# Genetic Architect: Discovering Genomic Structure with Learned Neural Architectures

**Laura Deming**,[*] **Sasha Targ**,[*] **Nate Sauder, Diogo Almeida, Chun Jimmie Ye**
Institute for Human Genetics
University of California, San Francisco, CA 94143, USA

Enlitic
San Francisco, CA 94111, USA
`ldeming.www@gmail.com, sasha.targ@ucsf.edu`
`nate@enlitic.com, diogo@enlitic.com, jimmie.ye@ucsf.edu`

## Abstract

Each human genome is a 3 billion base pair set of encoding instructions. Decoding the genome using deep learning fundamentally differs from most tasks, as we do not know the full structure of the data and therefore cannot design architectures to suit it. As such, architectures that fit the structure of genomics should be learned not prescribed. Here, we develop a novel search algorithm, applicable across domains, that discovers an optimal architecture which simultaneously learns general genomic patterns and identifies the most important sequence motifs in predicting functional genomic outcomes. The architectures we find using this algorithm succeed at using only RNA expression data to predict gene regulatory structure, learn human-interpretable visualizations of key sequence motifs, and surpass state-of-the-art results on benchmark genomics challenges.

## 1   Introduction

Deep learning demonstrates excellent performance on tasks in computer vision, text and many other fields. Most deep learning architectures consist of matrix operations composed with non-linearity activations. Critically, the problem domain governs how matrix weights are shared. In convolutional neural networks – dominant in image processing – translational equivariance ("edge/color detectors are useful everywhere") is encoded through the use of the convolution operation; in recurrent networks – dominant in sequential data – temporal transitions are captured by shared hidden-to-hidden matrices. These architectures mirror human intuitions and priors on the structure of the underlying data. Genomics is an excellent domain to study how we might learn optimal architectures on poorly-understood data because while we have intuition that local patterns and long-range sequential dependencies affect genetic function, much structure remains to be discovered.

The genome is a very challenging data type, because although we have tens of thousands of whole genome sequences, we understand only a small subset of base pairs within each sequence. While the genetic code allows us to annotate the 5% of the genome encoding proteins (∼20,000 genes in the human genome), we do not have a "grammar" for decoding the rest of the non-coding sequences (90-95% of the mouse and human genomes) important for gene regulation, evolution of species and susceptibility to diseases. The availability of a wealth of genomic assays (PBM, CHIP-seq, Hi-C) allows us to directly measure the function of specific regions of the genome, creating an enormous opportunity to decode non-coding sequences. However, the overwhelming volume of new data makes our job as decoders of the genome quite complex. The design and application of new domain-specific

---

[*]Equal Contribution

architectures to these datasets is a promising approach for automating interpretation of genomic information into forms that humans can grasp.

## 2 Related Work

Inspired by human foveal attention where global glances drive sequential local focus, attention components have been added to neural networks yielding state-of-the-art results on tasks as diverse as caption generation, machine translation, protein sublocalization, and differentiable programming. There are two main architectural implementations: hard attention, where the network's focus mechanism non-differentiably samples from the available input, and soft attention, where the component outputs an expected glimpse using a weighted average. Beyond biological inspiration, these components enable improved performance and excellent intepretability. Other techniques have been applied for interpreting neural networks without changing their architectures (Simonyan et al. [2013], Zeiler and Fergus [2014], Springenberg et al. [2014]), but these are simply heuristics for finding the relevant regions of an input and do not work with all existing modern neural network components.

Previous groups have demonstrated excellent progress applying deep learning to genomics. Both Alipanahi et al. [2015] and Lanchantin et al. [2016] provide initial results on the task of learning which sequences a transcription factor (a biological entity which affects gene expression) can bind using convolutional architectures. This problem appears suited for convolution, as motifs determining binding are expected to be modular ($\sim$7-10 base pair units) and the setup of the task (preselected input sequences of fixed short length) does not allow for learning significant long-term dependencies. In particular, Alipanahi et al. [2015] demonstrated that a single-layer convolutional neural network, DeepBind, outperformed 26 other tested machine learning approaches in predicting probe intensities on protein binding microarrays from the DREAM5 PBM challenge, and then showed that the same architecture generalized to the related task of predicting transcription factor binding sites (TFBSs) from sequencing measurements of bound DNA. Subsequently, Lanchantin et al. [2016] showed that a deeper network with the addition of highway layers improved on DeepBind results in the majority of cases tested [Srivastava et al., 2015]. In addition, Basset [Kelley et al., 2015], an architecture trained to predict motifs of accessible DNA from sequencing regions of open chromatin, was able to map half of the first layer convolutional filters to human TFBSs.

