

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

GERV: A Statistical Method for Generative Evaluation of Regulatory Variants for Transcription Factor Binding

Haoyang Zeng¹, Tatsunori Hashimoto¹, Daniel D Kang¹, David K Gifford^{1,2,*}

¹Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02142, USA

²Department of Stem Cell and Regenerative Biology, Harvard University and Harvard Medical School, Cambridge, MA 02138, USA

* Email: gifford@mit.edu

Abstract

The majority of disease-associated variants identified in genome-wide association studies (GWAS) reside in noncoding regions of the genome with regulatory roles. Thus being able to interpret the functional consequence of a variant is essential for identifying causal variants in the analysis of GWAS studies. We present GERV (Generative Evaluation of Regulatory Variants), a novel computational method for predicting regulatory variants that affect transcription factor binding. GERV learns a k-mer based generative model of transcription factor binding from ChIP-seq and DNase-seq data, and scores variants by computing the change of predicted ChIP-seq reads between the reference and alternate allele. The k-mers learned by GERV capture more sequence determinants of transcription factor binding than a motif-based approach alone, including both a transcription factor's canonical motif as well as associated co-factor motifs. We show that GERV outperforms existing methods in predicting SNPs associated with allele-specific binding. GERV correctly predicts a validated causal variant among linked SNPs, and prioritizes the variants previously reported to modulate the binding of FOXA1 in breast cancer cell lines. Thus, GERV provides a powerful approach for functionally annotating and prioritizing causal variants for experimental follow-up analysis. The implementation of GERV and related data are available at <http://gerv.csail.mit.edu/>

1 Introduction

Genome-Wide Association Studies (GWAS) have revealed genetic polymorphisms that are strongly associated with complex traits and diseases ([21, 20, 27, 12]). Missense and nonsense variants that occur in protein coding sequences are simple to characterize. However, many GWAS detected variants reside in non-coding regions with regulatory function ([12, 7]). The influence of non-coding variation on gene expression and other cellular functions is not well understood. Previous work has observed that non-coding DNA changes in the recognition sequences of transcription factors can affect gene expression and cellular phenotypes ([31]). Thus predicting the effect of genomic variants on TF binding is an essential part of interpreting the role of non-coding variants in pathogenesis. Most of existing computational approaches to predict the effect of SNPs on TF binding such as sTRAP and HaplogReg are based on quantifying the difference between the presented reference and alternate alleles in the context of canonical TF binding motifs ([1, 17, 19, 24, 30, 22, 28]). Recent work ([14]) uses k-mer weights learned from a gapped-kmer SVM ([9]) to score the effect of variants, taking into account the frequency of k-mer occurrences but not the spatial effect of k-mers.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Here we present GERV (Generative Evaluation of Regulatory Variants), a novel computational model that learns the spatial effect of k-mers on TF binding *de novo* from whole-genome ChIP-seq and DNase-seq data, and scores variants by the change in predicted ChIP-seq read counts between the reference and alternate alleles. GERV improves on existing models in three ways. First, GERV doesn't assume the existence of a canonical TF binding motif. Instead it models transcription factor binding by learning the effects of specific k-mers on observed binding. This allows GERV to capture more subtle sequence features underlying transcription factor binding including non-canonical motifs. Second, GERV accounts for the spatial effect of k-mers and learns the effect of cis-regulatory regions outside of the canonical TF motif. This enables us to model the role of important auxiliary sequences in transcription factor binding, such as cofactors. Third, GERV incorporates chromatin openness information as a covariate in the model which boosts the accuracy of the predicted functional consequence of a variant.

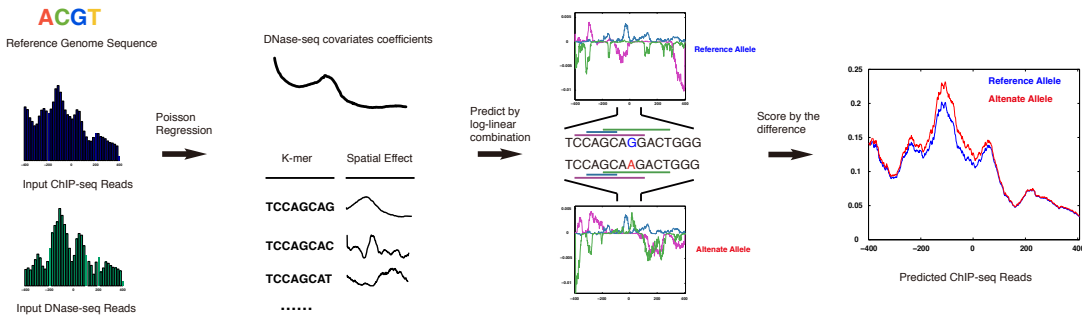


Figure 1: The schematic of GERV. The spatial effects of all the k-mers and the coefficients for DNase-seq covariates are learned from the reference genome sequence and ChIP-seq, DNase-seq datasets. Then the spatial effects (purple, cyan and green) of the k-mers underlying the reference (blue) and alternate (red) allele for a variant are aggregated with DNase-seq covariates by log-linear combination to yield a spatial prediction of local ChIP-seq reads for the two alleles. GERV scores the variant by the l^2 -norm of the predicted change of reads.

We first demonstrate the power of GERV on the ChIP-seq data for transcription factor NF- κ B. We show that GERV learns a vocabulary of k-mers that accurately predicts held-out NF- κ B ChIP-seq data and captures the canonical NF- κ B motifs as well as associated sequences such as known cofactors. Applying GERV to six transcription factors on which allele-specific binding (ASB) analysis is available, we show GERV outperforms existing approaches in prioritizing SNPs associated with allele-specific binding. We demonstrate the application of GERV in post-GWAS analysis by scoring risk-associated SNPs and their linked SNPs for breast cancer, and show that GERV trained on FOXA1 ChIP-seq data achieves superior performance in prioritizing SNPs previously reported to modulate FOXA1 binding in breast cancer cell lines.

