

Supplemental Materials

for

Prioritization of regulatory variants with tissue-specific function in the non-coding regions of human genome

Shengcheng Dong¹ and Alan P. Boyle^{1,2*}

1. Department of Computational Medicine and Bioinformatics, University of
Michigan, Ann Arbor, MI
2. Department of Human Genetics, University of Michigan, Ann Arbor, MI

Supplemental Tables

Supplemental Table S1. Number of allele specific TF binding (ASB) training SNVs in 6 cell lines.

Cell type name	# of ASB SNVs	# of control SNVs
GM12878	1,848	12,382
A549	215	1,357
H1hESC	1,464	6,525
HepG2	2,717	24,479
K562	767	12,635
MCF7	653	8,034

Supplemental Table S2. Feature list in random forest models of TURF.

Generic features	
TF binding sites from ChIP-seq	Binary variable
DNase I hypersensitive sites from DNase-seq	Binary variable
DNase footprints	Binary variable
DNase footprints with matched TF ChIP-seq peaks	Binary variable
TF motifs from PWM matching	Binary variable
TF motifs from PWM matching with matched TF ChIP-seq peaks	Binary variable
Information content change of two alleles in PWM matching	Numerical variable
Information content change of two alleles in PWM matching with matched TF ChIP-seq peaks	Numerical variable
eQTLs	Binary variable
Quantiles (25%,50%,75% and 100%) and variance of ChIP-seq signals across all available ChIP-seq experiments from ENCODE	Numerical variables
Functional significance score from <i>DeepSEA</i>	Numerical variable
Tissue-specific features (in final ensemble model)	
H3K4me1 peaks from ChIP-seq	Binary variable
H3K4me3 peaks from ChIP-seq	Binary variable
H3K27ac peaks from ChIP-seq	Binary variable
H3K36me3 peaks from ChIP-seq	Binary variable
H3K27me3 peaks from ChIP-seq	Binary variable
DNase I hypersensitive sites from DNase-seq	Binary variable
DNase footprints	Binary variable

Supplemental Table S3. Feature importance of TURF generic scores. The importance is calculated based on mean decrease of impurity from the random forest model.

Feature name	Importance
CHIP	0.124607
DNASE	0.033935
PWM	0.018686
FOOTPRINT	0.008609
EQTL_2	0.004167
PWM_matched	0.005055
FOOTPRINT_matched	0.001454
IC_change	0.050415
IC_matched_change	0.021396
funsig	0.122206
ChIP_var	0.149313
ChIP_max	0.167037
ChIP_quantile1	0.097849
ChIP_quantile2	0.090261
ChIP_quantile3	0.105011

Supplemental Table S4. Comparison of performance on cell type-specific predictions for MPRA variants and ASB SNVs.

Method	AUPR	AUROC	Correlation
Test on MPRA variants in E116 (GM12878)			
TURF_Tissue-Specific	0.548	0.731	0.450
TURF_Generic	0.536	0.740	0.447
GenoNet	0.536	0.723	0.442
DeepSEA	0.451	0.671	0.115
GenoSkyline	0.394	0.673	0.352
CADD	0.240	0.566	0.082
TURF_Ensemble (trained from ASB SNVs)	0.484	0.724	0.393
Test on MPRA variants in E118 (HepG2)			
TURF_Tissue-Specific	0.571	0.765	0.423
TURF_Generic	0.502	0.733	0.372
GenoNet	0.572	0.757	0.429
DeepSEA	0.401	0.667	0.166
GenoSkyline	0.326	0.628	0.225
CADD	0.332	0.607	0.152
TURF_Ensemble (trained from ASB SNVs)	0.479	0.728	0.357
Test on MPRA variants in E123 (K562)			
TURF_Tissue-Specific	0.441	0.681	0.344
TURF_Generic	0.400	0.664	0.294
GenoNet	0.420	0.706	0.340
DeepSEA	0.290	0.607	0.099

GenoSkyline	0.284	0.659	0.274
CADD	0.212	0.539	0.046
TURF_Ensemble (trained from ASB SNVs)	0.329	0.676	0.253

Test on ASB SNVs in GM12878

TURF_Ensemble	0.511	0.894	0.538
TURF_Tissue-Specific	0.510	0.882	0.537
TURF_Generic	0.284	0.822	0.390
GenoSkyline	0.246	0.743	0.340
DeepSEA	0.222	0.688	0.190
GenoNet	0.129	0.685	0.158