## 3 Development of Genetic Architect

Deep learning algorithm development is often dependent on the knowledge of human domain experts. Researchers in domains such as computer vision and natural language processing have spent much more time tuning architectures than in genomics. The challenge in genomics is that our insufficient understanding of biology limits our ability to inform architectural decisions based on data. Early genomic deep learning architectures have shown promising results but have undertaken only limited exploration of the architectural search space over possible components. In addition, not all components work well together, and there is evidence optimal component choice is highly dependent on the domain. Accordingly, we design a novel road-map for applying deep learning to data on which we have limited prior understanding, by developing an iterative architecture search over standard and cutting-edge neural net building blocks.

Prior approaches to architecture search focus on finding the best architecture in a single step, rather than sequentially learning more about the architecture space and iteratively improving models (Bergstra et al. [2011], Bergstra and Bengio [2012], Snoek et al. [2012]). Our framework understands the results allowing us to sequentially narrow the search space and learn about which combinations of components are most important. Since our algorithm limits the most important hyperparameters to their best ranges, they no longer dominate the search space and we discover additional hyperparameters that are most important and can help us create a highly tuned architecture. The sequential nature allows us to fork our architectural search into independent subspaces of coadapted components, thus enabling further search in each parallel branch to be exponentially more efficient than considering the union of all promising architectures.

The heart of the framework is an interactive visualization tool (Figure 4). Given any hyperparameter optimization run, it produces common patterns for the best few datapoints and presents this information in highly-interpretable decision trees showing effective architectural subspace and plots of the

Table 1: Mean and median AUC of models and percentage of datasets on which each model outperforms DeepBind or DeepMotif.

| model | mean AUC | median AUC | vs DeepBind | vs DeepMotif |
|---|---|---|---|---|
| DeepBind | 0.904 | 0.932 | - | - |
| DeepMotif | 0.927 | 0.949 | 85.2 | - |
| AttentionNet | 0.933 | 0.952 | 92.6 | 67.6 |

interactions between the most significant hyperparameters, informing general domain intuition and guiding future experiments. The framework is general enough to be applied to other domains, and is orthogonal to existing hyperparameter optimization algorithms. These algorithms can be applied in the inner loop of the sequential search of our tool, which then interprets the results and informs the user about the domain and how to manually prune the search space.

We employ Genetic Architect to discover an optimal architecture for a novel genome annotation task, regression to predict lineage-specific gene expression based on genomic sequence inputs, for which six stages of architecture search were required. Figure 1A shows the sequential process of architecture search, the most important findings at each stage of the process, and tool-guided division of the search into two separate promising architectures. By splitting effective architectures into separate branches for further optimization, our tool identifies high-performing but architecture-specific choices that may be difficult to notice when architectures are mixed together.

The application of our tool demonstrates the power in refining architectural components that dominate results to uncover additional hyperparameter combinations that perform well together. Several examples we encounter during use of the tool for design of architectures for genomics follow: 1) removal of batch normalization demonstrated clear superiority of exponential linear units, 2) dimensionality reduction in the middle of the convolutional network module was beneficial to the recurrent-based architectures (perhaps since it decreased the distance of long-range dependencies), and 3) in contrast, non-recurrent architectures required wider layers (likely to enable processing of long-range dependencies in final dense layers). In our search over architectures using soft attention, we found that fully-connected layers were preferred to convolutional layers as it made processing global information more important. Finally, only by proceeding through several steps of optimization did we find the unintuitive result that bidirectional LSTMs did not help with attentional models (perhaps because the preceding layer effectively attends to a single location, making it difficult to combine information from both directions).

The final models learned by Genetic Architect consist of several initial layers of convolutions, residual blocks, an LSTM layer in the case of the PromoterNet architecture, and an attention-based dimensionality reducing step followed by fully-connected layers. Previous approaches to genome annotation use convolutional networks, which are ideal for detecting local features. However, more closely approximating the structure of genomic information would take into account that a real genome is a sequence, not a disjointed set, of local features – an input type on which recurrent architectures generally excel. In addition, with larger sequences to analyze (identifiable promoter sequences reach hundreds of base pairs in length), a neural network must learn to focus on the most important parts of the sequence and integrate new information derived from each part with the contextual information of the previously-seen sequence. As such, long genomic sequences seem an ideal fit for the recurrent attentional models learned.