2 Methods

2.1 GERV Model Overview

GERV is a fully generative model of ChIP-seq reads. We assume that the genome is a long regulatory sequence containing k-mer “code words” that induce invariant spatial effects on proximal transcription factor binding. We use the level of chromatin openness in a region as a functional prior to predict the magnitude of a sequence-induced binding signal. Following this assumption, we model the read counts produced by transcription factor ChIP-seq at a given base as the log-linear combination of the DNase-seq signal on nearby bases and the spatial effect of a set of learned k-mers whose effect range covers that base.

The GERV procedure of variant scoring consists of the following three steps: (Figure 1)

1. GERV first learns the spatial effect of all k-mers ($k=1$ to 8) and the covariate coefficient of local DNase-seq signal over a spatial window of ± 200 base pairs (bp) *de novo* from ChIP-seq data using regularized Poisson regression
2. GERV then computes the predicted ChIP-seq read counts for the reference and alternate allele of a variant from the log-linear combination of the local DNase-seq signal and spatial effect of the learned k-mers.
3. GERV predicts the effect of a genomic variant on transcription factor binding by the l^2 -norm of the change of predicted reads between two alleles

2.2 Learning the Spatial Effect of K-mers

The effect profile of a k-mer is defined as a real-valued vector of length $2M$ that corresponds to a spatial window of $[-M, M - 1]$ relative to the start position of the k-mer. Specifically, the j -th entry of the profile for a k-mer is the expected log-change in read counts at the j -th base relative to the start of the k-mer. Here we consider k-mers with k from 1 to 8 ($k_{max} = 8$) as this is the maximum that would fit in memory in an Amazon EC2 c3.8 xlarge instance. Larger K-mers tested on a larger memory machine did not perform substantially better than 8-mers. As ChIP-seq signals are relatively sparse and spikey, we chose an effect range of ± 200 bp for each k-mer ($M = 200$).

For notational convenience we use i for genomic coordinate, k for k-mer length, and j for coordinate offset from the start of a k-mer. We assume that the genome consists of one large chromosome with coordinate 0 to N . In practice we will construct this by concatenating chromosomes with the telomeres acting as a spacer. We represent the effect vector of all k-mers of length k as a parameter matrix θ^k of size $4^k \times 2M$. For any particular k-mer of length k starting at base i on the reference genome, we define g_i^k as its row index in θ^k . So $\theta_{(g_i^k, j)}^k$ would denote the effect of this kmer at offset $j \in [-M, M - 1]$. Additionally, a special parameter θ_0 is used to set the average read rate of the genome globally.

The DNase-seq covariate κ is defined as a vector of length N that is aligned with each base of the genome, and we assume that ChIP-seq reads can be predicted with this covariate and the contributions from surrounding k-mers. The regression coefficients for the covariate are defined as β and have length $2L$. The regression coefficients β can be thought of as analogous to the k-mer effect θ , but occurring everywhere, and scaled by the covariate κ . In all the experiments in this analysis, we chose an $L = 200$ to balance between computational complexity and prediction power.

Given these definitions, we define a generative model for ChIP-seq reads on the genome. Observed counts at position i on the genome are generated from a Poisson distribution with rate parameter λ_i which is defined as:

$$\lambda_i = \exp \left(\left(\sum_{k \in [1, k_{max}]} \sum_{j \in [-M, M-1]} \theta_{(g_{i+j}^k, -j)}^k \right) + \left(\sum_{l \in [0, 2L]} \beta_l \times \kappa_{i+l-L} \right) - \theta_0 \right) \quad (1)$$

The problem we solve is a regularized Poisson regression. Particularly, we would like to maximize the following:

$$\max_{\theta, \beta} \left\{ \sum_i c_i \log(\lambda_i) - \lambda_i - \eta \sum \|\theta^k\|_1 \right\} \quad (2)$$

To efficiently optimize this objective function, we performed an accelerated gradient descent method. The detail of implementation can be found in the supplementary data (Supplementary Text S1).

2.3 ChIP-seq Signal Prediction for Reference and Alternate Allele

In step 2, given the effect profile of all the k-mers and coefficients of the DNase-seq covariate trained from step 1, we first predict the ChIP-seq count λ at each position across the reference genome by

162 combining the effect of proximal k-mers and DNase-seq level into the log-linear model using equa-
163 tion 1. Then in similar manner, we predict the read counts λ' of the alternate allele after replacing
164 the k-mers that are affected by the variant. If we assume a Single Nucleotide Polymorphism (SNP),
165 at most $\frac{4}{3} \times (4^{k_{max}} - 1)$ k-mers will change.

167 2.4 Variant Scoring

168
169 In step 3, we score a SNP at locus on the genome by the square root of the sum of squared per-base
170 change (l^2 -norm of the change) of binding signal at all bases within the effect range of any k-mers
171 affected by the variant:

$$172 \quad 173 \quad 174 \quad 175 \quad 176 \quad 177 \quad s_i = \sqrt{\sum_{j \in [-M - k_{max} + 1, M - 1]} (\lambda'_{i+j} - \lambda_{i+j})^2} \quad (3)$$

178 2.5 Collapsing GERV k-mers into a position weight matrix

179 We interpret the active k-mers captured by GERV with a post-processing framework that aggregates
180 similar k-mers into position weight matrixes (PWMs):

- 181 1. We filter k-mers based on the sum of spatial effect to eliminate inactive k-mers.
- 182 2. We calculate the Levenshtein distance (number of single character edits) between the re-
183 maining k-mers.
- 184 3. We perform UPGMA hierarchical clustering over the candidate k-mers until the minimal
185 distance among clusters is larger than 2.
- 186 4. For each cluster, we define its key k-mer as the one with the largest sum of spatial effect.
187 We obtain the position weight matrix for this cluster by aligning all k-mers in the cluster
188 against the key k-mer.
- 189 5. All the clusters are ranked by the average sum of spatial effect of all the k-mers in the
190 cluster.