Test on ASB SNVs in HepG2

TURF_Ensemble	0.355	0.857	0.406
TURF_Tissue-Specific	0.353	0.831	0.401
TURF_Generic	0.221	0.797	0.323
DeepSEA	0.226	0.611	0.138
GenoSkyline	0.212	0.582	0.133
GenoNet	0.161	0.654	0.134

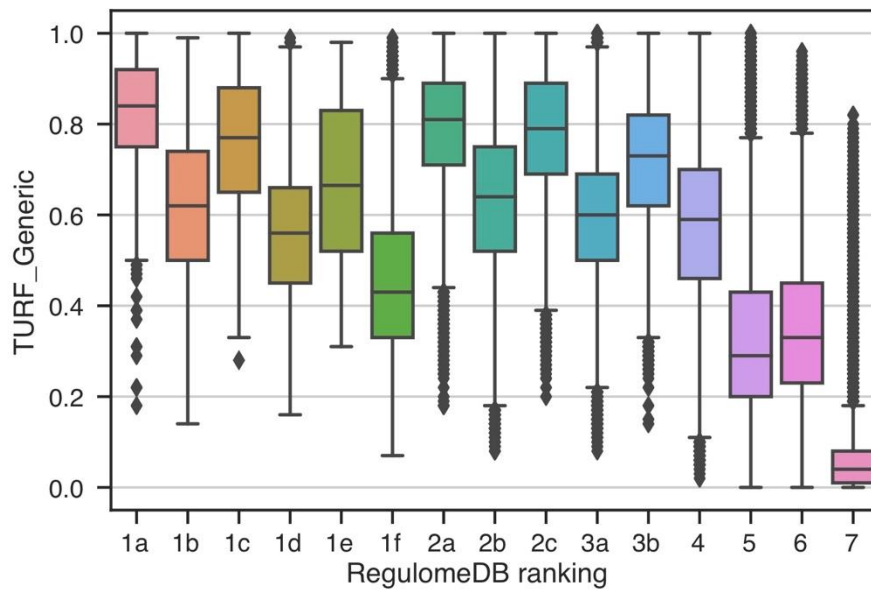
Test on ASB SNVs in K562

TURF_Ensemble	0.316	0.847	0.376
TURF_Tissue-Specific	0.313	0.815	0.366
TURF_Generic	0.193	0.782	0.301
GenoNet	0.140	0.644	0.113
GenoSkyline	0.090	0.628	0.144
DeepSEA	0.062	0.572	0.054

Test on ASB SNVs in A549

TURF_Ensemble	0.313	0.872	0.369
TURF_Tissue-Specific	0.315	0.833	0.369
TURF_Generic	0.158	0.796	0.256
DeepSEA	0.052	0.760	0.119
GenoNet	0.030	0.690	0.094
GenoSkyline	0.022	0.632	0.049

Supplemental Figures



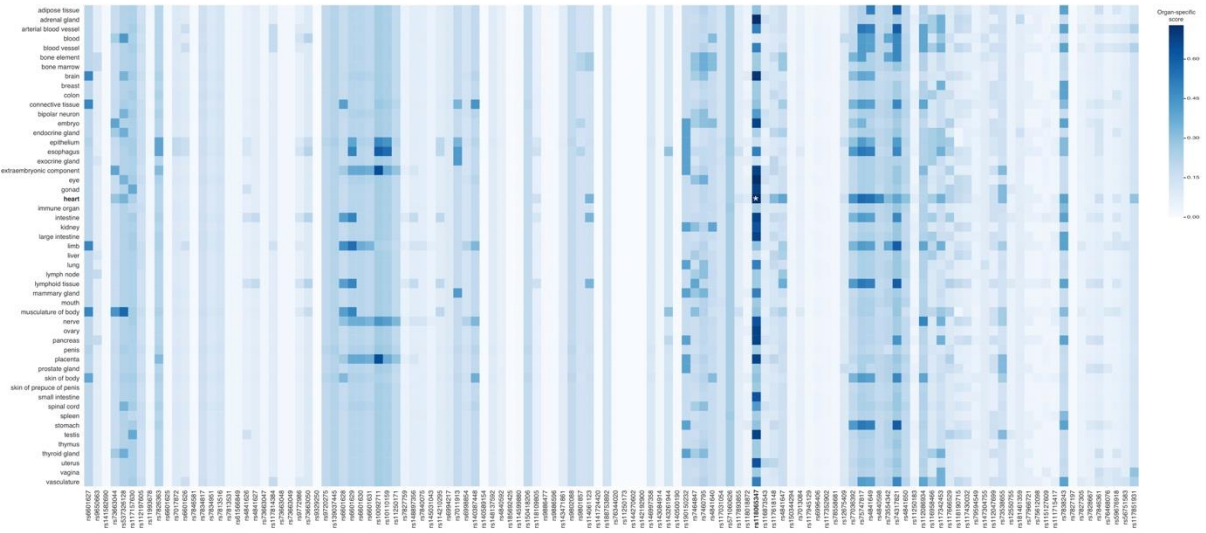
Supplemental Figure S1. Boxplot of TURF generic scores VS RegulomeDB ranking scores on 10,422,004 common SNVs from dbSNP153. X axis represents the original ranking scores from RegulomeDB v1.1, y axis represents the TURF generic scores from random forest model.