## 4 Experimental Results

### 4.1 Tasks

#### 4.1.1 Transcription factor binding site (TFBS) classification

The TFBS binary classification task was proposed by Alipanahi et al. [2015] and used as a benchmark by [Lanchantin et al., 2016]. The basic motivation is to learn a classifier that correctly predicts, from empirical binding data on a training sample of short DNA sequences, which sequences in a separate test set are TFBS (likely to be bound by biological entities, in this case, a given transcription factor protein).
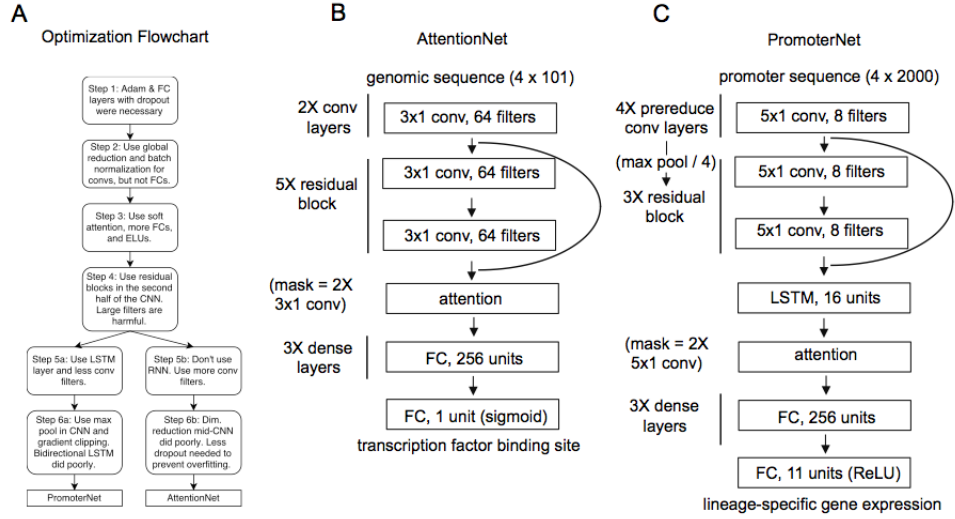
Figure 1: Schematic of hyperparameter optimization and final architecture designs. A) Overview of steps taken in hyperparameter optimization to generate AttentionNet and PromoterNet. B) Attention-Net architecture. C) PromoterNet architecture.
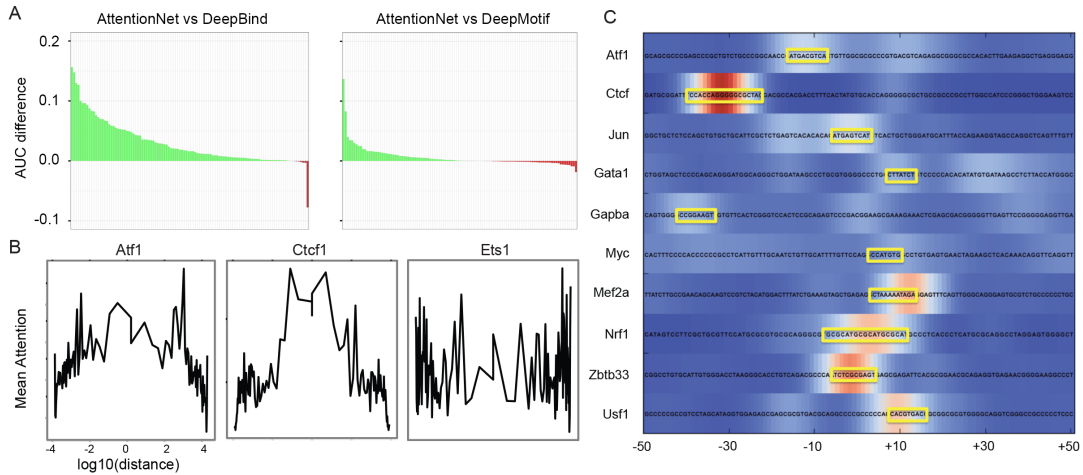


Figure 2: Results of AttentionNet on transcription factor binding site (TFBS) task. A) AttentionNet models outperform DeepMotif and DeepBind models trained on corresponding datasets. Each bar represents the difference in AUC for one of 108 different datasets. B) Mean of attention mask over all sequences in experiment. C) Recovery of transcription factor motifs by visualization of attention masks produced by AttentionNet over example sequences.