191 2.6 ChIP-seq Peak Prediction Comparison

192 Gapped-kmer SVM was downloaded from [http://www.beerlab.org/gkmsvm/](http://www.beerlab.org/gkmsvm/index.html)
193 index.html. To match with the training data for GERV, the positive training set for gapped-kmer
194 SVM consists of the all the NF- κ B ChIP-seq peaks on chr1-13 of GM12878 from ENCODE, and
195 the negative training set consists of the same number of randomly sampled regions of similar size
196 on chr1-13. The default parameter set (“-d 3”) was used. Both GERV and gapped-kmer SVM were
197 evaluated on the same test set. The positive test set consists of all the NF- κ B ChIP-seq peaks on
198 chr14-22 of GM12878 from ENCODE, and the negative test set consists of the same number of
199 randomly sampled regions of similar size on chr14-22.

200 2.7 Benchmark the performance in prioritizing SNPs with allele-specific binding

201 2.7.1 deltaSVM

202 deltaSVM source code was downloaded from [http://](http://www.beerlab.org/deltasvm/)
203 www.beerlab.org/deltasvm/. For each transcription factor included in the benchmarking, a gapped-
204 kmer SVM model was trained using ChIP-seq peaks of that factor on chr1-13 of GM12878 from
205 ENCODE as positive sets and the same number of randomly sampled region of similar size on chr1-
206 13 as negative sets. The default parameter set (“-d 3”) was used. As instructed by the software, the
207 gapped-kmer SVM model was then used to score all the possible 10-mers, the result of which was
208 input as the kmer-weight parameter to deltaSVM.

209 2.7.2 sTRAP

210 We used the R version of sTRAP downloaded from the website
211 (http://trap.molgen.mpg.de/download/TRAP_R_package/) for scalability. The built-in JASPAR and
212

216 TRANSFAC motif data included in the package were used. Specifically, MA0105.1, MA0105.2,
217 MA0105.3, MA0107.1, MA0061.1, V\$NFKAPPAB_01, V\$NFKB_Q6, V\$NFKAPPAB65_01,
218 V\$NFKAPPAB50_01, V\$P50_Q6, V\$NFKB_C and V\$RELA_Q6 were used for NF- κ B. MA0139.1,
219 MA0531.1, V\$CTCF_01, V\$CTCF_02 were used for CTCF. MA0099.1, MA0099.2, MA0476.1
220 and V\$CFOS_Q6 were used for FOS. MA0059.1, MA0058.1, MA0058.2, PB0043.1, PB00147.1,
221 V\$MAX_01, V\$MAX_04, V\$MAX_Q6, V\$MYCMAX_01, V\$MYCMAX_02, V\$MYCMAX_03
222 and V\$MYCMAX
223 _B were used for MAX. MA0059.1, MA0147.1, MA0147.2, V\$CMYC_01, V\$CMYC_02,
224 V\$MYC_01, V\$MYCMAX_01, V\$MYCMAX_02, V\$MYC
225 MAX_03 and V\$MYCMAX_B were used for MYC. None of the JUND motifs were included in the
226 built-in motif database of sTRAP. For each variant, the scores from different matrices of the same
227 factor were combined by taking the highest one.

228 229 **3 Materials**

230 231 **3.1 ChIP-seq Data**

232
233 ChIP-seq data for all the factors used in this analysis were downloaded from ENCODE. The full list
234 of GEO accession numbers can be found in Supplementary Table S1.

235 236 **3.2 DNase-seq Data**

237
238 DNase-seq data of GM12878 were downloaded from ENCODE (GEO accession GSM816665)

239 240 **3.3 Allele-Specific Binding (ASB) SNPs**

241
242 As a gold standard for SNPs that affect TF binding, we used the list of SNPs that are reported to
243 induce allele-specific binding (ASB) of NF- κ B, CTCF, FOS, JUND, MAX and MYC in GM12878.
244 The NF- κ B ASB SNPs are collected from [25] and [13]. The ASB SNPs data for all other tran-
245 scription factors are collected from [25]. The ASB SNPs were further filtered on allele frequency
246 (≥ 0.01) to keep only the common SNPs.

247 248 **4 Results**

249 250 **4.1 GERV learns a vocabulary of k-mers that regulate factor binding**

251
252 We first tested if GERV could predict held-out ChIP-seq data. We trained a GERV model on EN-
253 CODE NF- κ B ChIP-seq data and DNase-seq data from chromosomes 1-13 of GM12878, and com-
254 pared the predicted ChIP-seq signal from GERV to actual ChIP-seq reads on the held-out chro-
255 somes 14-22. The predicted ChIP-seq signals are very similar to actual ChIP-seq reads (Fig-
256 ure 2A,B), with a chromosome-wide Pearson's correlation of 0.76. We measured correlation after
257 smoothing predicted and actual reads over 400 bp windows since actual reads are insufficiently
258 sampled to produce base-pair resolution measurements. To further examine the ability of GERV to
259 model ChIP-seq peaks, we used the GERV model trained above to score a positive set of regions
260 defined as all the ENCODE GM12878 NF- κ B ChIP-seq peaks on chr14-22, and a negative set of
261 regions defined as same number of randomly sampled region of similar length on chr14-22. Each
262 region was scored by the sum of predicted signal in the region. We compared GERV with a previ-
263 ously published kmer-based model for TF peak prediction by training a gapped-kmer SVM ([9]) on
264 ENCODE NF- κ B peaks and same number of randomly sampled region of similar length on chr1-
265 13 of GM12878, and then performing the same scoring task on the same positive and negative set.
266 We quantified the performance of these two models in prioritizing positive regions over negative
267 regions by calculating the area under receiver operating characteristic curve (AUROC)(Figure 2C).
268 Our model achieved a better AUROC of 0.972 than that of 0.949 for gapped-kmer SVM. Thus GERV
269 learns a vocabulary of k-mers that can accurately predict the ChIP-seq data.

