

Cross-species regulatory sequence activity prediction

David R. Kelley
Calico Life Sciences
South San Francisco, CA
drk@calicolabs.com

June 4, 2019

1 Abstract

2 Machine learning algorithms trained to predict the regulatory activity of nucleic acid sequences have revealed
3 principles of gene regulation and guided genetic variation analysis. While the human genome has been
4 extensively annotated and studied, model organisms have been less explored. Model organism genomes offer
5 both additional training sequences and unique annotations describing tissue and cell states unavailable in
6 humans. Here, we develop a strategy to train deep convolutional neural networks simultaneously on multiple
7 genomes and apply it to learn sequence predictors for large compendia of human and mouse data. Training on
8 both genomes improves gene expression prediction accuracy on held out sequences. We further demonstrate
9 a novel and powerful transfer learning approach to use mouse regulatory models to analyze human genetic
10 variants associated with molecular phenotypes and disease. Together these techniques unleash thousands of
11 non-human epigenetic and transcriptional profiles toward more effective investigation of how gene regulation
12 affects human disease.

13 Introduction

14 Predicting the behavior of any nucleic acid sequence in any environment is a primary objective of gene
15 regulation research. In recent years, machine learning approaches to directly tackle this problem have
16 achieved significant accuracy gains predicting transcription factor (TF) binding, chromatin features, and
17 gene expression from input DNA sequence (1–6). These models have then been fruitfully applied to study
18 genetic variation in populations and generate mechanistic hypotheses for how noncoding variants associated
19 with human disease exert their influence (3, 4, 7). Estimates for how mutations influence regulatory activity
20 have also revealed insights into regulatory evolution and the robustness of genes to such mutations (6).

21 The human genome’s ~3 billion nucleotides provide ample training data for highly expressive deep convo-
22 lutional neural networks, which have achieved state of the art performance for many regulatory sequence
23 activity prediction tasks (1, 3–6). The complexity of mammalian gene regulation and these models impres-
24 sive but imperfect predictions suggest room for improvement remains. In particular, distal regulation by
25 enhancers is incompletely captured by existing models. Obtaining more training data is a reliable strategy
26 to improve model accuracy. The research field continues to generate new functional genomics profiles, but
27 these merely deliver additional labels for the existing sequence data; fitting more expressive and accurate
28 models would benefit more from entirely new training sequences. Individual human genomes differ only
29 slightly from each other, so acquiring functional profiles for more humans is unlikely to provide this boost.
30 Artificially designed sequences can offer more data for specific tasks, but only short sequences can be ef-
31 fectively manipulated and their profiling is limited to cell lines that cannot represent the full complexity of
32 human tissues (8–12).

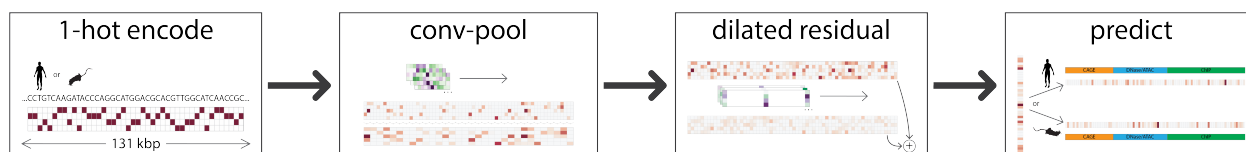


Figure 1: Predicting regulatory sequence activity for human and mouse genomes. We predict the regulatory activity of DNA sequences for multiple genomes in several stages (Methods). The model takes in 131,072 bp DNA sequences, encoded as a binary matrix of four rows representing the four nucleotides. We transform this representation with seven iterated blocks of convolution and max pooling adjacent positions to summarize the sequence information in 128 bp windows. To share information across the long sequence, we apply eleven dilated residual blocks, consisting of a dilated convolution with exponentially increasing dilation rate followed by addition back into the input representation. Finally, we apply a linear transform to predict thousands of regulatory activity signal tracks for either human or mouse. All parameters are shared across species except for the final layer.

33 Non-human species offer a potential source of this desired additional training data. Regulatory sequence
34 evolves rapidly, but TF binding preferences are highly conserved due to the drastic effect that modifying
35 affinity for many thousands of binding sites would confer on the organism (13–15). Thus, we hypothesized
36 that regulatory programs across related species have enough in common to benefit machine learning sequence
37 activity. To demonstrate the concept, we chose the mouse as a distant mammal with substantial functional
38 genomics data available (16). In addition to serving as a source of more genomic sequence, mouse experiments
39 can explore biological states that are challenging or unethical to acquire in humans, e.g. profiling mouse
40 development, disease, and genome modifications. If context-specific regulatory programs are sufficiently
41 conserved across species, then models learned to predict these data in the mouse may be applicable to
42 impute human genome profiles to study human regulatory sequences and genetic variation.

43 In this work, we trained a deep convolutional neural network to jointly learn the complex regulatory pro-
44 grams that determine TF binding, DNA accessibility, and transcription using the ENCODE and FANTOM
45 compendia of thousands of functional genomics profiles from hundreds of human and mouse cell types. We
46 introduce a novel model architecture that better captures long range interactions by applying residual con-
47 nections between layers. We benchmarked single versus joint training and found that jointly training on
48 human and mouse data leads to more accurate models for both species, particularly for predicting CAGE
49 RNA abundance. We demonstrated that mouse regulatory programs can be transferred across species to
50 human where they continue to make accurate tissue-specific predictions. Applying this procedure to predict
51 human genetic variant effects revealed significant correspondence with eQTL statistics and proved insightful
52 for studying human disease.

53 Results

54 Multi-genome training improves gene expression prediction accuracy

55
56 We applied the Basenji software and framework to predict functional genomics signal tracks from only DNA
57 sequence (4). The neural network takes as input a 131,072(= 2^{17}) bp sequence, transforms its representation
58 with iterated convolution layers, and makes predictions in 128 bp windows across the sequence for the
59 normalized signal derived from many datasets (Figure 1, Methods). We introduced a novel architecture
60 that uses residual connections to alleviate the strain of vanishing gradients in deep network optimization
61 to improve generalization accuracy (Supplementary Figure 1) (17). Training on multiple genomes required
62 several further developments (Methods). Most importantly, we modified the train/valid/test split of the
63 genomic sequences to ensure that homologous regions from different genomes did not cross splits (Methods);
64 without this extra care, we might overestimate generalization accuracy.

65 We assembled training data consisting of 6,956 human and mouse quantitative sequencing assay signal tracks
66 from the ENCODE and FANTOM consortiums (Methods). These data describe regulatory activity across
67 tissues and isolated cell types using several techniques—DNase and ATAC-seq to measure DNA accessibility,
68 which typically mark TF-bound sites, and ChIP-seq to map TF binding sites and histone modification
69 presence (18, 19). The FANTOM data consists of RNA abundance profiling with CAGE, where the 5' end of
70 the transcript is sequenced (20). These 5' RNA profiles are independent of splicing and allow us to provide
71 DNA sequence without gene annotations, which would not be the case for RNA-seq (4). In addition, we
72 added several mouse datasets describing cell states that are unavailable for humans: (1) a single cell ATAC-
73 seq atlas from 13 tissues clustered to 78 distinct profiles (21) and (2) several TF and chromatin profiles
74 obtained over 24 hour time courses in the liver to study circadian rhythms (Supplementary Table 1).