The input and target data for the TFBS classification task consists of 108 datasets with an average of ~31,000 sequences of 101 characters per dataset. Each sequence is a string of base pairs (A, C, G, or T) and is transformed into an array with one-hot encoding. Each sequence has an associated label (1 or 0) which indicates if this sequence is a TFBS. Each dataset represents a different chromatin immunoprecipitation sequencing (ChIP-seq) experiment with a specified transcription factor, and each sequence in the dataset a potential binding site. For each positive example, a negative example is generated. The data included in the TFBS classification task derive from ENCODE CHIP-seq experiments performed in K562 transformed human cell lines [Consortium, 2012].
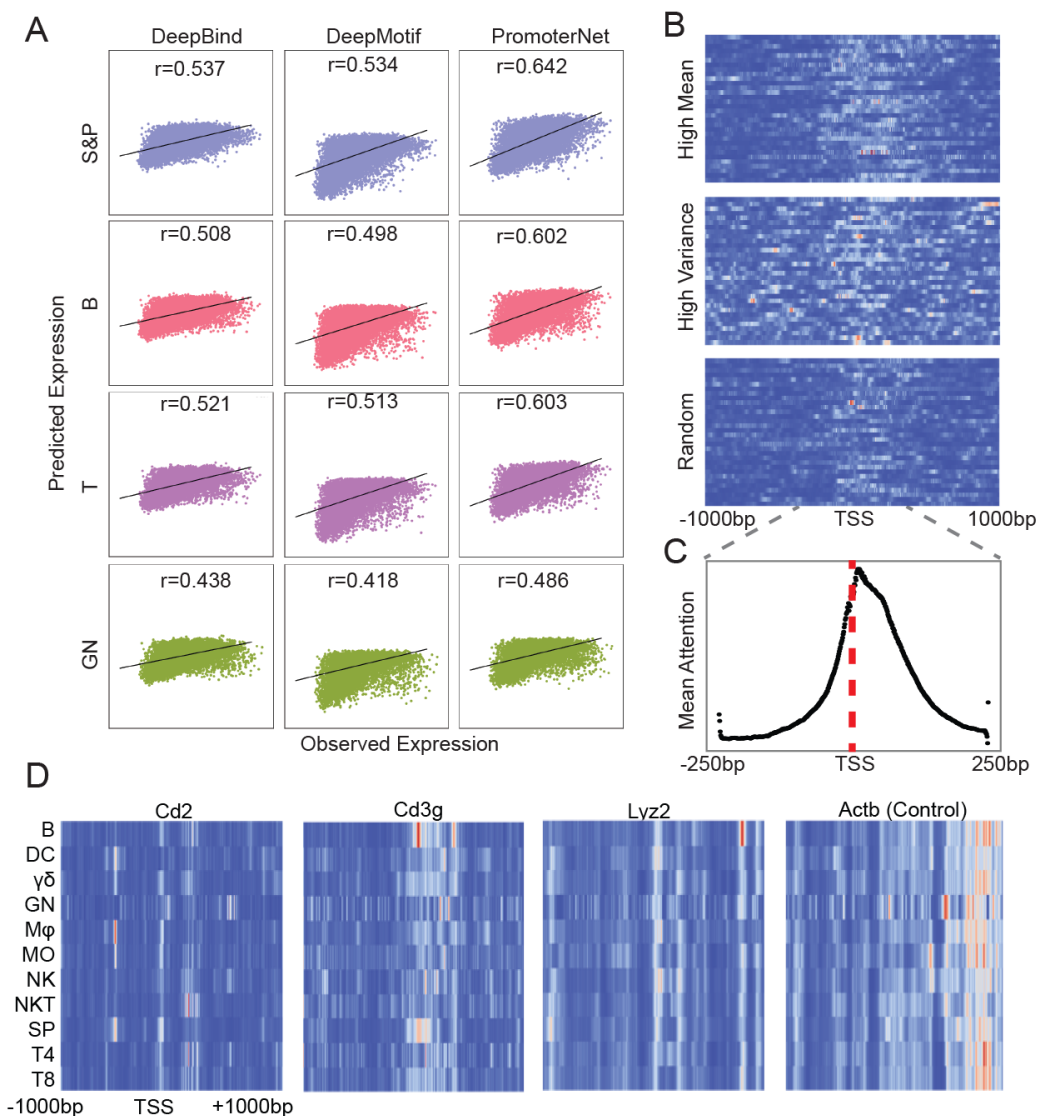
4

Figure 3: Results of PromoterNet on ImmGen lineage-specific expression prediction (ILSEP) task. A) Comparison of predicted versus observed gene expression for DeepBind, DeepMotif, and PromoterNet architectures. B) Visualization of attention mask over selected promoter sequences. C) Mean attention mask over all promoters. D) Visualization of attention masks learned by models trained on data from single lineages.