269 Although GERV fits a model with a potentially large parameter space (± 200 bp window for 87380
k-mers when $k_{max} = 8$), it uses sparsifying regularization to avoid overfitting and to limit the

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

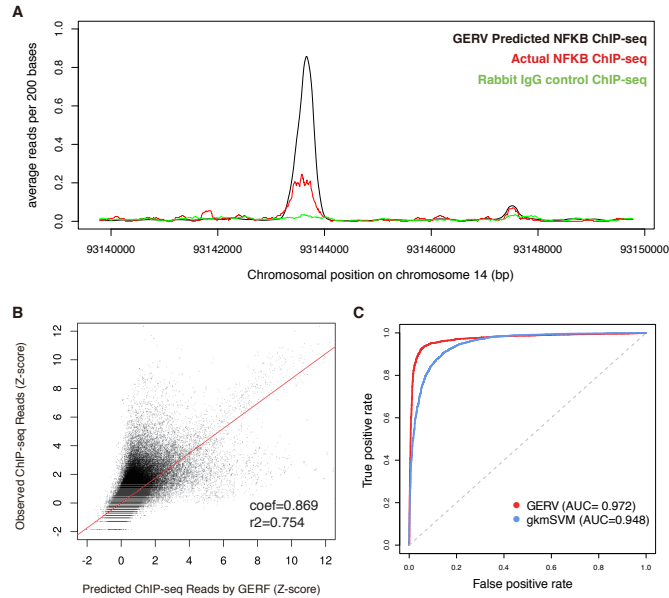


Figure 2: (A) Example held-out genomic region on chromosome 14 showing GERV-predicted NF- κ B reads (black), actual NF- κ B ChIP-seq reads (red), and rabbit IgG control ChIP-seq reads (green). (B) Comparison of GERV-predicted (x-axis) and observed (y-axis) NF- κ B ChIP-seq reads in binned regions of held-out chromosome 14-22. The coefficient and r^2 of a linear regression on predicted and actual z-score is plotted. (C) ROC curve for discriminating NF- κ B peaks from negative control sets using GERV and gapped-kmer SVM (gkmSVM).

number of active k-mers (Equation 2). For example, in the NF- κ B GERV model, most of the l^1 -norm of the parameter matrix is contained in the top 1% of the 87380 k-mers (Supplementary Figure S1). GERV is also robust to the choice of the window size for a k-mer's spatial effect and DNase-seq covariates (Supplementary Table S2).

4.2 GERV captures the binding sequence of a TF and its co-factors

We then examined if GERV learned the sequence features important for transcription factor binding. We trained a GERV model on DNase-seq data and NF- κ B ChIP-seq data combined from 10 LCL individuals. Position weight matrices were generated for visualization purposes by hierarchical clustering of the active k-mers in GERV (Section 2.5) and matched to known TF motifs in JASPAR and TRANSFAC using STAMP ([18]). With a threshold of significant matching at $1e^{-7}$, many clusters of the active k-mers correspond to known motifs (Table 1). The top two k-mer clusters for NF- κ B were matched to motifs from NF- κ B family (Supplementary Figure S2A), indicating that GERV correctly learned the strongest expected sequence features for the binding. Moreover, many of the other k-mer clusters learned by GERV correspond to transcription factors which have been associated with NF- κ B regulation (Supplementary Figure S2B), including ETS1, AP1, IRF1 and SP1 ([26, 8, 2, 29]). To validate the role of these transcription factors in NF- κ B binding, we performed co-factor analysis on the same NF- κ B data using GEM ([10]) to search for transcription factors that have spatially binding constraint with NF- κ B. This analysis identified AP-1 and IRF1 as the strongest co-factors of NF- κ B binding. Interestingly, some of the active-kmer clusters in GERV were matched to transcription factors such as ELF1, ERF2, CTCF and SUT1 which have not been associated with NF- κ B binding in previous studies.

To further interpret the role of the transcription factors whose motifs were matched to an active-kmer clusters in the NF- κ B GERV model, we performed motif analysis on the SNPs known to alter transcription factor binding. Allele-specific binding (ASB) studies have identified SNPs associated with significantly imbalanced binding events on heterozygous sites ([25, 13]). Therefore we collected a list of 56 ASB SNPs for NF- κ B, and use HaploReg ([31]) to query for the motifs that these ASB SNPs altered (Supplementary Table S3). Among the 56 ASB SNPs tested, only 16 (29%) were

Cluster	Matched Motif	Motif Database	Matched TF	E-value
PWM1	M00053 MA0101.1	TRANSFAC JASPAR	REL REL	5.1842e-08 1.2145e-09
PWM2	M00053 MA0101.1	TRANSFAC JASPAR	REL REL	7.0388e-14 1.1385e-12
PWM3	M00495 MA0099.2	TRANSFAC JASPAR	Bach1 API	8.0813e-13 9.4186e-10
PWM4	M01111	TRANSFAC	RBP-Jkappa	6.0791e-08
PWM5	M00339 MA0080.2	TRANSFAC JASPAR	ETS1 SPI1	1.3508e-11 3.6387e-10
PWM7	M01057 MA0123.1	TRANSFAC JASPAR	ERF2 abi4	2.0724e-08 1.9650e-08
PWM12	MA0139.1	JASPAR	CTCF	1.1289e-08
PWM15	M00062 MA0050.1	TRANSFAC JASPAR	IRF1 IRF1	2.7655e-10 3.1444e-09
PWM18	MA0399.1	JASPAR	SUT1	2.7097e-08
PWM20	M00722	TRANSFAC	core-binding	3.6831e-09
PWM22	MA0242.1	JASPAR	run_Bgb	3.1286e-11
PWM23	M01066	TRANSFAC	BLIMP1	1.4886e-09
PWM27	MA0453.1	JASPAR	nub	4.2762e-09
PWM32	M00345	TRANSFAC	GAMYB	8.8633e-08
PWM33	MA0344.1	JASPAR	NHP10	2.3002e-09
PWM38	MA0403.1	JASPAR	TBF1	5.6499e-08
PWM43	M00181	TRANSFAC	E2	2.9261e-08
PWM49	MA0152.1	JASPAR	NFATC2	5.3905e-11

Table 1: TF motifs matched to active-kmer clusters in NF- κ B GERV model using STAMP with E-value cutoff of $1e-07$. For each cluster, only the strongest match in each motif database (TRANSFAC and JASPAR) is shown. PWMs are ordered by the average sum of spatial effect of all the k-mers in the corresponding cluster.