75 To measure the influence of multi-genome training on generalization accuracy, we trained three separate
76 models on these data: one jointly fit to both human and mouse, one to human data alone, and one to
77 mouse data alone. For each scenario, we fit the same model architecture and hyperparameters. We allowed
78 each model to train until 30 epochs had passed without improvement on the validation set, which provides
79 considerable slack to ensure that each model has reached its full potential.

80 The joint training procedure improved test set accuracy for 94% of human CAGE and 98% of mouse CAGE
81 datasets (binomial test p-values $1e-16$ and $1e-16$), increasing the average Pearson correlation by .013 and
82 .026 for human and mouse respectively (Figure 2a,c). For DNase, ATAC, and ChIP, joint training improved
83 predictions by a lesser margin relative to single genome training; average test set correlation increased for
84 55% of human and 96% of mouse datasets (binomial test p-values $3e-11$ and $1e-16$) (Figure 2b,d).
85 Datasets where single genome accuracy exceeded joint did not show any interesting pattern and are likely
86 just attributable to noise from the stochastic training procedure. CAGE has several properties that may
87 explain the observed extra benefit of having more training data from multiple genomes. CAGE signal
88 has a larger dynamic range than the other data, spanning orders of magnitude, fewer relevant sites in the
89 genome, and more sophisticated transcriptional regulatory mechanisms that often involve distant sequences.
90 Altogether, these results demonstrate that regulatory programs are sufficiently similar across the 90 million
91 years of independent evolution separating human and mouse so that their annotated genomic sequences
92 provide informative multi-task training data for building predictive models for both species.

93 **Regulatory sequence activity models transfer across species**

94

95 Regulatory program conservation across related species has been observed in genome-wide functional profiles
96 of TF binding and histone modifications (13–15). In matched tissue samples, similar TFs are typically present
97 and those TFs have highly conserved motif preferences (15, 22). These findings suggest that a regulatory
98 sequence activity model trained to predict for one species will also make usefully accurate predictions for
99 matched samples from the other. To quantify this proposition, we selected several diverse and representative
100 tissues and cell types for which we could unambiguously match across species—cerebellum, liver, and CD4+
101 T cells. We extracted CAGE gene expression measurements from the transcription start sites (TSS) for
102 all human genes outside the training set and computed predictions for human and mouse versions of these
103 tissues and cell types (Figure 3a). For this exercise, and those to follow, we used the jointly trained multi-task
104 model and sliced out predictions of interest.

105 Across human gene TSSs, we observed high cross-species prediction accuracy of 0.73 Pearson correlation for
106 mouse predictions to human observed signal averaged across these samples, relative to 0.75 correlation for
107 human predictions to human observed signal. To assess whether the model further captures and transfers
108 tissue specificity, we normalized each TSSs data or predictions by its mean across all CAGE datasets. Mean
109 normalization removes correlation driven by accurate prediction of global cross-tissue activity. Pearson
110 correlation for normalized signal remained high for mouse predictions to human data for the matched samples
111 (mean 0.40, Figure 3b,c). In contrast, normalized predictions compared to data from distinct tissues/cell
112 types resulted in negative correlations (Figure 3c). Thus, the models have learned tissue and cell type
113 specificity beyond a baseline level and are able to transfer that knowledge across species.

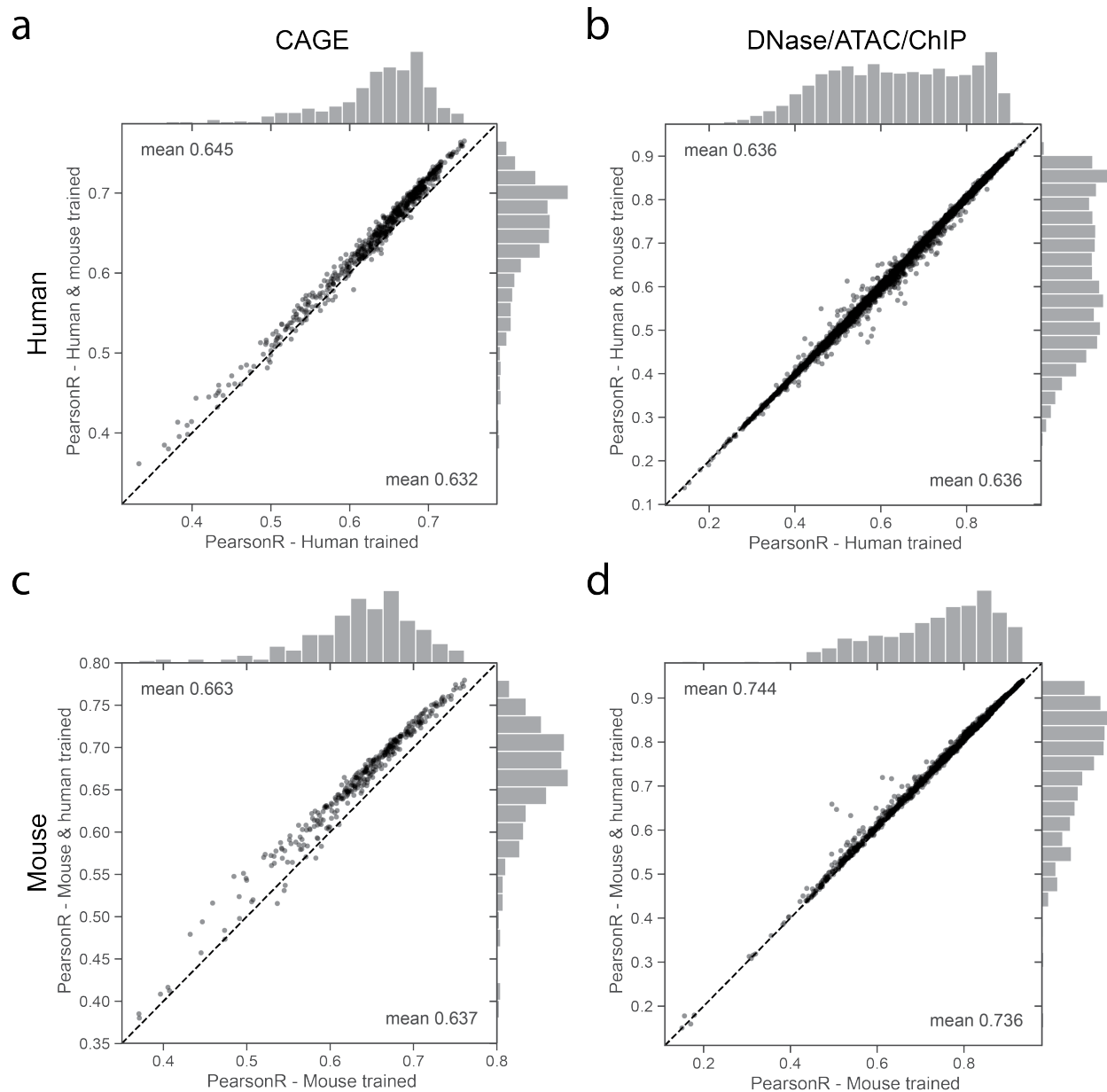
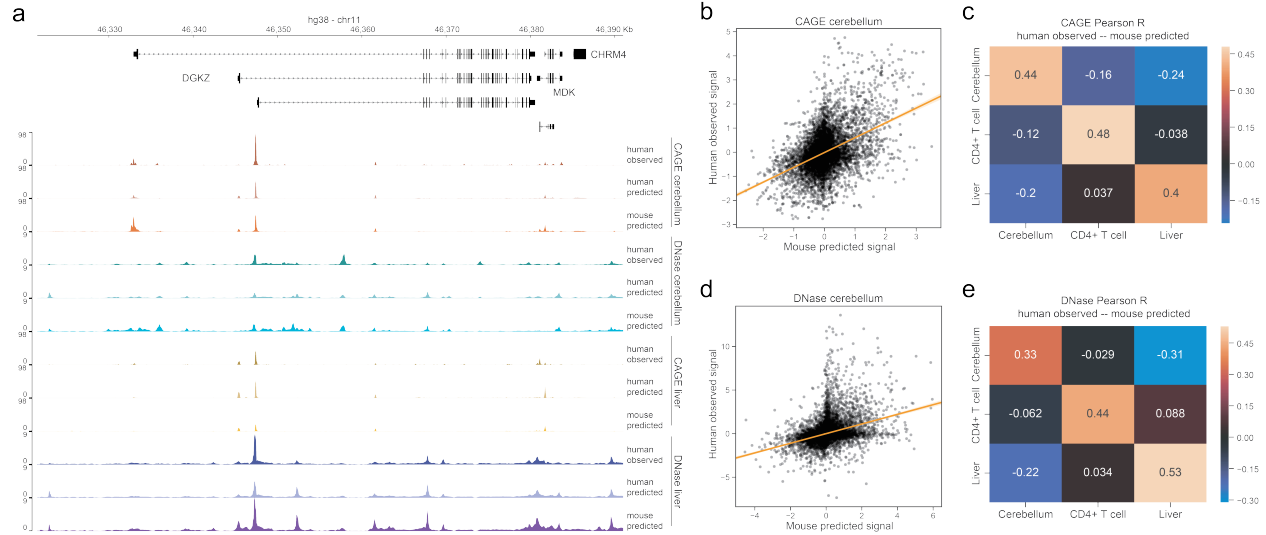


