

## 1 SUPPLEMENTARY NOTE 1

2

3 **Regularized linear models are poor models for predicting mean abundance from core promoter**  
4 **sequences (10-fold cross-validated average oos- $r^2 < 0$ ).** To assess how simple linear models compare  
5 to “deep learning” models, we fitted linear models by minimizing a regularized empirical squared loss,  
6 with a L2 (squared euclidean norm) penalty, using stochastic gradient descent from python’s sci-kit  
7 learn (sklearn). The input of the class B models (non-specific annotated DNA model; 1D matrix with  
8 32 channels) was flattened prior to modelling, i.e. the mean abundance each gene was modelled as a  
9 linear function of 128,000 variables (most were binary variables). For brevity, we limited our analysis  
10 to skeletal muscle, the tissue with the largest number of samples in GTEx. The model was fit using the  
11 partial fit function from SGDRegressor with sklearn using default parameters; the same exit criteria as  
12 described in **Online Methods** were applied. We noted that 10-fold cross-validated average oos- $r^2$  was  
13 consistently below zero, with performance considerably worse than the peaBrain model. The gap in  
14 performance was considerable (precludes visualization). We also repeated this analysis using a single  
15 dense neural network layer with linear activations, and noted that the 10-fold cross-validated average  
16 oos- $r^2$  was below 0; a single dense layer with linear activations is equivalent to a simple linear model.

17

18 **Fully connected neural networks are slower and more memory-intensive, with slightly worse**  
19 **performance than convolutional neural networks; non-linear activations for convolutional layers**  
20 **improve performance over linear activations.** We compared fully connected dense neural networks  
21 to the convolutional neural networks at the heart of peaBrain (**Figure 1**) for class B skeletal muscle  
22 model. We replaced each pair of convolutional-pooling layers with a single dense layer; the number of  
23 neurons in the three dense layers was limited only by GPU memory. The penultimate layer and single  
24 output neuron were kept consistent between the models. Using the skeletal muscle Class B input, we  
25 noted that fully connected performed slightly worse (oos- $r^2 = 0.42$ ) than the convolutional neural  
26 networks (oos- $r^2 = 0.46$ ). Increasing the number of layers allows fully connected neural network to  
27 reach performance parity with CNNs, at the cost of increased memory and computational cost (not  
28 feasible for Stage 2 peaBrain analyses). We subsequently wanted to assess the importance of the non-

29 linear activation for the convolutional layers of peaBrain. We constructed an identical model, but  
30 replacing all CNN activations with linear functions and noted that performance for this model was  
31 consistently worse ( $\text{oos-r}^2 = 0.43$ ) than the classical peaBrain architecture ( $\text{oos-r}^2 = 0.46$ ). It is important  
32 to note that this is not an exhaustive search of the ideal set of parameters, but an exploratory analysis to  
33 begin to understand peaBrain's performance.

34

35 **DNA sequence, annotated with experimentally-derived TFBS, from core promoter sequences are**  
36 **insufficient to predict mean abundance with high accuracy – epigenetic/histone markers contain**  
37 **the bulk of the information and are not readily accessible from the DNA sequence alone.** We were  
38 interested in determining the contribution of epigenetic/histone makers, alongside more general  
39 genomic annotations (such as coding sequences), in predicting the mean abundance of genes. In  
40 particular, we wanted to explore whether the DNA sequence alone was sufficient to predict expression  
41 in skeletal muscle. We noted that increasing the number of convolutional layers or the number of filters  
42 did not improve model performance (**Figure 1**). Explicitly incorporating TFBS into the model (i.e.  
43 annotating the DNA only and explicitly with TFBS) only improved performance slightly ( $\text{oos-r}^2 =$   
44 23%), and was still considerably worse than the full class B model with epigenetic/histone marker  
45 annotations ( $\text{oos-r}^2 = 46\%$ ; **Figure 1**). (Class-A DNA-only models had an average  $\text{oos-r}^2$  of 16% for  
46 skeletal muscle; class-C models annotated with tissue-specific information had an average  $\text{oos-r}^2$  of  
47 57%.) The TFBS were collected from the Gene Transcription Regulation Database (GTRD) v17.4 with  
48 data on 476 human transcription factors and included peak calling with four different software (MACS,  
49 SISSRs, GEM, and PICS). In addition to including the processed peak calls, we also incorporated  
50 clusters (i.e. peaks merged for the same transcription factor but under different experimental conditions)  
51 and meta-clusters (i.e. non-redundant peaks synthesized from all four methods). This absence of  
52 improvement suggests that peaBrain model already recognizes many of the TFBS; identified by the  
53 convolutional filters inherent to the model architecture. These results indicate that experimentally-  
54 derived epigenetic and genomic annotations add information to that contained in the DNA sequence  
55 alone. As described in the main text, this is broadly consistent with the observation that other