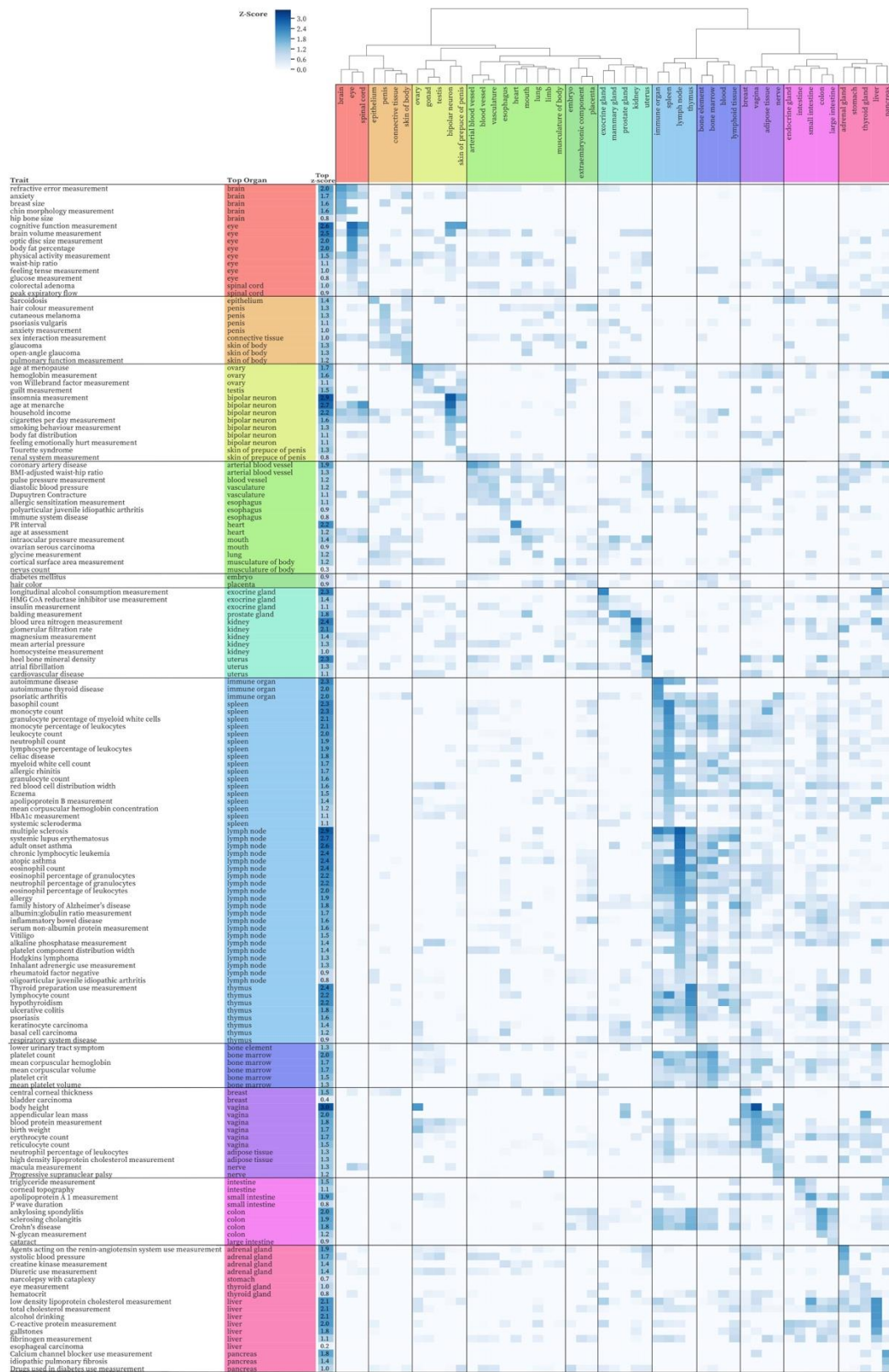
Supplemental Figure S2. Pearson correlation of labels and tissue-specific features in three MPRA datasets (E116: GM12878; E118: HepG2; E123: K562). Blue bars represent positive correlations while red bars represent negative correlations.