Figure 2: Training on human and mouse data improves generalization accuracy. We trained three separate models with the same architecture on human data alone, mouse data alone, and both human and mouse data jointly. For each model, we computed the Pearson correlation of test set predictions and observed experimental data for thousands of datasets from various experiment types. Points in the scatter plots represent individual datasets, with single genome training accuracy on the x-axis and joint training accuracy on the y-axis. For CAGE, training on multiple genomes increases test set accuracy on nearly all datasets for both human and mouse. For DNase/ATAC/ChIP-seq, test set accuracy improves by a smaller average margin.



114 We repeated these analyses with DNase accessibility profiles for the same tissues and cell types to assess how
 115 general this transferability is for different data. Because most sites lack activity, we selected the top 10%
 116 most variable. We observed the same statistical trends for accessibility—high correlation between mouse
 117 predictions and human data for matched samples (mean 0.84) and specificity for scaled comparisons (Figure
 118 3d,e).

119 Mouse-trained models elucidate human genetic variant effects

120

121 A driving goal of regulatory sequence modeling is to predict the effect of human genetic variants on gene
 122 expression and downstream phenotypes. For any biallelic variant, we can predict signal across the surround-
 123 ing genomic sequence for each allele and derive a summary score for the variant effect (Figure 4a). Here, we
 124 sum the signal across the sequence and take the difference between alleles. We can compute this score for
 125 every dataset using two forward passes of the convolutional neural network.

126 Models trained on mouse data allow one to predict the difference between how two human alleles would
 127 behave if they were present in the regulatory environment of mouse cells. Given the evidence that analogous
 128 human and mouse cells largely share regulatory programs, we hypothesized that models trained on mouse
 129 data would be insightful towards understanding human regulatory variants function. To test this hypothesis,
 130 we studied the Gene-Tissue Expression (GTEx) release v7a data of genotypes and gene expression profiles
 131 for hundreds of humans across dozens of tissues (23). In previous work, we showed that variant scores derived

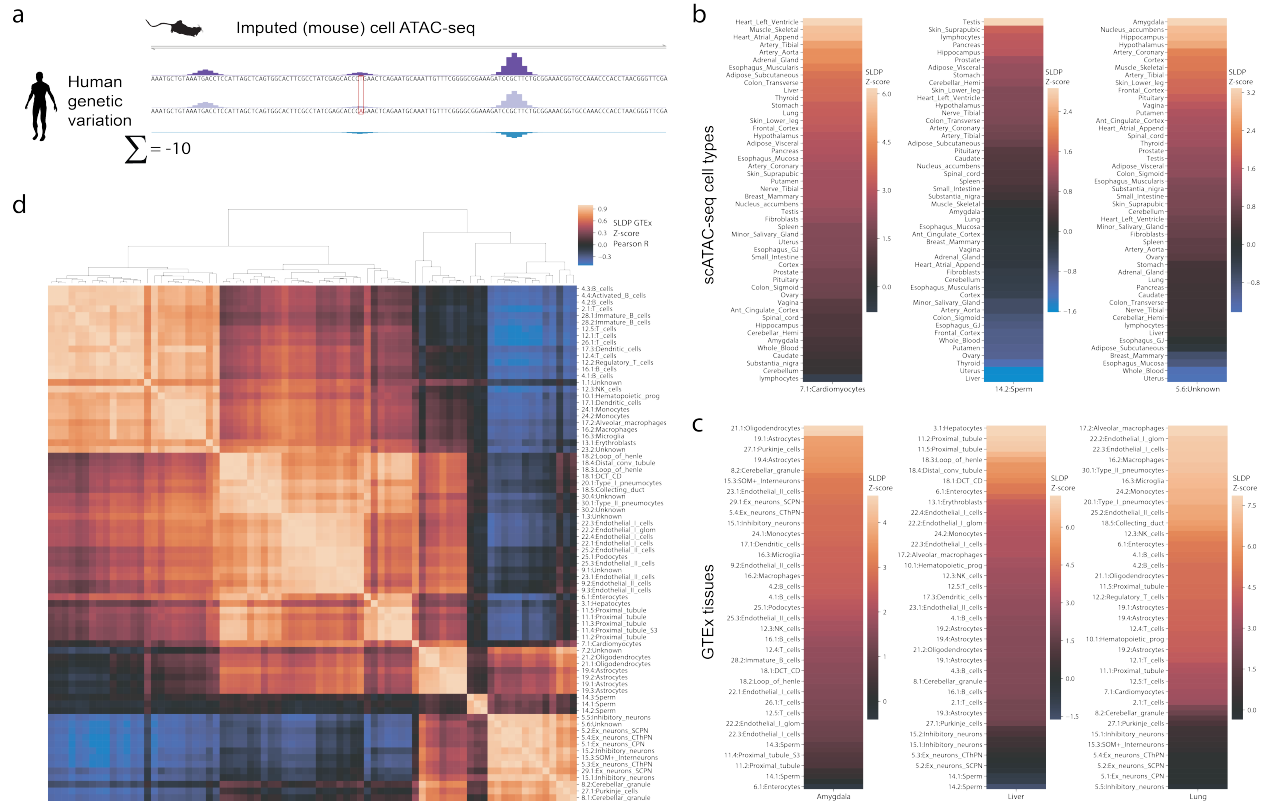


Figure 4: Mouse cell type accessibility predictions show a strong and specific statistical relationship with human eQTLs. (a) We predicted the effect of human genetic variants on imputed regulatory signal trained on mouse single cell ATAC-seq (scATAC) cluster profiles. We scored variants by subtracting the signal from the minor allele from that of the major and summing across the sequence. (b) We used signed linkage disequilibrium profile (SLDP) regression to compare the cell type-specific variant effect predictions to tissue-specific eQTL summary statistics from GTEx. Cell type profiles correspond best with the expected tissues. (c) GTEx tissues correspond best with the expected cell types. (d) Clustering scATAC cell types by their z-scores across GTEx tissues reveals the expected structure.

