

METHOD

Open Access



# ExplaiNN: interpretable and transparent neural networks for genomics

Gherman Novakovsky<sup>1†</sup>, Oriol Fornes<sup>1†</sup>, Manu Saraswat<sup>1,2,3†</sup>, Sara Mostafavi<sup>4</sup> and Wyeth W. Wasserman<sup>1\*</sup>

<sup>†</sup>Gherman Novakovsky, Oriol Fornes, and Manu Saraswat contributed equally to this work and are considered co-first authors.

\*Correspondence: wyeth@cmmt.ubc.ca

<sup>1</sup> Department of Medical Genetics, Centre for Molecular Medicine and Therapeutics, BC Children's Hospital Research Institute, University of British Columbia, Vancouver, BC, Canada

<sup>2</sup> Division of Computational Genomics and Systems Genetics, German Cancer Research Center (DKFZ), Heidelberg, Germany

<sup>3</sup> European Molecular Biology Laboratory (EMBL), Genome Biology Unit, Heidelberg, Germany

<sup>4</sup> Paul G. Allen School of Computer Science and Engineering, University of Washington (UW), Seattle, USA

## Abstract

Deep learning models such as convolutional neural networks (CNNs) excel in genomic tasks but lack interpretability. We introduce ExplaiNN, which combines the expressiveness of CNNs with the interpretability of linear models. ExplaiNN can predict TF binding, chromatin accessibility, and de novo motifs, achieving performance comparable to state-of-the-art methods. Its predictions are transparent, providing global (cell state level) as well as local (individual sequence level) biological insights into the data. ExplaiNN can serve as a plug-and-play platform for pretrained models and annotated position weight matrices. ExplaiNN aims to accelerate the adoption of deep learning in genomic sequence analysis by domain experts.

**Keywords:** Deep learning, Gene regulation, Genomics, Model interpretation, Explainable artificial intelligence, Transcription factors

## Background

High-throughput genomics methods such as ATAC-seq [1] or ChIP-seq [2], which respectively assess genome-wide accessibility of chromatin and binding of transcription factors (TFs), allow functional annotation of DNA elements. The sheer scale of the data generated by these methods precludes manual analyses. Machine and deep learning have become pervasive in large-scale genomic analyses due to their ability to identify meaningful features in massive datasets (reviewed in [3, 4]). For instance, deep learning models have shown superior performance in predicting genome folding [5], chromatin accessibility states [6–9], gene expression levels [10–12], and TF binding sites (TFBSs; reviewed in [13]).

While increasingly complex models are now feasible, some are opaque; they do not readily reveal the features and properties that underlie their predictions [14]. For genomics, there are multiple approaches to improve the interpretability of deep learning models (reviewed in [15]), of which we focus on filter visualization and attribution methods.



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

For convolutional neural networks (CNNs) and related models, a powerful interpretation approach is to visualize the filters of the first convolutional layer as position weight matrices (PWMs), a common format in bioinformatics to represent TFBS patterns [16], by separately aligning the set of sequences activating each filter [17]. The resulting filter PWMs can be assigned a biological annotation by comparing them to TF binding profiles from reference databases like JASPAR [18]. While filter visualization provides an overview of the genomic features that the model has learned in the first layer, the success of this approach is highly dependent on model architectural choices such as max pooling or filter sizes [19]. Moreover, it is difficult to interpret how the model combines the learned motifs within the deeper layers. Furthermore, to gain interpretability at a global level (i.e., how each filter influences the model's predictions), the filters must be nullified sequentially, which is both computationally intensive and dependent on arbitrary thresholds [8, 9].

Attribution methods quantify the importances of individual nucleotides in the input sequences using forward- (e.g., *in silico* mutagenesis [6, 20]) or back-propagation [21, 22]. These importance scores can be the basis for further clustering of activating subsequences into PWMs [23], which, as with filter visualization, can in turn be annotated biologically by comparing them to known TF motifs. Although attribution methods provide local interpretability by identifying important nucleotides in the input sequences, quantifying the contribution of each feature to the overall model's predictions (i.e., global interpretability) remains challenging [24]. Noteworthy, neither approach is transparent as to how the model makes predictions.