56 convolutional neural networks models like DeepSEA are better at predicting TFBS (median AUC =  
57 0.958) than at predicting histone modifications (median AUC = 0.856)<sup>1</sup>.

58

59 **peaBrain score out-performs existing measures in predicting allele-specific transcription factor**

60 **binding.** As with **tasks A and B** (described in the **Main Text**), we compared the performance of the

61 non-tissue-specific peaBrain score to predictions by CADD and EIGEN in predicting allele-specific

62 binding, after accounting for allele frequency and evolutionary conservation. We assessed performance

63 of the three non-coding metrics across 6675 sites in core promoter regions after filtering for duplicate

64 sites<sup>2</sup>; 1896 of which exhibited allele-specific binding at an unadjusted binomial  $p < 0.05$  (see **Online**

65 **Methods**). We noted that only peaBrain impact score was significantly predictive of allele-specific

66 binding sites (coefficient = 35.38 [12.00, 58.67];  $p = 0.003$ ; see **Table 1 in Main Text**); relaxing the

67 binomial p-value threshold (i.e. increasing the number of sites considered as allele-specific) brings the

68 other non-coding metrics to significance. peaBrain's discriminative ability to identify allele-specific

69 binding sites is consistent with our earlier observation that explicitly adding TFBS annotations did not

70 improve the model. Notably, peaBrain's ability indicates that average expression of all genes in a single

71 tissue and the reference genome is sufficient to learn both TFBS and allele-specific binding.

72

73 To further investigate peaBrain's ability to identify allele-specific binding sites, we compared peaBrain

74 impact scores to predictions by methods specifically designed to predict TFBS, including two neural-

75 network methods (DeepBind<sup>3</sup> and DeepSEA<sup>1</sup>), two kmer-based variant scoring methods (gkmSVM<sup>4</sup>

76 and GERV<sup>5</sup>), and three position-weighted matrices (PWM)-related methods<sup>2</sup>. These methods depend

77 on modelling TF ChIP-seq data in various ways and may have multiple models for the same TF. After

78 confirming the predictive ability of these methods to identify allele-specific binding sites, we noted that

79 peaBrain scores positively correlated only with GERV measures, a kmer-based variant scoring

80 algorithm (**Figure 2**). Unlike the other methods, peaBrain (and GERV) do not assume the existence of

81 canonical motifs and learn TFBS by modelling sequences (or kmers) directly (i.e. not simply by

82 modelling the absence or presence of a ChIP-seq peak). In contrast, for both DeepBind and DeepSEA,

83 we noted positive correlation with at least one PWM-method. These methods generally assume the

84 existence of canonical TF binding sites and predictions are based on the extent of perturbation of those  
85 motifs. While this comparison is limited to variants for which data was available, the peaBrain results  
86 suggest that explicitly characterizing TF motifs is not necessary to understand the consequences of  
87 sequence variation on TF binding and transcriptional dysregulation.