132 from Basenji predictions corresponded significantly with GTEx summary statistics (4). Here, we conducted
 133 a similar analysis using signed linkage disequilibrium profile (SLDP) regression to measure the statistical
 134 concordance between signed variant effect predictions and GTEx summary statistics (Methods) (7). SLDP
 135 distributes a signed annotation (i.e. our scores) according to a given population's LD structure and compares
 136 it to a set of summary statistics. Using a permutation scheme, the method produces a signed Z-score that
 137 specifies the direction and magnitude of the relationship and a p-value describing its significance.

138 We focused on a dataset unique to the mouse—a single cell ATAC-seq atlas from 13 adult mouse tissues,
 139 from which 85 distinct cell type patterns were identified (21). We sliced predictions for these datasets from
 140 the model trained jointly on all human and mouse data. We first asked whether coverage tracks derived from
 141 clustering single cell assays are amenable to Basenji modeling. Predictions for held out sequences achieved
 142 Pearson correlation ranging from 0.43-0.84 in 128 bp windows for these 85 profiles, which is in line with
 143 predictions for bulk DNase/ATAC-seq.

144 Human variant predictions for these models generally exhibited a strong, positive effect on GTEx summary
 145 statistics, in line with prior observations that increased accessibility typically increases gene expression. Fur-
 146 thermore, cell type predictions aligned well with anatomical expectations. For example, variant predictions
 147 for cardiomyocytes have the strongest correlation with GTEx measurements in the heart and skeletal muscle
 148 (Figure 4b). From the opposite direction, GTEx measurements for the liver have the strongest correlation

149 with variant predictions for hepatocytes (Figure 4c). These results further support the claim that human
150 and mouse cells share relevant regulatory factors and that our procedure can project these factors across
151 species from mouse experiments to human variants.

152 For each pair of mouse ATAC cell types, we computed the correlation between their SLDP Z-scores across
153 GTEx tissues (Figure 4d). The correlations revealed expected structure, with clusters representing the
154 blood, endothelial cells, neurons, among others. The original authors abstained from annotating 9 of the 85
155 clusters. Through this procedure, we can suggest high-level annotations for several of the unknown clusters.
156 For example, 5.6 appears similar to various neuron subtypes due to the strong statistical relationship between
157 variant predictions and the GTEx brain tissue summary statistics (Figure 4b,d).

158 **Mouse-trained models highlight mutations relevant to human neurodevelopmental disease**

159

160

161 Having established the relevance and specificity of mouse dataset predictions for expression phenotypes, we
162 asked whether these data could provide insight into the genetic basis of human disease. Mouse data has
163 proven valuable for studying human genetic variants in previous work (16, 21), but these analyses were
164 limited to studying variants in homologous sequences in their mouse genome context. Given the substantial
165 regulatory sequence turnover between these genomes, this limitation is severe. The predictive framework
166 here avoids this limitation by mapping the learned mouse regulatory program to the human genome setting
167 for all variants.

168 To explore the utility of this procedure for studying human disease, we retrieved a recent dataset of 1902
169 quartet families from the Simons Simplex Collection (24) with whole genome sequencing of a mother, father,
170 child affected by autism, and unaffected sibling. In these data, the offspring have an average of 67 de novo
171 mutations, which have a slight enrichment in promoters (25). Recent work demonstrated that variant effect
172 predictions further differentiate autism cases from their unaffected sibling controls (26). We hypothesized
173 that predictions using models trained on mouse data would also distinguish the disease and perhaps provide
174 additional insight via novel developmental profiles.

175 We applied the model to predict how each de novo mutation would influence signal in 357 mouse CAGE
176 profiles of tissues and cell types throughout the body. Mann-Whitney U (MWU) tests revealed significantly
177 more negative predictions in the case versus control variant sets for 246 CAGE profiles at FDR ≤ 0.1 (Figure
178 5a). Appreciating the correlations in these data, we also transformed the variants by predictions matrix with
179 PCA to represent each variant by its first principal component score (which explained 51% of the variance).
180 In principal component space, the MWU test comparing case and control variants was significant with p-
181 value 0.002. Most leading datasets described brain regions and cell types; the 76 brain dataset p-values were
182 less than non-brain data with p-value 1×10^{-10} by MWU test.

183 Highly negative predictions indicate mutations that disturb active regulatory elements. For example, a case
184 variant upstream of ZNF644 modifies a critical nucleotide in a consensus motif for the transcription factor
185 YY1, which the model identifies as active and relevant (Figure 5b). ZNF644 has considerable evidence for
186 intolerance to loss of function mutations in the Genome Aggregation Database v2.1.1 (gnomAD) with prob-
187 ably 0.999 of intolerance (27). YY1 has been implicated in processes that determine the three-dimensional
188 positioning of promoters and enhancers (28). Thus, we hypothesize that the variant modifies the enhancer
189 regulation of this critical protein.

190 Perhaps unexpectedly, 15 datasets describing the developing heart also emerged from this analysis (Figure
191 5a). This result is supported by whole genome sequencing of congenital heart disease probands, which has
192 revealed affected gene sets that overlap significantly with those observed in neurodevelopmental sequencing
193 efforts like this one (29, 30). In addition to the brain and heart, whole body profiles from the embryo and
194 neonate stages also have p-values among the lowest.