### 4.1.2 ImmGen lineage-specific expression prediction (ILSEP) regression

In addition to the TFBS classification problem, neural network architectures could be extended to treat a much broader and complex variety of problems to do with interpreting biological data. Here, we develop a novel genomic benchmark task, ILSEP, which requires regression to predict empirically-determined related target data, namely, prediction of the amount of various biological entities produced in different cellular contexts given an input genomic sequence.

The input dataset for the ILSEP task is 14,116 one-hot encoded (4,2000) input promoter sequences and corresponding (243,) floating point gene expression outputs ranging between 2.60 and 13.95 (see appendix for details). We split the dataset using 10-fold cross validation to obtain predictions for all promoter gene expression pairs.

## 4.2 Results on TFBS

### 4.2.1 Model performance

We benchmark the performance of AttentionNet models learned by hyperparameter optimization described above against published state-of-the-art neural network models on the TFBS task, DeepBind [Alipanahi et al., 2015] and DeepMotif [Lanchantin et al., 2016]. To compare the architectures, we train models for each of 108 datasets, as in Lanchantin et al. [2016]. In a head-to-head comparison on each dataset, AttentionNet outperforms DeepMotif in 67.6% of cases and the mean AUC across datasets for AttentionNet is 0.933, improving over both DeepMotif (0.927) and DeepBind (0.904) (Table 1).

### 4.2.2 Prediction and visualization of genomic information

Interpretable information about sequence features is an important consideration for genomic learning tasks where fundamental understanding of biology is as important as prediction power. We hypothesize that a net which performed well on the TFBS classification task would be able to make biologically meaningful inferences about the sequence structure. We show that the mean attention weights across all positive sequences show a distinct "footprint" of transcription factor (TF) binding consistent with known nucleotide preferences within each sequence (Figure 2B). Further, visualizing the attention mask (with the addition of Gaussian blur) across input sequences for 10 representative TFs showed the net focusing its attention on parts of the sequence known to be regulatory (Figure 2C).

To see if we could directly obtain motif sequences from the net, we took 10 nucleotides surrounding the position with highest attention for each of the top 100 sequences of a TF and averaged across the motifs. We took the maximum score for each nucleotide per position and queried the results against JASPAR, the "gold standard" TFBS database (with $q < 0.5$) [Mathelier et al., 2015]. 30/57 motifs possible to check (i.e. in JASPAR) were correct, and 39/57 corresponded to at least one transcription factor. By additionally searching the top 3 recurring sequences attended to for each TF, we recover a total of 42/57 correct motifs.

## 4.3 Results on ILSEP

### 4.3.1 Model performance

The PromoterNet architecture demonstrates a marked gain in performance over DeepBind and DeepMotif architectures adapted to the ILSEP regression task, achieving an average Pearson r correlation value of 0.587 between out-of-sample predictions and target expression values across lineages, compared to 0.506 and 0.441 for DeepBind [Alipanahi et al., 2015] and DeepMotif [Lanchantin et al., 2016] respectively (Figure 3A). We also train PromoterNet architectures on single task regression with a separate model for each of the 11 lineages and on cell type specific multi-task regression with one output unit for each of 243 cell types, which obtains similar improvements in average Pearson r correlation value of 0.592 over 0.502 for DeepBind and 0.498 for DeepMotif.

### 4.3.2 Promoter element recovery and visualization of proximal regulatory elements

Visualization of attention mask weights from the PromoterNet model reveals attended locations over promoter sequences of 32 genes selected for highest mean expression across lineages are enriched directly adjacent to the TSS, suggesting that properties of the core promoter sequence constitute the most informative features for genes that do not show differences in expression across lineages (Figure 3B) (see appendix for list of genes). In contrast, attended locations over promoter sequences of 32 genes with maximal variance in expression across lineages span a much greater range of positions. This indicates that in genes with the greatest degree of lineage-specific expression, informative sequence features can occur throughout the promoter sequence. This observation merits follow up given previous reports that the performance of (non-deep) classifiers for cell type specific expression tasks trained only on TSS proximal promoter information is close to that of a random classifier [Natarajan et al., 2012]. Consistent with accepted understanding that TSS proximal regions contain genomic elements that control gene expression levels, we observe the maximum of average attention mask weights across all promoters occurs at the center of input sequences, which corresponds to