88

89 **Neural activations of penultimate layer of peaBrain model can be used to construct an embedding**

90 **from the genes that encodes correlation information.** Having demonstrated the predicative ability

91 of the peaBrain model (see **Main text** for details), we were subsequently interested in using the

92 activations from the penultimate layer of the model as a continuous (and compressed) representation of

93 the genes. These neural activations capture both the annotated DNA (input) and its additive

94 contributions to tissue- and phenotype-specific abundance (output) in a compressed form amenable to

95 downstream analyses. Furthermore, as these vectors were obtained from a regression model, they

96 readily capture only the salient portions of DNA abundance encoded in the annotated-genome (the

97 weights of the model corresponding to the transcription factor that regulate and interact with this

98 genome). Because of model choice, the mean abundance of each gene was encoded as a linear

99 combination of the vector elements, *i.e.* the output of the regression model. As with our earlier analyses,

100 for brevity, we limited our analysis of the properties of the embeddings to class B models for skeletal

101 muscle. We observed, for the skeletal muscle embeddings, that pairwise cosine similarity between these

102 dense gene representations corresponded to the measured RNAseq correlation between the gene pair.

103 After excluding self-correlations and weakly correlated genes (RNAseq  $\rho < 0.5$ ), we noted that the

104 cosine similarity of the embedding was significantly correlated (Spearman's  $\rho = 0.18$ ;  $p < 2.2 \times 10^{-16}$ )

105 to the experimentally RNAseq-derived correlation. This suggests the annotated-DNA model, without

106 supervision, imposes a linear structure on this vector space: the angle between the vectors corresponds

107 to the co-regulation of the gene pair.

108

109 **peaBrain-derived gene embeddings also encode membership to pathways and other curated gene**

110 **sets.** We were interested in further exploring the utility of these embeddings in other applications. We

111 noted that this dense representation from the class B skeletal muscle pea Brain model encodes

112 membership to the MSigDB Hallmark curated gene sets (average 10-fold cross-validated for all  
113 pathways AUC = ~0.70, **Table 1**), suggesting that the representations themselves, not only encode  
114 abundance and regulatory information, but also functional relationships. (We filtered pathway sets not  
115 relevant to the tissues, such as “PANCREAS\_BETA\_CELLS”, “COMPLEMENT”, or  
116 “SPERMATOGENESIS”). Taken all together, this suggests the gene embeddings capture both the  
117 annotated DNA (input) and its additive contributions to tissue-specific abundance (output) in a  
118 compressed form amenable to downstream analyses (e.g. network-based analyses).

119 **Supplementary Note 1 – Table 1.** Tabulated 10-fold cross-validated AUC for genomewide pathway  
 120 membership predictions using class B MuscleSkeletal Embeddings.

Hallmark Gene Set	10-fold cross-validated average auc
MYC_TARGETS_V1	0.80
MYC_TARGETS_V2	0.79
G2M_CHECKPOINT	0.77
UNFOLDED_PROTEIN_RESPONSE	0.76
OXIDATIVE_PHOSPHORYLATION	0.76
MTORC1_SIGNALING	0.76
EPITHELIAL_MESENCHYMAL_TRANSITION	0.74
MITOTIC_SPINDLE	0.74
E2F_TARGETS	0.73
REACTIVE_OXIGEN_SPECIES_PATHWAY	0.73
TNFA_SIGNALING_VIA_NFKB	0.72
PROTEIN_SECRETION	0.72
TGF_BETA_SIGNALING	0.71
UV_RESPONSE_DN	0.71
PI3K_AKT_MTOR_SIGNALING	0.71
DNA_REPAIR	0.71
HYPOXIA	0.70
P53_PATHWAY	0.69
APOPTOSIS	0.68
APICAL_JUNCTION	0.68
ADIPOGENESIS	0.67
MYOGENESIS	0.66
IL2_STAT5_SIGNALING	0.66
ANGIOGENESIS	0.66
GLYCOLYSIS	0.65
PANCREAS_BETA_CELLS	0.65
ANDROGEN_RESPONSE	0.65
KRAS_SIGNALING_DN	0.64
HEME_METABOLISM	0.64
CHOLESTEROL_HOMEOSTASIS	0.64
HEDGEHOG_SIGNALING	0.63
APICAL_SURFACE	0.63
UV_RESPONSE_UP	0.63
INTERFERON_GAMMA_RESPONSE	0.63
ESTROGEN_RESPONSE_EARLY	0.62
INTERFERON_ALPHA_RESPONSE	0.62
ESTROGEN_RESPONSE_LATE	0.61
NOTCH_SIGNALING	0.61
KRAS_SIGNALING_UP	0.61
INFLAMMATORY_RESPONSE	0.61
COAGULATION	0.60
SPERMATOGENESIS	0.60
WNT_BETA_CATENIN_SIGNALING	0.58