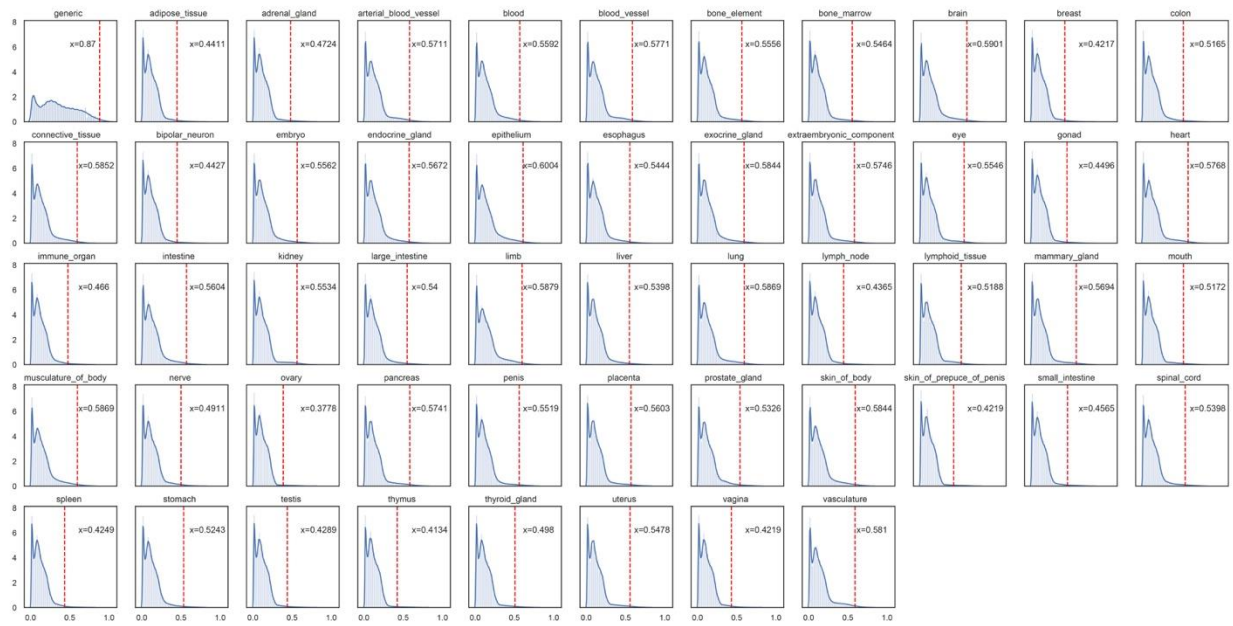
Supplementary Figure 3. Pearson correlation of labels and tissue-specific features in 6 ASB datasets. Blue bars represent positive correlations while red bars represent negative correlations.



Supplemental Figure S4. Organ-specific scores of variants in the *GATA4* locus. rs118065347 has the top organ-specific score in heart compared to other candidate SNPs found through genome-wide association scan on bicuspid aortic valve cases. rs118065347 was shown to have heart-specific regulatory functionality from a previous study Yang et al. 2017.



Supplemental Figure S5. Enrichment of regulatory variants with high organ-specific scores over variants associated with diverse traits. Only organ-trait pairs with z-scores higher than 0 and passing multiple test correction (FDR threshold of 5%) are shown.



Supplemental Figure S6. Distribution of TURF generic and organ-specific scores on 3,102,261 GM12878 SNVs. The red vertical dashed line represents the 99th percentile score in each subplot.