195 This significant enrichment indicates that variant effect predictions may help classify disease at the individual

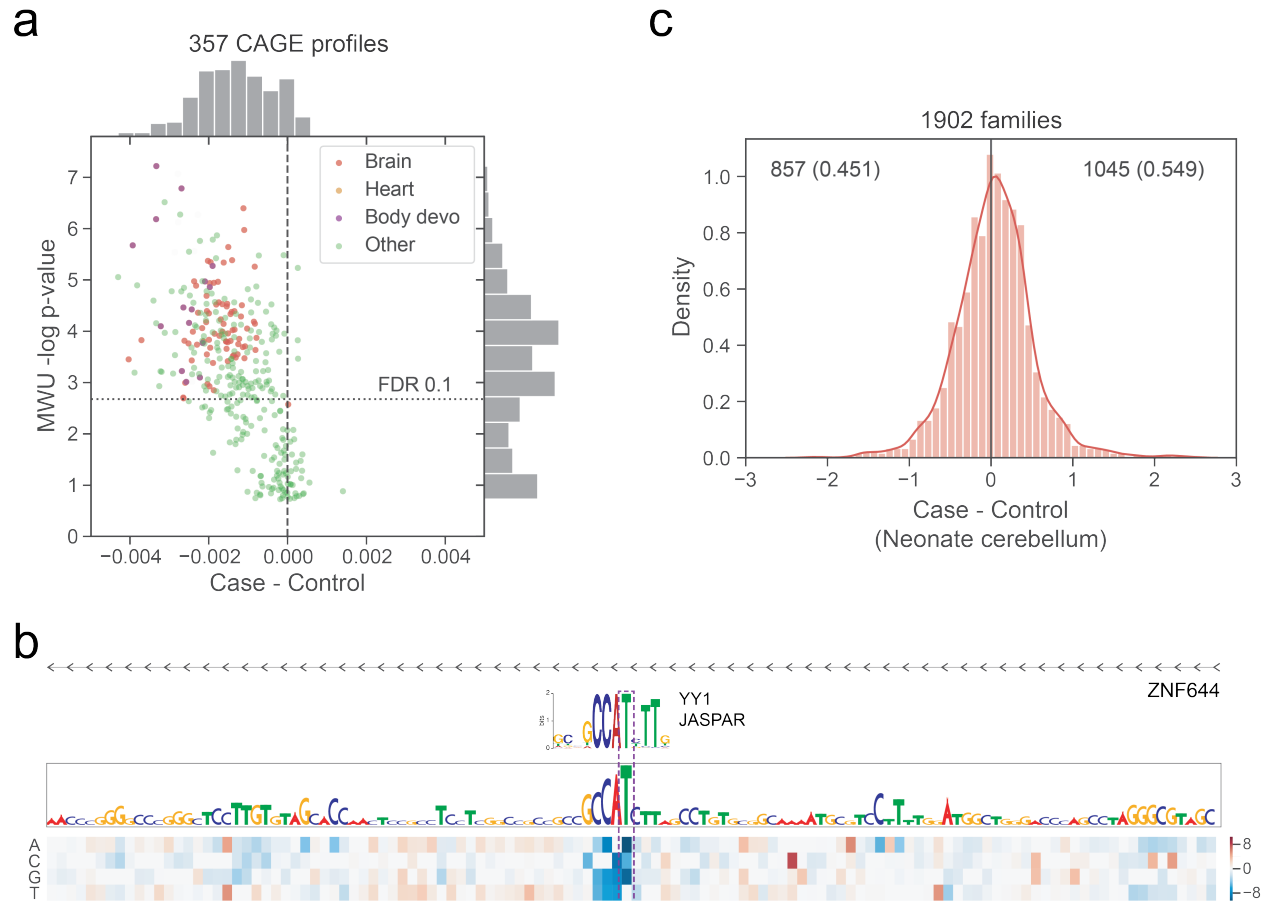


Figure 5: Human de novo variant predictions for mouse data enrich for autism cases versus controls. (a) We predicted the influence of 234k de novo variants split between cases and controls on 357 CAGE datasets in mouse. For each dataset, we computed a Mann-Whitney U (MWU) test between case and control sets and corrected for multiple hypotheses using the Benjamini-Hochberg procedure. Predictions for many datasets were enriched for more negative values in the cases, driven largely by brain, heart, and whole body developmental profiles. Each datasets x-axis position is the mean inverse hyperbolic sine over case variants minus the equivalent over control variants. The mean inverse hyperbolic sine transformation is similar to logarithm, but gracefully performs the symmetric transformation for negative values. (b) A case variant at chr1:91021795 modifies a critical T in a YY1 motif to an A in the promoter region of ZNF644. (c) At the individual level, a simple score summing all negative predictions for the leading dataset describing neonate cerebellum significantly separates cases from their matched controls. The x-axis position represents the log ratio between case and control sums.

196 level. For each individual, we computed a simple risk score by summing the negative predictions in the
 197 neonate cerebellum dataset. This score suggests more deleterious de novo variants for 54.9% of the cases
 198 versus their controls (binomial test p-value 9×10^{-6}) (Figure 5c). Thus, this approach is a strong candidate
 199 for inclusion with complementary feature sources from coding mutations and structural variation to continue
 200 to characterize this incompletely understood disorder.

201 Discussion

202 In this work, we developed a novel convolutional neural network architecture and multiple species training
203 procedure to enable one model to train on 6956 functional genomics signal tracks annotating the human
204 or mouse genomes. We observed that training jointly on both species produced models that make more
205 accurate predictions on unseen test sequences relative to models trained on a single species. Regulatory
206 sequence activity predictions for human sequences in mouse tissues correlate well with datasets describing
207 the corresponding human tissues. Model predictions for altered regulatory activity of human genetic vari-
208 ants made with respect to mouse datasets have a strong statistical concordance with tissue-specific human
209 eQTL measurements. Mouse machine learning models can be used to study human disease, exemplified by
210 enrichment of deleterious predictions among de novo autism variants relative to control sets.

211 We focused here on human and mouse because both species have been comprehensively studied with genome-
212 wide functional genomics. Our observation that joint training on these two genomes improves prediction
213 accuracy opens the possibility of more complex schemes for training on larger numbers of genomes. Given
214 the substantial evolutionary distance between human and mouse, regulatory annotations for all mammalian
215 genomes are likely to provide similarly useful training data. Primate genomes will be particularly interesting
216 to explore; their tissues and cell types will more closely match those of human, but their sequences are far
217 more similar. Prediction accuracy improved more for CAGE gene expression measurements than accessibility
218 or ChIP-seq, which suggests that the number of events and their regulatory complexity are relevant features
219 for determining whether multiple genome training will be worthwhile. Efforts to predict spatial contacts
220 between chromosomes as mapped by Hi-C and its relatives likely fit this criteria, and we hypothesize that
221 training sequence-based models on human and mouse data together will be fruitful.

222 Much prior work has revealed the similarity of regulatory programs across species, but transferring knowledge
223 gleaned from an accessible model organism (such as mouse) to another of interest (such as human) has
224 remained challenging. Existing approaches rely on whole genome alignments to transfer annotations from
225 one genome to the other (21, 31). These approaches are constrained by the quality of the alignment, which
226 is a notoriously challenging bioinformatics problem (32), and the limited proportion of each genome that
227 aligns (40% for human and 45% for mouse). Here, we demonstrated an alternative approach where a machine
228 learning model trained on the model organism data compresses the relevant knowledge into its parameters,
229 which can then be applied to make predictions for sequences from the genome of interest. Substantial
230 research in transfer learning with neural networks for natural language processing motivates and supports
231 the viability of this procedure (e.g. (33)). The strong tissue-specific statistical relationship between human
232 genetic variant predictions from model parameters trained to predict mouse annotations and GTEx tissue-
233 specific eQTLs highlights the successful nucleotide resolution of our mouse to human transfer learning. The
234 Gene Expression Omnibus (GEO) contains tens of thousands of mouse functional genomics profiles, many
235 describing experiments impossible in humans. For example, we included dozens of datasets describing mouse
236 liver profiles over 24 hour time courses to study the circadian rhythms of gene expression and chromatin.
237 Models trained to predict all datasets, as well as open source software to compute these predictions and
238 train new models on users own data, are available in the Basenji software package (34).

239 Methods

240 Functional genomics data

241
242 In this work, we studied quantitative sequencing assays performed on human and mouse samples. Specifically,
243 we focused on DNase and ATAC-seq profiling DNA accessibility, ChIP-seq profiling TF binding or histone
244 modifications, and CAGE profiling RNA abundance derived from 5 transcription start sites. Preprocessing
245 these data effectively is critical to successful machine learning. Our primary preprocessing objective is to
246 denoise these data to the relevant signal at nucleotide-resolution.

247 We largely followed the preprocessing pipeline described in prior research introducing the Basenji framework
248 (4). The standard pipeline through which experimental data flowed follows:

- 249 1. Trim raw sequencing reads using fastp, which can automatically detect and remove unwanted adapter
250 nucleotides (35).
- 251 2. Align reads using BWA to hg38 or mm10 and requesting 16 multi-mapping read positions (36).
- 252 3. Estimate nucleotide-resolution signal using an open source script from the Basenji software that dis-
253 tributes multi-mapping reads, normalizes for GC bias, and smooths across positions using a Gaussian
254 filter (4).