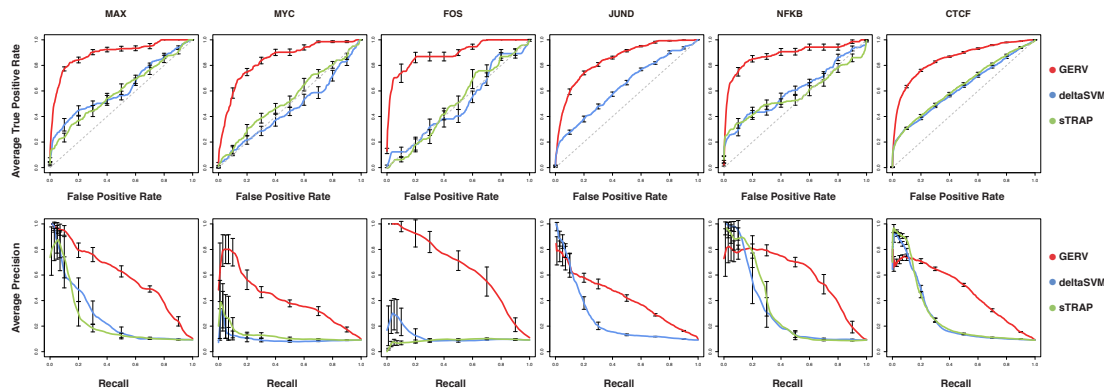


Figure 3: ROC (first row) and PRC (second row) curves for discriminating ASB SNPs from the second type of negative variant set (10 times of the size of positive set) using GERV (red), deltaSVM (blue), and sTRAP (green). Grey dashed line in ROC curves indicates random chance. In each figure, 95% confidence intervals of the true positive rate (for ROC) or precision (for PRC) are plotted. The performance of sTRAP on JUND is not measurable as JUND motif is not included in its built-in motif database.

found to alter the canonical motif of NF- κ B, while another 11 (20%) were found to alter the TF motif matched to other active-kmer clusters in the GERV model. Thus, GERV captures the sequence context of factor binding, which provides additional descriptive power and biological insight for auxiliary elements in TF binding.

4.3 GERV outperforms existing approaches in prioritizing ASB SNPs

To demonstrate the power of GERV in detecting regulatory variants, we compared GERV's performance against existing approaches in discriminating ASB SNPs from negative control variants. We collected ASB SNPs with known differential binding for NF- κ B, CTCF, JUND, MAX, MYC and FOS from previous studies ([25, 13]) as positive sets, resulting in a total of 56 SNPs for NF- κ B, 1225 SNPs for CTCF, 26 SNPs for FOX, 233 SNPs for JUND, 71 SNPs for MAX and 69 SNPs for MYC (Section 3.3). Note that these ASB SNPs were completely held-out in the training process of any model compared in this analysis, and were only used as the test set.

For each of the six transcription factors, we constructed two types of negative SNP sets that we assume do not exhibit differential factor binding. Both kinds of negative sets are subsets of 1000 Genome Project (1KG) common (minor allele frequency $\geq 1\%$) SNPs. In the first case we randomly sampled 100 negative samples for each positive sample, to get a reasonable sample of the background while making analyses computationally tractable. The second set is a fine-mapping task which is an important topic in post-GWAS analysis where a list of lead SNPs and their linked SNPs are under interrogation for regulatory consequence. To simulate such tasks, this second set was constructed as random selection of 1KG common SNPs within 10kb from any ASB SNP. To reflect the number of SNPs typically in a single LD block, we calculated LD information from phased genotype data in the 1KG pilot release using PLINK ([23]). With a r^2 cutoff of 0.8, the median number of linked SNPs for a variant is 10 (Supplementary Figure S3). Thus in this set we sampled 10 negative samples for each positive sample. For both types of negative sets, we sampled 10 sets with replacement so that we could obtain the mean and confidence intervals. For each of the 10 negative sets, we constructed a paired positive set, same size as the corresponding ASB SNP set, by sampling with replacement from the ASB SNPs.

For each transcription factor, we evaluated the performance of GERV as well as two published regulatory variant scoring methods sTRAP ([19]) (motif-based), and deltaSVM ([14]) (kmer-based) in discriminating the positive set from each of the two negative sets. The other motif-based methods are not included due to either the inability to produce numerical scores for the queried variants, or the low throughput that can't scale up to thousands of SNPs. For each factor, a GERV model was trained on ENCODE ChIP-seq data from chr1-13 of GM12878, and a deltaSVM model was trained on ENCODE ChIP-seq peaks and same number of random regions of similar length on chr1-13 of GM12878. The built-in JASPAR and TRANSFAC motif dataset was used for sTRAP, which includes the motif for all the factors but JUND (Section 2.7).

We show the averaged receiver operating characteristic (ROC) curves and precision recall curves (PRC) (Supplementary Figure S4 for the first control set, Figure 3 for the second control set) of all the methods for different transcription factors and negative sets. We evaluated two aspects of the curves. The first metric is the area under curve (AUC) (Supplementary Table S4) which summarizes the overall performance in prioritizing the positive set over negative set. The second metric is the true positive rate at low false positive rate (for ROC) or the recall at high precision (for PRC), which reflects the practical need for low false discovery rate in post-GWAS analysis where thousands of lead and linked SNPs are tested for regulatory consequence. The ROC curves for GERV consistently dominated the competing methods for all factors and control scenarios, with much better AUC and higher true positive rate at low false positive rates. In PR curves because of the small size of the positive set, the confidence intervals of precision when the recall is low tend to be large, making the left-most part of the curves less informative for comparison. For transcription factor FOS, MAX and MYC, GERV achieved a PRC curve clearly superior to the others, without overlapping in the confidence interval. For factor JUND and NF- κ B, GERV had a similarly high precision for low recall, but outperformed the other methods with consistently high precision for larger recall. For CTCF, the competing methods achieved higher precision when recalling less than 10% of the positives, but their precision dropped dramatically afterwards resulting in much lower AUC than that of GERV. Given the fact that CTCF has a motif (19 bp) more than twice as long as the maximum length of k-mer (8 bp) learnable for GERV (Section 2.2), the competitive performance on CTCF demonstrates the strong descriptive power of GERV in modeling TF binding.