Unlike deep learning models, linear models are interpretable and transparent: the basis on which they make predictions can be readily understood by evaluating the learned feature coefficients. Agarwal and colleagues recently introduced neural additive models (NAMs), which combine features of deep learning and linear models [25]. NAMs compute predictions as a linear combination of outputs from independent deep learning models, each tuned to one input feature, resulting in the levels of explainability much appreciated in linear models without compromising on accuracy.

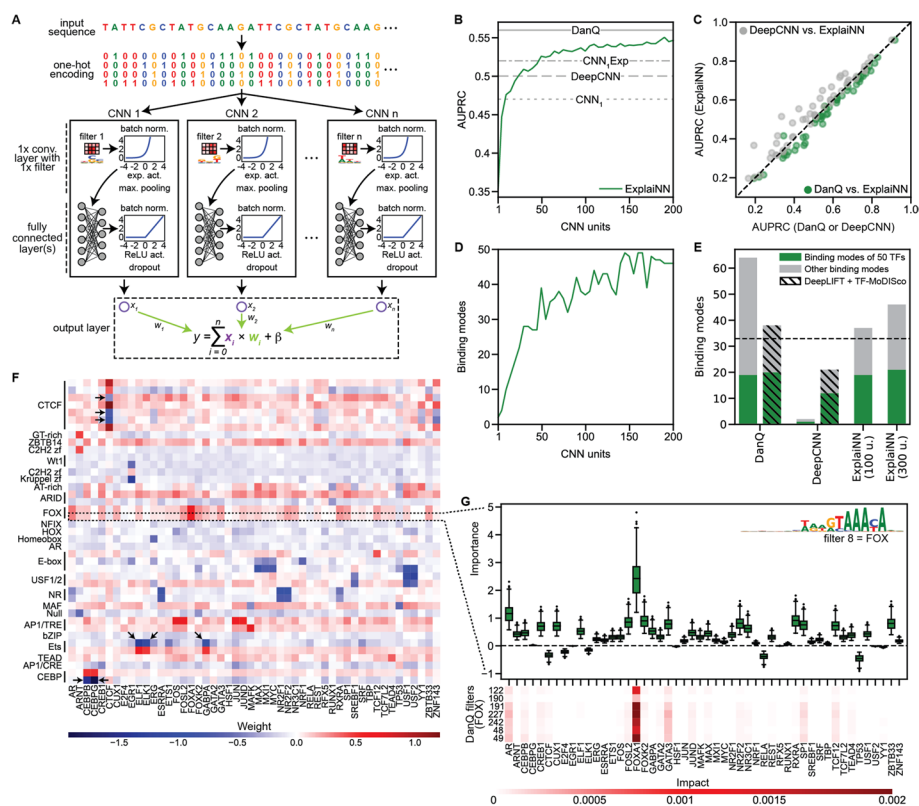
In this study, we present ExplaiNN (explainable neural networks), a fully interpretable and transparent, sequence-based deep learning model for genomic tasks inspired by NAMs. We evaluate ExplaiNN on different tasks, demonstrating that it performs as well as state-of-the-art models while yielding similar interpretations to more complex approaches, both locally and globally, but in less time and in a more intuitive and simple manner. Next, we show that the motifs learned by the convolutional filters of ExplaiNN are equivalent to those discovered by *de novo* tools on the same data. Finally, we apply ExplaiNN as a plug-and-play platform for JASPAR PWMs and pre-trained deep learning models with which to interpret genomic data, such as by distinguishing the individual roles of TFs with highly similar DNA-binding specificities.

## Results

### ExplaiNN is a glass box deep learning model for genomics

ExplaiNN is a deep learning approach for genomic tasks trained on one-hot encoded sequences inspired by NAMs [25]. Predictions are computed as a linear combination of multiple independent CNNs (hereafter referred to as units), each consisting of one

convolutional layer with a single filter followed by exponential activation, which has been seen to improve the motif representations learned by CNN filters [26], and two fully connected layers (Fig. 1 (A)). ExplainNN provides global interpretability, as the filter