PEROXISOME	0.58
IL6_JAK_STAT3_SIGNALING	0.57
ALLOGRAFT_REJECTION	0.57
COMPLEMENT	0.57
BILE_ACID_METABOLISM	0.56
XENOBIOTIC_METABOLISM	0.56
FATTY_ACID_METABOLISM	0.56

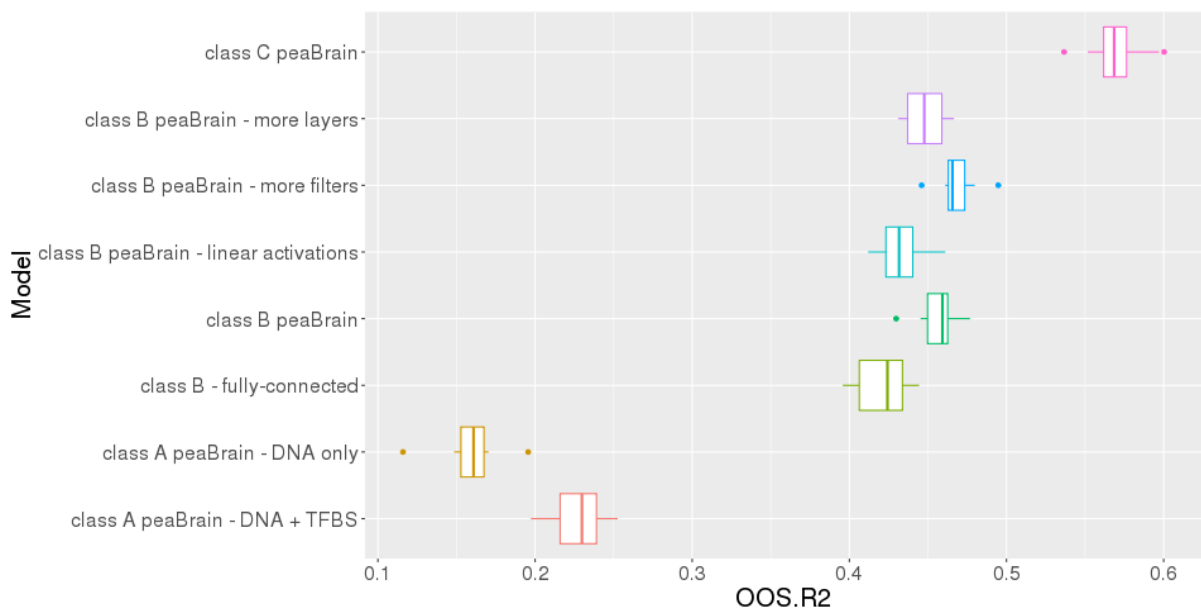
---

121

122

123 **Supplementary Note 1 – Figure 1.** Boxplots of 10-fold cross-validated oos-r<sup>2</sup>, as assessed in skeletal  
124 muscle, for class A models (labelled as “class A peaBrain – DNA only”), class A with TFBS annotations  
125 (labelled as “class A peaBrain – DNA+TFBS”), class B models with tissue-agnostic annotations (“class  
126 B peaBrain – CNNs”), fully connected neural networks (“class B – fully-connected”), class B models  
127 with linear activation functions (“class B peaBrain – linear activations”), class B models with increased  
128 number of layers (“class B peaBrain – more layers”), class B models with increased number of filters  
129 (“class B peaBrain – more filters”), and class C models with tissue-specific annotations (“class C  
130 peaBrain”).