6

Table 2: Search space explored for AttentionNet and PromoterNet architectures, including techniques from Maas et al. [2013], Graham [2014], Shah et al. [2016], Ioffe and Szegedy [2015], He et al. [2015], Hochreiter and Schmidhuber [1997], Kingma and Ba [2014], Sutskever et al. [2013], Srivastava et al. [2014].

| Hyperparameter | Values |
|---|---|
| conv filter size | 3, 5, 7, 9 |
| nonlinearity | ReLU, Leaky ReLU, Very Leaky ReLU, ELU |
| using batch normalization for convs | True, False |
| number of conv filters | 8, 16, 32, 64 |
| number of convs before dim. reduction | 1, 2, 3, 4, 5 |
| using residual blocks before dim. reduction | True, False |
| type of dim. reduction | None, Max Pool, Mean Pool, Strided Conv |
| dim. reduction stride | 2, 4, 8 |
| number of convs after dim. reduction | 1, 2, 3, 4, 5 |
| using residual blocks after dim. reduction | True, False |
| number of RNN layers | 0, 1, 2 |
| number of units in RNN | 16, 32, 64 |
| RNN type | Simple RNN, LSTM |
| using bidirectional RNNs | True, False |
| RNN gradient clipping | 0, 1, 5, 15 |
| global reduction type | None, Attention, Max Pool, Mean Pool |
| number of FC layers | 0, 1, 2, 3 |
| number of units in FC | 64, 128, 256, 512, 1024 |
| using batch normalization for FCs | True, False |
| dropout probability for FCs (after) | 0., 0.1, 0.2, 0.3, 0.4, 0.5, 0.6 |
| L2 regularization | 0, 1e-3, 1e-4, 1e-5 |
| optimizer | Adam, SGD w/ Nesterov Momentum |
| learning rate scale | 0.03, 0.1, 0.3, 1.0, 2.0, 10. |
| batch size | 25, 50, 125, 250 |

core promoter elements required for recruitment of transcriptional machinery [Maston et al., 2006] (Figure 3C).

PromoterNet models trained for multi-task regression result in a global attention mask output across all lineages. To investigate whether the PromoterNet architecture is capable of learning distinct features for each lineage, we also visualize attention weights for a given promoter sequence from separate models, each trained on expression data for a single lineage. We find that genes selected for maximal variance in expression demonstrate distinct patterns of learned attention across lineages, while a shared pattern of attention is learned for a control gene with high mean expression in all lineages even when each lineage was trained on a separate model (Figure 3D).

# 5   Conclusion

We tackle the problem of discovering architectures on datasets where human priors are not available. To do so we create a novel architecture search framework that is domain agnostic, is capable of sequential architectural subspace refinement and informing domain understanding, and is composable with existing hyperparameter optimization schemes. Using this search algorithm, we create state-of-the art architectures on significant challenges in the domain of genomics utilizing a combination of standard and cutting-edge components. In particular, the learned architecture is capable of simultaneous discovery of local and non-local patterns, important subsequences, and sequential composition thereby capturing substantial genomic structure.

Table 3: AttentionNet architecture

| Layer |
| --- |
| 3x1 conv 64 filters + BN + ELU |
| 3x1 conv 64 filters + BN + ELU |
| residual block (w/ 3x1 conv 64 filters + BN + ELU + 3x1 conv 64 filters + BN) |
| residual block (w/ 3x1 conv 64 filters + BN + ELU + 3x1 conv 64 filters + BN) |
| residual block (w/ 3x1 conv 64 filters + BN + ELU + 3x1 conv 64 filters + BN) |
| residual block (w/ 3x1 conv 64 filters + BN + ELU + 3x1 conv 64 filters + BN) |
| residual block (w/ 3x1 conv 64 filters + BN + ELU + 3x1 conv 64 filters + BN) |
| attention (w/ 3x1 conv 64 filters + BN + tanh + 3x1 conv 1 filter + BN + softmax) |
| FC 256 units + ELU + 0.2 dropout |
| FC 256 units + ELU + 0.2 dropout |
| FC 256 units + ELU + 0.2 dropout |
| FC 1 unit + sigmoid |