**Fig. 1** The ExplainNN model and its application to TFBS prediction in OCRs. **A** ExplainNN takes as input one-hot encoded DNA sequences. Architecturally, it is composed of multiple independent CNNs (i.e., units), each of which comprising one convolutional layer with a single filter, batch normalization, exponential activation, and max pooling, and two fully connected layers with batch normalization, ReLU activation, and dropout. The final linear layer of ExplainNN (i.e., the output) combines the outputs from each unit (denoted here using  $X_s$ ). **B** Performances (AUPRC; y-axis) of ExplainNN models trained using increasing numbers of units (x-axis) on predicting the binding of 50 TFs in OCRs (green line). The performances of DanQ [7], a deep CNN with 3 convolutional layers (i.e., DeepCNN), and two shallow CNNs with 1 convolutional layer (i.e.,  $CNN_1$  and  $CNN_1Exp$  featuring an exponential activation function instead of ReLU) are provided as baselines (gray lines). **C** Pairwise comparison of the individual performances (AUPRC) of the 50 TFs from the previous dataset between the ExplainNN model trained using 100 units and either DanQ (green dots) or the DeepCNN (gray dots). **D** Number of binding modes (y-axis) detected with ExplainNN using increasing numbers of units (x-axis). **E** Number of binding modes detected with DanQ, the DeepCNN, and ExplainNN trained using 100 and 300 units on the previous dataset using either filter visualization (no pinstripes) or TF-ModISco [23] clustering on DeepLIFT [22] attribution scores (pinstripes). The 50 TFs in the dataset are represented by 33 unique binding modes (dashed line), some of which are detected using different combinations of models and interpretive approaches (green); other detected binding modes (i.e., different from those 33) are shown in gray. **F** Heatmap of the final linear layer weights of the ExplainNN model trained using 100 units, with rows representing units with assigned biological annotations based on their Tomtom [27] similarity to known TF profiles from the JASPAR database [18] and columns representing the 50 TFs predicted by the model. More than one filter can learn the same TF motif representation, but some may not contribute to the model's predictions (black arrows). **G** (top) Visualization of importance scores for a unit annotated as FOXA1 from the ExplainNN model trained using 100 units. This unit contributes the most to the prediction of FOXA1 binding. (bottom) Filter nullification analysis for the seven DanQ filters annotated as FOX TFs. The results are consistent with the unit importance scores. AUPRC, area under the precision-recall curve; CNN, convolutional neural network; OCR, open chromatin region; ReLU, rectified linear unit; TF, transcription factor; TFBS, TF binding site

of each unit can be mapped to a TF profile from JASPAR using Tomtom [27] (see the “Methods” section), thus assigning a biological annotation to that unit. The main difference between ExplainNN and a standard CNN is that the learned weights of each unit from the final linear layer (i.e., the coefficients) are interpretable, similar to a linear model. In addition, ExplainNN can also provide local interpretability by multiplying the output of each unit by the weight of that unit for each input sequence (hereafter referred to as unit importance scores; see the “Methods” section).

As a proof of concept, we trained ExplainNN to predict the binding of 50 TFs to open chromatin regions (OCRs) from a reference dataset describing the binding of 163 TFs to >1.8 M 200-bp long OCRs that we repurposed for this study [28] (see the “Methods” section). A key hyperparameter in ExplainNN is the number of independent units to be used. To assess the impact of this hyperparameter on model performance, we trained multiple ExplainNN models using increasing numbers of units (from 1 to 200). As expected, the performance of ExplainNN improved with the number of units used, plateauing at around 100 units (Fig. 1 (B)). For comparison, we evaluated four additional models on the same dataset (see the “Methods” section): a CNN with one convolutional layer ( $CNN_1$ ); a  $CNN_1$  with exponential activation function ( $CNN_1Exp$ ); a deep CNN with three convolutional layers (DeepCNN); and DanQ [7], a hybrid deep learning model with convolutional and recurrent layers. ExplainNN outperformed all three CNN models as measured by average area under the precision-recall curve (AUPRC) and, when using more than 100 units, nearly reached the performance of the more complex DanQ (Fig. 1 (B)). Focusing on the ExplainNN model trained with 100 units, it outperformed the DeepCNN model for most TFs, performing only slightly worse than DanQ (Fig. 1 (C)).