131

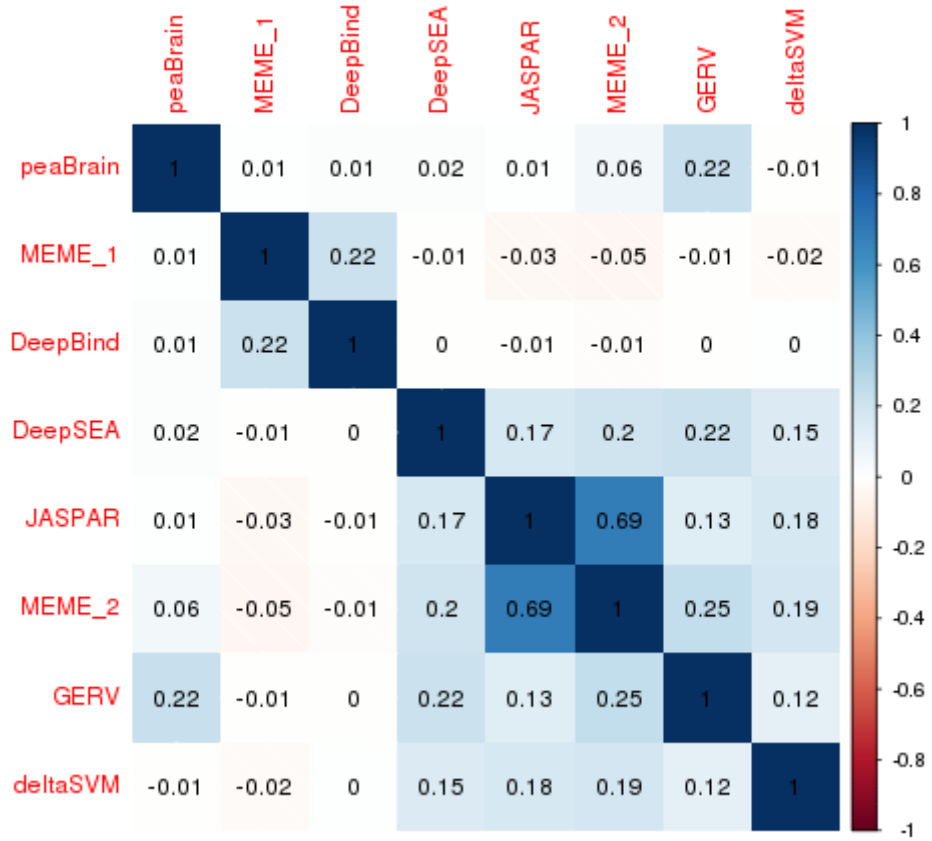


132

133



134 **Supplementary Note 1 – Figure 2.** Rank correlation plot for TF-binding algorithms and the peaBrain  
 135 impact score. JASPAR, MEME\_1 and MEME\_2 are PWM-approaches.



136  
 137  
 138  
 139  
 140  
 141  
 142  
 143  
 144  
 145  
 146  
 147

148 **REFERENCES**

149

150 1 Zhou, J. & Troyanskaya, O. G. Predicting effects of noncoding variants with deep learning–  
151 based sequence model. *Nature methods* **12**, 931 (2015).

152 2 Wagih, O., Merico, D., DeLong, A. & Frey, B. J. Allele-specific transcription factor binding as  
153 a benchmark for assessing variant impact predictors. *bioRxiv*, 253427 (2018).

154 3 Alipanahi, B., DeLong, A., Weirauch, M. T. & Frey, B. J. Predicting the sequence specificities  
155 of DNA-and RNA-binding proteins by deep learning. *Nature biotechnology* **33**, 831 (2015).

156 4 Lee, D. *et al.* A method to predict the impact of regulatory variants from DNA sequence. *Nature*  
157 *genetics* **47**, 955 (2015).

158 5 Zeng, H., Hashimoto, T., Kang, D. D. & Gifford, D. K. GERV: a statistical method for  
159 generative evaluation of regulatory variants for transcription factor binding. *Bioinformatics* **32**,  
160 490-496 (2015).

161