255 However, we varied from this standard pipeline for all data available from the ENCODE consortium website,
256 which is 4506 human and 1019 mouse experiments. These data have been thoughtfully processed using open
257 source pipelines and are available for download at several stages, including log fold change signal tracks in
258 BigWig format (37). Rather than reprocess these data without full knowledge of how replicate and control
259 experiments match, we chose to use these signal tracks directly. The Seattle Organismal Molecular Atlas
260 (SOMA) server provides a single cell mouse ATAC-seq atlas (21). These data are also available in log fold
261 change BigWig format, and we similarly chose to use these rather than reprocess the single cell data. We
262 clipped negative values in all such BigWig tracks to zero.

263 We applied several transformations to these tracks to protect the training procedure from large incorrect
264 values. First, we collected blacklist regions from ENCODE and added all RepeatMasker satellite and sim-
265 ple repeats (38), which we found to frequently collect large false positive signal (39). We further defined
266 unmappable regions of >32 bp where 24-mers align to >10 genomic sites using Umap mappability tracks
267 (40). We set signal values overlapping these regions to the 25th percentile value of that dataset. Finally, we
268 soft clipped high values with the function $f(x) = \min(x, t_c + \sqrt{\max(0, x - t_c)})$. Above the threshold
269 t_c (chosen separately for each experiment and source), this function includes only the square root of the
270 residual $x - t_c$ rather than the full difference.

271 When replicate experiments profiling the same or related samples were available, we averaged the signal
272 tracks. Altogether, the training data includes 638 CAGE, 684 DNase/ATAC, and 3991 ChIP datasets in
273 human and 357 CAGE, 228 DNase/ATAC, and 1058 ChIP datasets in mouse. Supplementary Table 1
274 describes all data with preprocessing parameters.

275 Model architecture

276

277 We modeled genomic regulatory sequence activity signal as a function of solely DNA sequence using a convo-
278 lutional neural network. Such deep learning architectures have excelled for many similar tasks (1, 3–5). We
279 follow our prior work in analyzing large 131 kbp sequences in order to consider long range interactions.

280 The first stage of the architecture aims to extract the relevant sequence motifs from the DNA sequence using
281 the following block of operations:

- 282 1. Convolution width 5 (or 15 in first layer)
- 283 2. Batch normalization
- 284 3. Gaussian Error Linear Unit (GELU) activation
- 285 4. Max pool width 2

286 We applied this block seven times so that each sequence position represents 128 bp, increasing the number
287 of filters from an initial 288 by 1.1776x each block to 768 filters by the end. The GELU activation slightly
288 outperformed the more common ReLU in our benchmarks (41).

289 The second stage of the architecture aims to spread information across the sequence to model long range
290 interactions. In prior work, we applied densely connected dilated convolutions for this task (4). Here, we

291 applied a related but more effective variation, which we refer to as a dilated residual block. Recent deep
292 learning research has revealed that skip connections between layers where one layers representation is directly
293 added to a subsequent layers representation relieve vanishing gradients and improve gradient descent training
294 (17). Thus, we applied the following series of operations:

- 295 1. GELU activation
- 296 2. Dilated convolution width 3, dilation rate d , 384 filters
- 297 3. Batch normalization
- 298 4. GELU activation
- 299 5. Convolution width 1, back to 768 filters
- 300 6. Batch normalization
- 301 7. Dropout probability 0.3
- 302 8. Addition with the block input representation before step 1.

303 We applied this block eleven times, increasing the dilation rate d by 1.5x each time.

304 In the final stage, we first transformed this 1024x768 (length x filters) representation of 128 bp windows with
305 an additional width 1 convolution block using 1536 filters and dropout probability 0.05. To make predictions
306 for either 5313 human or 1643 mouse datasets, we applied a final width one convolution followed by a
307 softplus activation to make all predictions positive. We attached a genome indicator bit to each sequence to
308 determine which final layer to apply.

309 We trained to minimize a Poisson log likelihood in the center 896 windows, ignoring the far sides where
310 context beyond the sequence is missing. The Poisson model is not technically appropriate for the log fold
311 change tracks. However, by clipping negative values to zero, the distribution of values resembles that from
312 our standard processing. On a subset of data, we observed that using the log fold change track did not
313 decrease accuracy or the utility of the model for genetic variant analysis.

314 We minimized with stochastic gradient descent (SGD) on batches of 4 sequences. We implemented the
315 network in TensorFlow and used automatic differentiation to compute gradients via back propagation (42).
316 We performed several grid searches to choose model and optimization hyper parameters for the following
317 sets: (1) SGD learning rate and momentum; (2) initial convolution filters and convolution filter multiplication
318 rate; (3) dilated convolution filters and dropout rate; (4) final convolution filters and dropout rate.

319 Data augmentation describes a suite of techniques to expand the implicit size of the training dataset from the
320 perspective of model training by applying transformations that preserve annotations to data examples. We
321 tiled the 131,072 bp sequences across the chromosomes by 65,599 bp, representing a 50% overlap minus 63
322 bp in order to also shift the 128 window boundaries and max pooling boundaries. During training, we cycled
323 over combinations of two transformations that maintain the relationship between sequence and regulatory
324 signal while changing the model input: (1) reverse complementing the sequence and reversing the signal; (2)
325 shifting the sequence 1-3 bp left or right. Both transformations improved test accuracy and reduce overfitting
326 in our benchmarks.

327 **Multi-genome training**

328

329 Training on multiple genomes containing orthologous sequence complicates construction of holdout sets.
330 Independently splitting each genomes sequences would allow training on a human promoter and testing on
331 its mouse orthologue. If the model memorized conserved elements of the sequence, rather than learning a
332 general function, we might overestimate generalization accuracy.

333 We used the following procedure to minimize occurrence of this potential issue:

- 334 1. Divide each genome into 1 mbp regions.
- 335 2. Construct a bipartite graph where vertexes represent these regions. Place edges between two regions
336 if they have >100 kbp of aligning sequence in a whole genome alignment.
- 337 3. Find connected components in the bipartite graph.
- 338 4. Partition the connected components into training, validation, and test sets.

339 We used the hg38-mm10 syntenic net format alignment downloaded from the UCSC Genome Browser site
340 (43). Using this procedure, we set aside approximately 12% of each genome into validation and test sets
341 respectively. Stricter parameter settings created a single large connected component that did not allow for
342 setting aside enough validation and test sequences.

343 Another complication of training on multiple genomes arises from imbalance between each genome’s se-
344 quences and datasets. We extracted 38.2k human and 33.5k mouse sequences for analysis. We assembled
345 batches of sequences from one genome or the other, chosen randomly proportional to the number of sequences
346 from each genome. The overall loss function comprises a term for every target dataset summed, which leads
347 to larger step magnitudes for batches of human sequences that are annotated with >3 times more datasets.
348 Explicit weighting could be applied to preference training towards a particular species, but we found this to
349 be unnecessary in our experiments for good mouse performance.

350 Jointly training on both human and mouse data constrains the model slightly more than is ideal. We found
351 that training several epochs on only one genome or the other after the full joint procedure improved validation
352 and test set accuracy.