Table 4: PromoterNet architecture

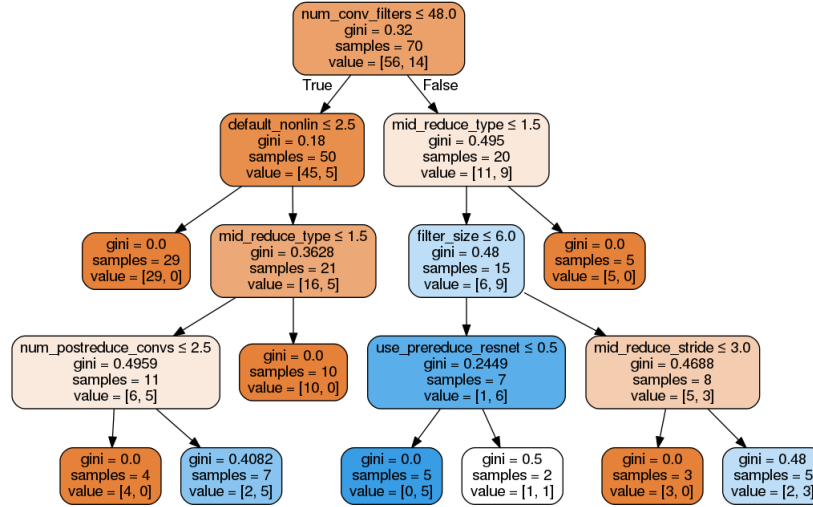| Layer |
| --- |
| 5x1 conv 8 filters + BN + ELU |
| 5x1 conv 8 filters + BN + ELU |
| 5x1 conv 8 filters + BN + ELU |
| 5x1 conv 8 filters + BN + ELU |
| 4x1 maxpool, stride 4 |
| residual block (w/ 5x1 conv 8 filters + BN + ELU + 5x1 conv 8 filters + BN) |
| residual block (w/ 5x1 conv 8 filters + BN + ELU + 5x1 conv 8 filters + BN) |
| residual block (w/ 5x1 conv 8 filters + BN + ELU + 5x1 conv 8 filters + BN) |
| LSTM 16 units |
| attention (w/ 5x1 conv 8 filters + BN + tanh + 5x1 conv 1 filter + BN + softmax) |
| FC 256 units + ELU + 0.3 dropout |
| FC 256 units + ELU + 0.3 dropout |
| FC 256 units + ELU + 0.3 dropout |
| FC 1 unit + sigmoid |



Figure 4: Example decision tree output from Genetic Architect visualization tool depicting significant hyperparameters for models with top 20% performance.

# 6 Appendix

## 6.1 ILSEP data processing

For the input sequences in the ILSEP task, we use sequences spanning 1 kilobase upstream and downstream of transcriptional start sites (TSS), the region of the promoter at which production of the gene is initiated, for 17,565 genes from the Eukaryotic Promoter Database [epd]. For the corresponding labels, we obtain expression data (log2 normalized microarray intensities) for each of these genes in each of 243 immune cell types from the ImmGen Consortium April 2012 release, which contains data for 21,755 genes x 243 immune cell types after quality control according to published methods [Ericson et al.]. Intersection of these two datasets by gene results in a dataset of 14,116 input promoter sequences and expression value target pairs.

To create lineage-specific gene expression value targets, we combine cell types into 11 groups following the lineage tree outlined in previous work: B cells (B), dendritic cells (DC), gamma delta T cells ($\gamma\delta$), granulocytes (GN), macrophages (M$\phi$), monocytes (MO), natural killer cells (NK), stem and progenitor cells (SP), CD4+ T cells (T4), and CD8+ T cells (T8), and average expression values across samples within each group [Jojic et al., 2013].

## 6.2 Selected genes (Figure 3B)

Highest mean expression across lineages: Rac2, Rpl28, Pfn1, Rpl9, Ucp2, Tmsb4x, Tpt1, Rplp1, Hspa8, Srgn, Rpl27a, Rpl13a, Cd53, Eef2, Rps26, Cfl1, Ppia, Gm9104, Rps2, Rps27, Actg1, Laptm5, Rpl21, Eef1a1, Rplp0, Gm15427, Pabpc1, B2m, Gapdh, Actb, Rpl17, Rps6