Next, we visualized the filters of each ExplainNN model and assigned them TF binding modes from JASPAR, which we defined based on the hierarchically clustered groups of DNA-binding profiles included in the database (Additional file 1; see the “Methods” section). As with performance, the number of binding modes recovered by the model increased with the number of units used (Fig. 1 (D)). For comparison, we provide the number of binding modes recovered by the DeepCNN and DanQ when applying filter visualization, as we did with ExplainNN, or when using TF-MoDISco [23] clustering on DeepLIFT [22] attribution scores (see the “Methods” section). Out of 33 different binding modes associated with the set of 50 TFs analyzed, ExplainNN models trained with 100 and 300 units recovered 19 and 21, respectively, which is a similar number as DanQ (19 when applying filter visualization and 20 when using DeepLIFT and TF-MoDISco) and greater than obtained for the DeepCNN model (Fig. 1 (E)).

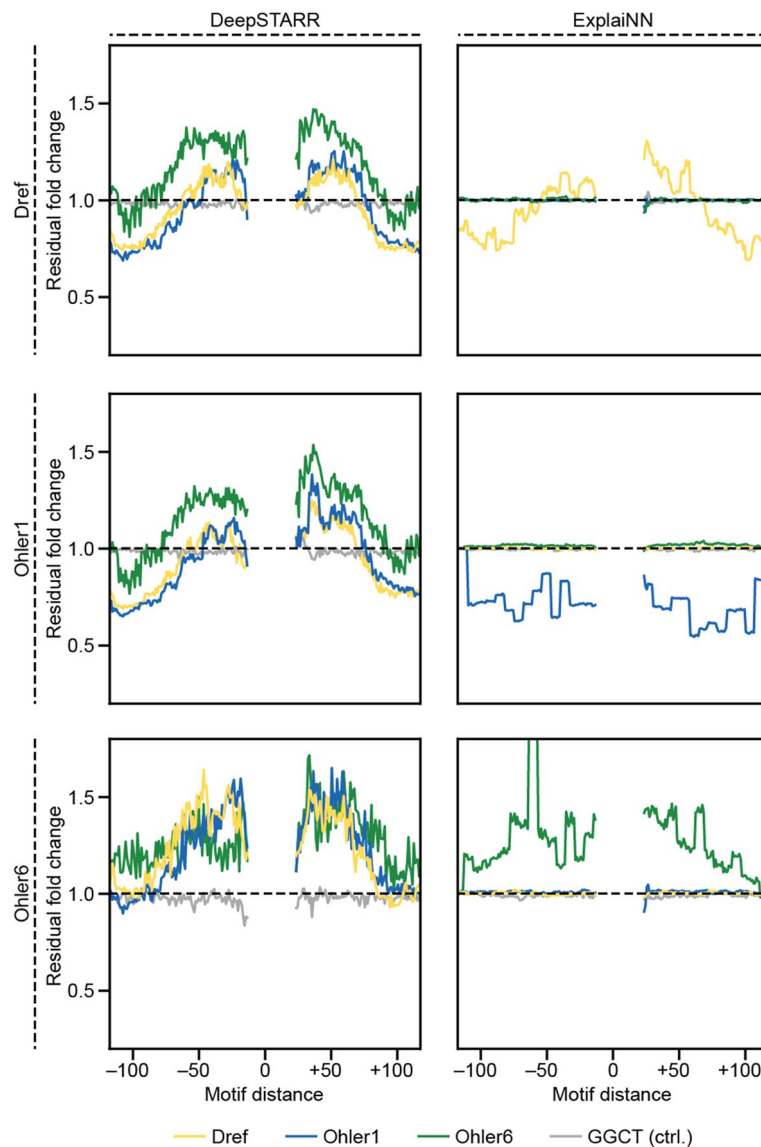
An advantageous feature of ExplainNN is that one can readily visualize the final linear layer weights for global interpretation purposes (Fig. 1 (F); see the “Methods” section). For example, units with filters annotated as FOX motifs had high positive weights for predicting the FOXA1 class. Similarly, CEBP-, CTCF-, and Ets-like units had high positive weights associated with predicting the classes of CEBP factors, CTCF, and Ets family members, respectively. However, some units had negative weights for predicting the class of their annotated TF (Fig. 1 (F); highlighted with arrows). To delve further into the contribution of each unit to the prediction of each class, we computed unit importance scores (see the “Methods” section). Visualization of the importance scores of a FOX-like