We found that for our second negative control setup, choosing 50 instead of 10 negative SNPs for each positive SNP caused a noticeable decline in the area under precision recall curves (Supplementary Figure S5). With this ratio of positive to negative SNPs, deltaSVM outperformed GERV at the

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

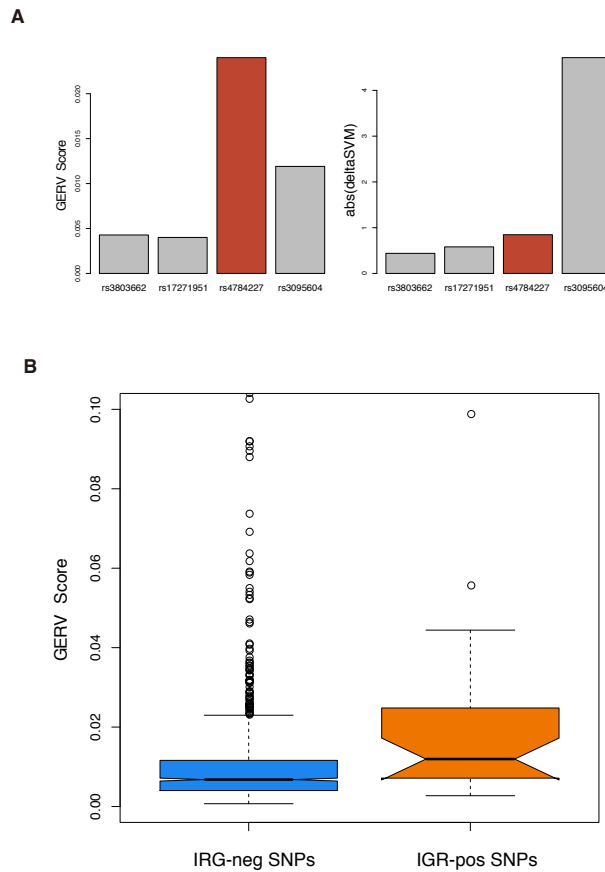


Figure 4: (A) GERV correctly predicted the effect of validated causal SNP rs4784227 on FOXA1 binding, while deltaSVM failed. (B) The 28 variants previously reported to modulate FOXA1 binding had significantly higher (Mann-Whitney U test $P=0.001$) GERV scores than the rest of the AVS ($n=1065$ after filtering for common SNPs)

10% recall point for JUND, NF- κ B, and CTCF, but by the 30% recall point GERV produced better precision-recall for these and all other factors.

4.4 GERV prioritizes linked-SNPs that modulate FOXA1 binding in breast cancer

To demonstrate the application of GERV in post-GWAS analysis, we applied GERV to a breast cancer associated variant set (AVS) collected from a previous study ([5]). It is composed of 44 risk-associated SNPs and 1,053 linked SNPs that are in strong LD with any risk-associated SNP. It has been shown that breast cancer associated SNPs are enriched for the binding sites of FOXA1, a pioneer transcription factor essential for chromatin opening and nucleosome positioning favorable to transcription factor recruitment ([11, 6, 3, 4, 16]). We trained a GERV model and a deltaSVM model on ENCODE FOXA1 ChIP-seq data from a breast cancer cell line T47D. The rs4784227 breast cancer associated SNP has been shown to disrupt the binding of FOXA1 with several lines of evidence ([5, 15]). GERV correctly predicted the effect of rs4784227 on FOXA1 binding among its linked SNPs, while deltaSVM failed (Figure 4A). Having probed a single risk-associated SNP, we then applied GERV to all the SNPs in the breast cancer AVS. The 28 variants previously reported to modulate FOXA1 binding ([5]) had significantly higher GERV scores than the rest of the AVS (Figure 4B, Mann-Whitney U test $P=0.001$, $AUC=0.69$). In contrast, deltaSVM couldn't distinguish the positive set from the rest of the AVS (Mann-Whitney U test $P=0.22$, $AUC=0.57$)

486 5 Discussion

487

488 Despite the recent substantial advances in characterizing the genome-wide transcription factor bind-
489 ing sites with ChIP-seq experiments, it remains a challenge to interpret variation in the noncoding
490 region of the genome and to determine variants that cause transcription factor binding changes in
491 post-GWAS analysis. Our work improves the prediction of causal non-coding variants when com-
492 pared to other contemporary methods.

493 As the first generative model that directly predicts the ChIP-seq signal, GERV achieved greater
494 accuracy than other methods in predicting ChIP-seq peaks. GERV models the spatial effect of all
495 the k-mers and thus captures the effect of the primary motif and auxiliary sequences on TF binding.
496 We have shown that many of these auxiliary sequences correspond to known binding cofactors,
497 while others were matched to transcription factors whose roles in the binding regulation have not
498 been previously characterized.

499 The generative nature of the GERV model scores each variant as the predicted change to a proximal
500 ChIP-seq signal. The analysis on six transcription factors NF- κ B, CTCF, FOS, JUND, MAX and
501 MYC demonstrated that GERV outperforms existing methods in discriminating variants known to
502 alter TF binding from negative control sets. In a few cases (Figure 3F, Supplementary Figure S4F),
503 the discriminative nature of the competing methods equipped them with higher precision for recall-
504 ing a small fraction of positives. However their inability to model auxiliary sequences led to the
505 dramatic precision decrease afterwards, while GERV achieved constantly high precision for larger
506 recall.

