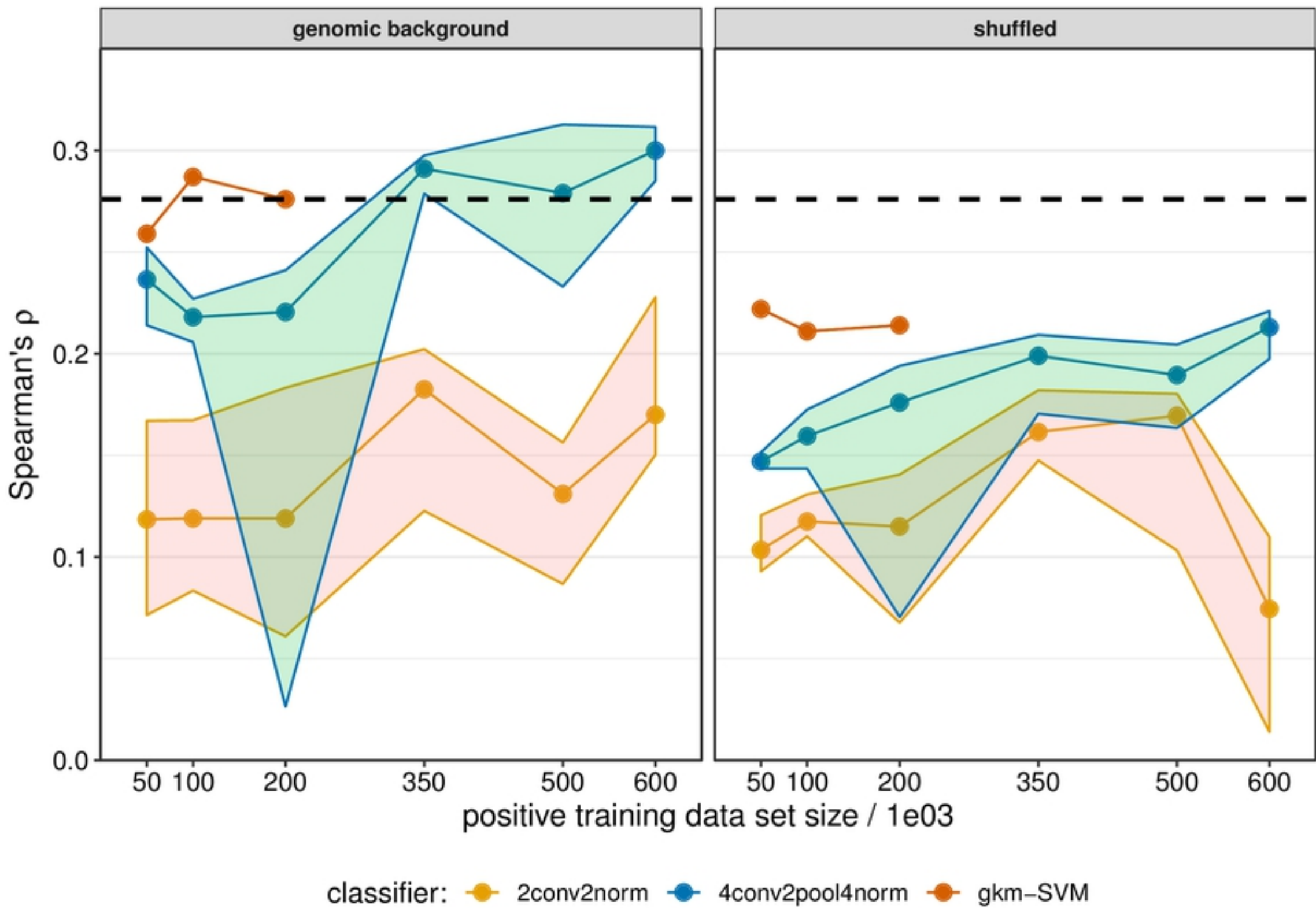


Fig 4