Highest variance in expression across lineages: Plbd1, Tlr13, Tyrobp, Ifitm2, Pld4, Pla2g7, Gda, Cd96, Gzma, Nkg7, Ctsh, Klrb1c, Ccl6, Prkcq, Itgam, Sfpi1, Itk, Ms4a4b, Alox5ap, Ly86, Cd2, Fcer1g, Gimap3, Il2rb, Gimap4, Ifitm6, Cybb, Ifitm3, Mpeg1, H2-Aa, Cd3g, Lyz2

Random control: Krt84, Lrrc8b, 8030411F24Rik, Syngr2, Spint3, Slc17a4, Slc22a23, Thoc6, AF529169, Phf5a, Yif1b, 4930467E23Rik, Pgam1, Pcdhb1, Bak1, Neu3, Plcb2, Fabp4, Srgap1, Olfr1339, Sox12, Atg7, Gdf10, 1810008A18Rik, 1700011A15Rik, Anks4b, Magea2, Pygb, Spc25, Rras2, Slc28a3, 9130023H24Rik

# References

Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.

Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Computer vision–ECCV 2014*, pages 818–833. Springer, 2014.

Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.

Babak Alipanahi, Andrew Delong, Matthew T Weirauch, and Brendan J Frey. Predicting the sequence specificities of dna-and rna-binding proteins by deep learning. *Nature biotechnology*, 2015.

Jack Lanchantin, Ritambhara Singh, Zeming Lin, and Yanjun Qi. Deep motif: Visualizing genomic sequence classifications. *arXiv preprint arXiv:1605.01133*, 2016.

Rupesh K Srivastava, Klaus Greff, and Jürgen Schmidhuber. Training very deep networks. In *Advances in Neural Information Processing Systems*, pages 2368–2376, 2015.

David R Kelley, Jasper Snoek, and John Rinn. Basset: Learning the regulatory code of the accessible genome with deep convolutional neural networks. *bioRxiv*, page 028399, 2015.

James S Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In *Advances in Neural Information Processing Systems*, pages 2546–2554, 2011.

James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *The Journal of Machine Learning Research*, 13(1):281–305, 2012.

Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959, 2012.

ENCODE Project Consortium. An integrated encyclopedia of dna elements in the human genome. *Nature*, 489(7414):57–74, 2012.

Anthony Mathelier, Oriol Fornes, David J Arenillas, Chih-yu Chen, Grégoire Denay, Jessica Lee, Wenqiang Shi, Casper Shyr, Ge Tan, Rebecca Worsley-Hunt, et al. Jaspar 2016: a major expansion and update of the open-access database of transcription factor binding profiles. *Nucleic acids research*, page gkv1176, 2015.

Anirudh Natarajan, Galip Gürkan Yardımcı, Nathan C. Sheffield, Gregory E. Crawford, and Uwe Ohler. Predicting cell-type–specific gene expression from regions of open chromatin. *Genome Research*, 22(9):1711–1722, 2012. doi: 10.1101/gr.135129.111. URL http://genome.cshlp.org/content/22/9/1711.abstract.

Glenn A. Maston, Sarah K. Evans, and Michael R. Green. Transcriptional regulatory elements in the human genome. *Annual Review of Genomics and Human Genetics*, 7(1):29–59, 2006. doi: 10.1146/annurev.genom.7.080505.115623. URL http://dx.doi.org/10.1146/annurev.genom.7.080505.115623. PMID: 16719718.

Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. ICML*, volume 30, page 1, 2013.

Benjamin Graham. Spatially-sparse convolutional neural networks. *arXiv preprint arXiv:1409.6070*, 2014.

Anish Shah, Eashan Kadam, Hena Shah, and Sameer Shinde. Deep residual networks with exponential linear unit. *arXiv preprint arXiv:1604.04112*, 2016.

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.

Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.

Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *Proceedings of the 30th international conference on machine learning (ICML-13)*, pages 1139–1147, 2013.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.

Eukaryotic promoter database. URL http://epd.vital-it.ch/mouse/mouse_database.php.

Jeff Ericson, Scott Davis, Jon Lesh, Melissa Howard, Diane Mathis, and Christophe Benoist. Immgen microarray gene expression data: Data generation and quality control pipeline. URL www.immgen.org/Protocols/ImmGenQCDocumentation_ALL-DataGeneration_0612.pdf.

Vladimir Jojic, Tal Shay, Katelyn Sylvia, Or Zuk, Xin Sun, Joonsoo Kang, Aviv Regev, Daphne Koller, Immunological Genome Project Consortium, et al. Identification of transcriptional regulators in the mouse immune system. *Nature immunology*, 14(6):633–643, 2013.