507 Applied to an associated variant set (AVS) of breast cancer, GERV correctly predicted the effect of
508 previous validated causal SNP rs4784227, and highly prioritized variants reported to affect FOXA1
509 binding in breast cancer cell line. With the superior performance exemplified in this task, we expect
510 GERV to play an important role in functionally annotating and prioritizing putative causal variants
511 for downstream experimental analysis

512

513 Acknowledgement

514

515 We thank Yuchun Guo for technical support in co-factor analysis using GEM. We also thank
516 Matthew Edwards for many helpful comments and discussions.

517

518 Funding

519

520 This work was supported by the National Institutes of Health [1U01HG007037 to D.K.G.]

521

522 References

523

- 524 [1] Malin C Andersen, Pär G Engström, Stuart Lithwick, David Arenillas, Per Eriksson, Boris
525 Lenhard, Wyeth W Wasserman, and Jacob Odeberg. In silico detection of sequence variations
526 modifying transcriptional regulation. *PLoS computational biology*, 4(1):e5, January 2008.
- 527 [2] Myriam Bartels, Aike Torben Schweda, Ursula Dreikhausen, Ronald Frank, Klaus Resch, Win-
528 fried Beil, and Mahtab Nourbakhsh. Peptide-mediated disruption of NFkappaB/NRF inter-
529 action inhibits IL-8 gene activation by IL-1 or Helicobacter pylori. *Journal of immunology*
530 (*Baltimore, Md. : 1950*), 179:7605–7613, 2007.
- 531 [3] Jason S Carroll, X Shirley Liu, Alexander S Brodsky, Wei Li, Clifford A Meyer, Anna J Szary,
532 Jerome Eeckhoutte, Wenlin Shao, Eli V Hestermann, Timothy R Geistlinger, Edward A Fox,
533 Pamela A Silver, and Myles Brown. Chromosome-wide mapping of estrogen receptor binding
534 reveals long-range regulation requiring the forkhead protein FoxA1. *Cell*, 122(1):33–43, July
535 2005.
- 536 [4] Jason S Carroll, Clifford A Meyer, Jun Song, Wei Li, Timothy R Geistlinger, Jérôme Eeck-
537 houte, Alexander S Brodsky, Erika Krasnickas Keeton, Kirsten C Fertuck, Giles F Hall,
538 Qianben Wang, Stefan Bekiranov, Victor Sementchenko, Edward A Fox, Pamela A Silver,
539 Thomas R Gingeras, X Shirley Liu, and Myles Brown. Genome-wide analysis of estrogen
receptor binding sites. *Nature Genetics*, 38(11):1289–1297, October 2006.

- 540 [5] Richard Cowper-Salari, Xiaoyang Zhang, Jason B Wright, Swneke D Bailey, Michael D Cole,
541 Jerome Eeckhoute, Jason H Moore, and Mathieu Lupien. Breast cancer risk-associated SNPs
542 modulate the affinity of chromatin for FOXA1 and alter gene expression. *Nature genetics*,
543 44(11):1191–8, November 2012.
- 544 [6] Jérôme Eeckhoute, Jason S Carroll, Timothy R Geistlinger, Maria I Torres-Arzayus, and Myles
545 Brown. A cell-type-specific transcriptional network required for estrogen regulation of cy-
546 clin D1 and cell cycle progression in breast cancer. *Genes & development*, 20(18):2513–26,
547 September 2006.
- 548 [7] Kelly a Frazer, Sarah S Murray, Nicholas J Schork, and Eric J Topol. Human genetic variation
549 and its contribution to complex traits. *Nature reviews. Genetics*, 10(4):241–51, April 2009.
- 550 [8] Shuichi Fujioka, Jiangong Niu, Christian Schmidt, M Guido, Bailu Peng, Tadashi Uwagawa,
551 Zhongkui Li, Douglas B Evans, James L Abbruzzese, Paul J Chiao, and Guido M Sclabas. NF- κ B
552 and AP-1 Connection: Mechanism of NF- κ B-Dependent Regulation of AP-1 Activity.
553 *Society*, 24(17):7806–7819, 2004.
- 554 [9] Mahmoud Ghandi, Dongwon Lee, Morteza Mohammad-Noori, and Michael A Beer. Enhanced
555 regulatory sequence prediction using gapped k-mer features. *PLoS computational biology*,
556 10(7):e1003711, July 2014.
- 557 [10] Yuchun Guo, Shaun Mahony, and David K Gifford. High resolution genome wide binding
558 event finding and motif discovery reveals transcription factor spatial binding constraints. *PLoS*
559 *computational biology*, 8(8):e1002638, January 2012.
- 560 [11] Housheng Hansen He, Clifford A Meyer, Hyunjin Shin, Shannon T Bailey, Gang Wei, Qianben
561 Wang, Yong Zhang, Kexin Xu, Min Ni, Mathieu Lupien, Piotr Mieczkowski, Jason D Lieb,
562 Keji Zhao, Myles Brown, and X Shirley Liu. Nucleosome dynamics define transcriptional
563 enhancers. *Nature genetics*, 42(4):343–7, April 2010.
- 564 [12] Lucia a Hindorff, Praveen Sethupathy, Heather a Junkins, Erin M Ramos, Jayashri P Mehta,
565 Francis S Collins, and Teri a Manolio. Potential etiologic and functional implications of
566 genome-wide association loci for human diseases and traits. *Proceedings of the National*
567 *Academy of Sciences of the United States of America*, 106(23):9362–7, June 2009.
- 568 [13] Konrad J Karczewski, Joel T Dudley, Kimberly R Kukurba, Rong Chen, Atul J Butte,
569 Stephen B Montgomery, and Michael Snyder. Systematic functional regulatory assessment
570 of disease-associated variants. *Proceedings of the National Academy of Sciences of the United*
571 *States of America*, 110(23):9607–12, June 2013.
- 572 [14] Dongwon Lee, David U Gorkin, Maggie Baker, Benjamin J Strober, Alessandro L Asoni,
573 Andrew S McCallion, and Michael a Beer. A method to predict the impact of regulatory
574 variants from DNA sequence. *Nature Genetics*, (June), 2015.
- 575 [15] Jirong Long, Qiuyin Cai, Xiao-Ou Shu, Shimian Qu, Chun Li, Ying Zheng, Kai Gu, Wenjing
576 Wang, Yong-Bing Xiang, Jiarong Cheng, Kexin Chen, Lina Zhang, Hong Zheng, Chen-Yang
577 Shen, Chiun-Sheng Huang, Ming-Feng Hou, Hongbing Shen, Zhibin Hu, Furu Wang, Sandra L
578 Deming, Mark C Kelley, Martha J Shrubsole, Ui Soon Khoo, Kelvin Y K Chan, Sum Yin
579 Chan, Christopher A Haiman, Brian E Henderson, Loic Le Marchand, Motoki Iwasaki, Yoshio
580 Kasuga, Shoichiro Tsugane, Keitaro Matsuo, Kazuo Tajima, Hiroji Iwata, Bo Huang, Jiajun
581 Shi, Guoliang Li, Wanqing Wen, Yu-Tang Gao, Wei Lu, and Wei Zheng. Identification of a
582 functional genetic variant at 16q12.1 for breast cancer risk: results from the Asia Breast Cancer
583 Consortium. *PLoS genetics*, 6(6):e1001002, June 2010.
- 584 [16] Mathieu Lupien, Jérôme Eeckhoute, Clifford A. Meyer, Qianben Wang, Yong Zhang, Wei Li,
585 Jason S. Carroll, X. Shirley Liu, and Myles Brown. FoxA1 Translates Epigenetic Signatures
586 into Enhancer-Driven Lineage-Specific Transcription. *Cell*, 132(6):958–970, March 2008.
- 587 [17] Geoff Macintyre, James Bailey, Izhak Haviv, and Adam Kowalczyk. is-rSNP: a novel tech-
588 nique for in silico regulatory SNP detection. *Bioinformatics (Oxford, England)*, 26(18):i524–
589 30, September 2010.
- 590 [18] Shaun Mahony and Panayiotis V. Benos. STAMP: A web tool for exploring DNA-binding
591 motif similarities. *Nucleic Acids Research*, 35:253–258, 2007.
- 592 [19] Thomas Manke, Matthias Heinig, and Martin Vingron. Quantifying the effect of sequence
593 variation on regulatory interactions. *Human mutation*, 31(4):477–83, April 2010.