353 **GTEx SLDP**

354

355 We predicted the effect of a genetic variant on various annotations by computing a forward pass through the
356 convolutional network using the reference and alternative alleles, subtracting their difference, and summing
357 across the sequence to obtain a single signed score for each annotation. We averaged scores computed using
358 the forward and reverse complement sequence and small sequence shifts to the left and right. We computed
359 scores for all 1000 Genomes SNPs, which we provide for download from [available upon publication].

360 Signed linkage disequilibrium profile (SLDP) regression is a technique for measuring the statistical concor-
361 dance between a signed variant annotation v and a genome-wide association study’s marginal correlations
362 between variants and a phenotype $\hat{\alpha}$ (7). The functional correlation between v and the true variant effects
363 on the phenotype describes how relevant the annotation is for the phenotype’s heritability. Our model
364 produces these signed variant annotations, and SLDP offers a validated approach to assessing their rele-
365 vance to human phenotypes. Briefly, the method estimates this functional correlation using a generalized
366 least-squares regression, accounting for the population LD structure. SLDP performs a statistical test for
367 significance by randomly flipping the the signs of entries in v in large consecutive blocks to obtain a null
368 distribution. We follow previous work in conditioning on minor allele frequency and binary annotations for
369 variant overlap with coding sequence (and 500 bp extension), 5’ UTR (and 500 bp extension), 3’ UTR (and
370 500 bp extension), and introns.

371 We downloaded GTEx v7a summary statistics for 48 tissues (23). We summarized each SNP’s effect on all
372 cis-genes using the following transformation suggested for SLDP analysis

$$\hat{\alpha}_m = \frac{1}{\sqrt{|G_m|}} \sum_{k \in G_m} \hat{\alpha}_m^{(k)}$$

373 where G_m is the set of all genes for which a cis-eQTL test was performed for variant m and $\hat{\alpha}_m^{(k)}$ is the
374 marginal correlation of SNP m and gene k expression (7). We passed $\hat{\alpha}_m$ to SLDP for analysis of variant
375 predictions.

376 Simons Simplex Collection

377

378 We downloaded 255,106 de novo variants derived from whole-genome sequencing of 1902 quartet families
379 with an autistic child from the Simons Simplex Collection from the supplement of An et al. (25). We filtered
380 these variants for SNPs and computed predictions as described above.

381 Acknowledgements

382 Jacob Kimmel, Leland Taylor, Geoff Fudenberg, Vikram Agarwal, and Han Yuan for valuable feedback.

383 References

- 384 [1] Babak Alipanahi, Andrew Delong, Matthew T Weirauch, and Brendan J Frey. Predicting the sequence
385 specificities of DNA- and RNA-binding proteins by deep learning. *Nature Biotechnology*, 33:831–838,
386 July 2015.
- 387 [2] Dongwon Lee, David U Gorkin, Maggie Baker, Benjamin J Strober, Alessandro L Asoni, Andrew S
388 McCallion, and Michael A Beer. A method to predict the impact of regulatory variants from DNA
389 sequence. *Nature Genetics*, 47:955–961, June 2015.
- 390 [3] David R Kelley, Jasper Snoek, and John L Rinn. Basset: learning the regulatory code of the accessible
391 genome with deep convolutional neural networks. *Genome research*, 26(7):990–999, July 2016.
- 392 [4] David R Kelley, Yakir A Reshef, Maxwell Bileschi, David Belanger, Cory Y McLean, and Jasper
393 Snoek. Sequential regulatory activity prediction across chromosomes with convolutional neural net-
394 works. *Genome Research*, 28(5):739–750, May 2018.
- 395 [5] Jian Zhou and Olga G Troyanskaya. Predicting effects of noncoding variants with deep learning–based
396 sequence model. *Nature methods*, 12(10):931, 2015.
- 397 [6] Jian Zhou, Chandra L Theesfeld, Kevin Yao, Kathleen M Chen, Aaron K Wong, and Olga G Troyan-
398 skaya. Deep learning sequence-based ab initio prediction of variant effects on expression and disease
399 risk. *Nature Genetics*, 464:1, July 2018.
- 400 [7] Yakir A Reshef, Hilary K Finucane, David R Kelley, Alexander Gusev, Dylan Kotliar, Jacob C Ulirsch,
401 Farhad Hormozdiari, Joseph Nasser, Luke OConnor, Bryce Van De Geijn, et al. Detecting genome-wide
402 directional effects of transcription factor binding on polygenic disease risk. *Nature genetics*, 50(10):1483,
403 2018.
- 404 [8] Alexandre Melnikov, Anand Murugan, Xiaolan Zhang, Tiberiu Tesileanu, Li Wang, Peter Rogov, Soheil
405 Feizi, Andreas Gnirke, Curtis G Callan, Justin B Kinney, Manolis Kellis, Eric S Lander, and Tar-
406 jei S Mikkelsen. Systematic dissection and optimization of inducible enhancers in human cells using a
407 massively parallel reporter assay. *Nature Biotechnology*, 30(3):271–277, February 2012.
- 408 [9] Eilon Sharon, Yael Kalma, Ayala Sharp, Tali Raveh-Sadka, Michal Levo, Danny Zeevi, Leeat Keren,
409 Zohar Yakhini, Adina Weinberger, and Eran Segal. Inferring gene regulatory logic from high-throughput
410 measurements of thousands of systematically designed promoters. *Nature biotechnology*, 30(6):521, 2012.
- 411 [10] Rupali P Patwardhan, Joseph B Hiatt, Daniela M Witten, Mee J Kim, Robin P Smith, Dalit May, Choli
412 Lee, Jennifer M Andrie, Su-In Lee, Gregory M Cooper, et al. Massively parallel functional dissection
413 of mammalian enhancers in vivo. *Nature biotechnology*, 30(3):265, 2012.

- 414 [11] Jamie C Kwasmieski, Ilaria Mogno, Connie A Myers, Joseph C Corbo, and Barak A Cohen. Com-
415 plex effects of nucleotide variants in a mammalian cis-regulatory element. *Proceedings of the National*
416 *Academy of Sciences*, 109(47):19498–19503, 2012.
- 417 [12] Alexander B Rosenberg, Rupali P Patwardhan, Jay Shendure, and Georg Seelig. Learning the Sequence
418 Determinants of Alternative Splicing from Millions of Random Sequences. *Cell*, 163(3):698–711, October
419 2015.
- 420 [13] Michael D Wilson, Nuno L Barbosa-Morais, Dominic Schmidt, Caitlin M Conboy, Lesley Vanes, Vic-
421 tor L J Tybulewicz, Elizabeth M C Fisher, Simon Tavaré, and Duncan T Odom. Species-Specific
422 Transcription in Mice Carrying Human Chromosome 21. *Science*, 322(5900):434–438, October 2008.
- 423 [14] Diego Villar, Paul Flicek, and Duncan T Odom. Evolution of transcription factor binding in metazoans
424 — mechanisms and functional implications. *Nature Reviews Genetics*, 15(4):221–233, March 2014.
- 425 [15] Samuel A Lambert, Arttu Jolma, Laura F Campitelli, Pratyush K Das, Yimeng Yin, Mihai Albu,
426 Xiaoting Chen, Jussi Taipale, Timothy R Hughes, and Matthew T Weirauch. The Human Transcription
427 Factors. *Cell*, 172(4):650–665, February 2018.
- 428 [16] Feng Yue, Yong Cheng, Alessandra Breschi, Jeff Vierstra, Weisheng Wu, Tyrone Ryba, Richard Sand-
429 strom, Zhihai Ma, Carrie Davis, Benjamin D Pope, et al. A comparative encyclopedia of dna elements
430 in the mouse genome. *Nature*, 515(7527):355, 2014.
- 431 [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
432 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 433 [18] ENCODE Project Consortium. An integrated encyclopedia of dna elements in the human genome.
434 *Nature*, 489(7414):57, 2012.
- 435 [19] Anshul Kundaje, Wouter Meuleman, Jason Ernst, Misha Bilenky, Angela Yen, Alireza Heravi-Moussavi,
436 Pouya Kheradpour, Zhizhuo Zhang, Jianrong Wang, Michael J Ziller, et al. Integrative analysis of 111
437 reference human epigenomes. *Nature*, 518(7539):317, 2015.
- 438 [20] Alistair RR Forrest, Hideya Kawaji, Michael Rehli, J Kenneth Baillie, Michiel JL De Hoon, Vanja
439 Haberle, Timo Lassmann, Ivan V Kulakovskiy, Marina Lizio, Masayoshi Itoh, et al. A promoter-level
440 mammalian expression atlas. *Nature*, 507(7493):462, 2014.
- 441 [21] Darren A Cusanovich, Andrew J Hill, Delasa Aghamirzaie, Riza M Daza, Hannah A Pliner, Joel B
442 Berletch, Galina N Filippova, Xingfan Huang, Lena Christiansen, William S DeWitt, et al. A single-cell
443 atlas of in vivo mammalian chromatin accessibility. *Cell*, 174(5):1309–1324, 2018.
- 444 [22] Kazuhiro R Nitta, Arttu Jolma, Yimeng Yin, Ekaterina Morgunova, Teemu Kivioja, Junaid Akhtar,
445 Korneel Hens, Jarkko Toivonen, Bart Deplancke, Eileen EM Furlong, et al. Conservation of transcription
446 factor binding specificities across 600 million years of bilateria evolution. *Elife*, 4:e04837, 2015.
- 447 [23] GTEx Consortium et al. Genetic effects on gene expression across human tissues. *Nature*, 550(7675):
448 204, 2017.
- 449 [24] Gerald D Fischbach and Catherine Lord. The simons simplex collection: a resource for identification of
450 autism genetic risk factors. *Neuron*, 68(2):192–195, 2010.
- 451 [25] Joon-Yong An, Kevin Lin, Lingxue Zhu, Donna M Werling, Shan Dong, Harrison Brand, Harold Z Wang,
452 Xuefang Zhao, Grace B Schwartz, Ryan L Collins, et al. Genome-wide de novo risk score implicates
453 promoter variation in autism spectrum disorder. *Science*, 362(6420):eaat6576, 2018.
- 454 [26] Jian Zhou, Christopher Y Park, Chandra L Theesfeld, Aaron K Wong, Yuan Yuan, Claudia Scheckel,
455 John J Fak, Julien Funk, Kevin Yao, Yoko Tajima, et al. Whole-genome deep-learning analysis identifies
456 contribution of noncoding mutations to autism risk. *Nature Genetics*, 2019.

- 457 [27] Konrad J Karczewski, Laurent C Francioli, Grace Tiao, Beryl B Cummings, Jessica Alföldi, Qingbo
458 Wang, Ryan L Collins, Kristen M Laricchia, Andrea Ganna, Daniel P Birnbaum, et al. Variation across
459 141,456 human exomes and genomes reveals the spectrum of loss-of-function intolerance across human
460 protein-coding genes. *BioRxiv*, page 531210, 2019.
- 461 [28] Abraham S Weintraub, Charles H Li, Alicia V Zamudio, Alla A Sigova, Nancy M Hannett, Daniel S
462 Day, Brian J Abraham, Malkiel A Cohen, Behnam Nabet, Dennis L Buckley, et al. Yy1 is a structural
463 regulator of enhancer-promoter loops. *Cell*, 171(7):1573–1588, 2017.
- 464 [29] Jason Homsy, Samir Zaidi, Yufeng Shen, James S Ware, Kaitlin E Samocha, Konrad J Karczewski,
465 Steven R DePalma, David McKean, Hiroko Wakimoto, Josh Gorham, et al. De novo mutations in
466 congenital heart disease with neurodevelopmental and other congenital anomalies. *Science*, 350(6265):
467 1262–1266, 2015.
- 468 [30] Sheng Chih Jin, Jason Homsy, Samir Zaidi, Qiongshi Lu, Sarah Morton, Steven R DePalma, Xue Zeng,
469 Hongjian Qi, Weni Chang, Michael C Sierant, et al. Contribution of rare inherited and de novo variants
470 in 2,871 congenital heart disease probands. *Nature genetics*, 49(11):1593, 2017.
- 471 [31] Rachelly Normand, Wenfei Du, Mayan Briller, Renaud Gaujoux, Elina Starosvetsky, Amit Ziv-Kenet,
472 Gali Shalev-Malul, Robert J Tibshirani, and Shai S Shen-Orr. Found in translation: a machine learning
473 model for mouse-to-human inference. *Nature methods*, 15(12):1067, 2018.
- 474 [32] Dent Earl, Ngan Nguyen, Glenn Hickey, Robert S Harris, Stephen Fitzgerald, Kathryn Beal, Igor
475 Seledtsov, Vladimir Molodtsov, Brian J Raney, Hiram Clawson, et al. Alignathon: a competitive
476 assessment of whole-genome alignment methods. *Genome research*, 24(12):2077–2089, 2014.
- 477 [33] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
478 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 479 [34] Kelley, David R. Basenji v0.3. <https://github.com/calico/basenji>, 2019.
- 480 [35] Shifu Chen, Yanqing Zhou, Yaru Chen, and Jia Gu. fastp: an ultra-fast all-in-one fastq preprocessor.
481 *Bioinformatics*, 34(17):i884–i890, 2018.
- 482 [36] Heng Li and Richard Durbin. Fast and accurate short read alignment with burrows–wheeler transform.
483 *bioinformatics*, 25(14):1754–1760, 2009.
- 484 [37] Carrie A Davis, Benjamin C Hitz, Cricket A Sloan, Esther T Chan, Jean M Davidson, Idan Gabdank,
485 Jason A Hilton, Kriti Jain, Ulugbek K Baymuradov, Aditi K Narayanan, et al. The encyclopedia of
486 dna elements (encode): data portal update. *Nucleic acids research*, 46(D1):D794–D801, 2017.
- 487 [38] AFA Smit, R Hubley, and P Green. Repeatmasker open-4.0. <http://www.repeatmasker.org>, 2015.
488 Accessed 2019-4-20.
- 489 [39] Joseph K Pickrell, Daniel J Gaffney, Yoav Gilad, and Jonathan K Pritchard. False positive peaks in
490 chip-seq and other sequencing-based functional assays caused by unannotated high copy number regions.
491 *Bioinformatics*, 27(15):2144–2146, 2011.
- 492 [40] Mehran Karimzadeh, Carl Ernst, Anshul Kundaje, and Michael M Hoffman. Umap and bismap: quan-
493 tifying genome and methylome mappability. *Nucleic acids research*, 46(20):e120–e120, 2018.
- 494 [41] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*,
495 2016.
- 496 [42] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin,
497 Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine
498 learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*,
499 pages 265–283, 2016.
- 500 [43] Scott Schwartz, W James Kent, Arian Smit, Zheng Zhang, Robert Baertsch, Ross C Hardison, David
501 Haussler, and Webb Miller. Human–mouse alignments with blastz. *Genome research*, 13(1):103–107,
502 2003.