- 594 [20] Teri A Manolio. Genomewide association studies and assessment of the risk of disease. *The*
595 *New England journal of medicine*, 2010.
- 596 [21] Mark I McCarthy, Gonçalo R Abecasis, Lon R Cardon, David B Goldstein, Julian Little, John
597 P a Ioannidis, and Joel N Hirschhorn. Genome-wide association studies for complex traits:
598 consensus, uncertainty and challenges. *Nature reviews. Genetics*, 9(5):356–69, May 2008.
- 599 [22] Ivan Molineris, Davide Schiavone, Fabio Rosa, Giuseppe Matullo, Valeria Poli, and Paolo
600 Provero. Identification of functional cis-regulatory polymorphisms in the human genome. *Hu-*
601 *man mutation*, 34(5):735–42, May 2013.
- 602 [23] Shaun Purcell, Benjamin Neale, Kathe Todd-Brown, Lori Thomas, Manuel A R Ferreira, David
603 Bender, Julian Maller, Pamela Sklar, Paul I W de Bakker, Mark J Daly, and Pak C Sham.
604 PLINK: a tool set for whole-genome association and population-based linkage analyses. *Amer-*
605 *ican journal of human genetics*, 81(3):559–75, September 2007.
- 606 [24] Alberto Riva. Large-scale computational identification of regulatory SNPs with rSNP-
607 MAPPER. *BMC genomics*, 13 Suppl 4(Suppl 4):S7, January 2012.
- 608 [25] Joel Rozowsky, Alexej Abyzov, Jing Wang, Pedro Alves, Debasish Raha, Arif Harmanci, Jing
609 Leng, Robert Bjornson, Yong Kong, Naoki Kitabayashi, Nitin Bhardwaj, Mark Rubin, Michael
610 Snyder, and Mark Gerstein. AlleleSeq: analysis of allele-specific expression and binding in a
611 network framework. *Molecular systems biology*, 7(522):522, January 2011.
- 612 [26] Marco Sgarbanti, Anna L Remoli, Giulia Marsili, Barbara Ridolfi, Alessandra Borsetti, Edvige
613 Perrotti, Roberto Orsatti, Ramona Ilari, Leonardo Sernicola, Emilia Stellacci, Barbara Ensoli,
614 and Angela Battistini. IRF-1 is required for full NF-kappaB transcriptional activity at the
615 human immunodeficiency virus type 1 long terminal repeat enhancer. *Journal of virology*,
616 82(7):3632–3641, 2008.
- 617 [27] Barbara E Stranger, Eli a Stahl, and Towfique Raj. Progress and promise of genome-wide
618 association studies for human complex trait genetics. *Genetics*, 187(2):367–83, February 2011.
- 619 [28] Mingxiang Teng, Shoji Ichikawa, Leah R. Padgett, Yadong Wang, Matthew Mort, David N.
620 Cooper, Daniel L. Koller, Tatiana Foroud, Howard J. Edenberg, Michael J. Econs, and Yunlong
621 Liu. Regsnps: A strategy for prioritizing regulatory single nucleotide substitutions. *Bioinfor-*
622 *matics*, 28(14):1879–1886, 2012.
- 623 [29] R S Thomas, M J Tymms, L H McKinlay, M F Shannon, a Seth, and I Kola. ETS1, NFkappaB
624 and AP1 synergistically transactivate the human GM-CSF promoter. *Oncogene*, 14:2845–
625 2855, 1997.
- 626 [30] Lucas D Ward and Manolis Kellis. HaploReg: a resource for exploring chromatin states,
627 conservation, and regulatory motif alterations within sets of genetically linked variants. *Nucleic*
628 *acids research*, 40(Database issue):D930–4, January 2012.
- 629 [31] Lucas D Ward and Manolis Kellis. Interpreting noncoding genetic variation in complex traits
630 and human disease. *Nature biotechnology*, 30(11):1095–106, November 2012.
- 631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647