unit in a heatmap confirmed its importance for predicting the FOXA1 and AR classes (Fig. 1 (G); top panel), in agreement with the observation that FOXA1 helps shape AR signaling in prostate cells [29]. Visualizing unit importance scores also revealed an aspect as to why several CTCF-like units had negative weights for predicting the CTCF class: the importances of these units for the CTCF class were negligible, suggesting that the model was not using them to make predictions for that class (Additional file 2: Fig. S1). The same was true for units annotated as CEBP and Ets with negative weights for those classes (Additional file 2: Fig. S1). Finally, we calculated the impact from nullifying the FOX-like filters in the DanQ model one at a time (Fig. 1 (G); bottom panel; see the “Methods” section) and, as expected, the impact scores from the filter nullification process were consistent with the unit importance scores.

Taken together, these analyses demonstrated that ExplaiNN performs comparably to more complex models, at least for TF binding prediction. In addition, ExplaiNN provided local and global interpretation quickly and readily compared to using DeepLIFT followed by TF-MoDISco or filter visualization and nullification.

#### **ExplaiNN captures homotypic but not heterotypic cooperative interactions**

Given the architecture of ExplaiNN, in which each unit filter is independent of the rest, we expected that it would not be able to learn nonlinear interactions such as heterotypic cooperativity between pairs of TF motifs. In contrast, DeepSTARR is a CNN trained on STARR-seq data to predict the activities of *Drosophila* developmental and housekeeping enhancers that can capture these types of interactions [30]. For each pair of motifs analyzed, the authors of DeepSTARR performed a distance dependence analysis by sliding one motif along randomly generated sequences in the center of which they embedded the second motif. Accounting for additive effects, they observed that the output of DeepSTARR increased when the two motifs were proximal, suggesting that the model had learned heterotypic cooperative interactions. To check whether this was also the case for ExplaiNN, we trained it on the same dataset and compared its performance to DeepSTARR by calculating the Pearson correlation coefficient (PCC) between predicted and actual enhancer activities (see the “Methods” section). ExplaiNN performed worse than DeepSTARR on both developmental (PCC=0.61 vs. 0.68) and housekeeping (PCC=0.71 vs. 0.74) enhancers, which we attributed to the greater presence of nonlinear interactions in this particular dataset. Next, following the specifications from DeepSTARR, we performed a distance dependence analysis between the housekeeping TFs Dref, Ohler1, and Ohler6 (see the “Methods” section). As a negative control, we slid the 5-mer GGGCT. As expected, DeepSTARR was able to learn distance dependencies between the three motifs (Fig. 2). This was not the case for ExplaiNN: during the analysis, the resulting model outputs using the three motifs were similar to sliding GGGCT (Fig. 2). While these results were expected given the architecture of ExplaiNN, we wondered if the fully connected layers would be able to capture homotypic cooperative interactions between pairs of instances of the same TF motif. We repeated the previous distance-dependent analysis but this time sliding a second instance of the motif that had been fixed at the center of the sequences. We observed that both ExplaiNN and DeepSTARR consistently captured homotypic cooperativity for Dref and Ohler6; however, for Ohler1, the distance dependencies obtained from the two methods were



**Fig. 2** Exploring the limitations of ExplainNN in capturing nonlinear interactions. Cooperativity (residual fold change; y-axis) plotted as a function of distance (x-axis) between the motifs of the housekeeping TFs Dref (top row; yellow), Ohler1 (middle row; blue), and Ohler6 (bottom row; green) for DeepSTARR [30] (left column) and ExplainNN (right column). The 5-mer GGCT is provided as a negative control (light gray). TF, transcription factor

non-concordant (Fig. 2). To confirm whether the fully connected layers were responsible for learning these interactions, we trained a second ExplainNN model on the DeepSTARR dataset in which the fully connected layers of each unit had been replaced with a global max pooling layer, after which we repeated the distance dependencies analysis. This model without fully connected layers was unable to learn any distance dependencies: results were similar to the negative control (Additional file 2: Fig. S2). Taken together, and as we had anticipated based on its architecture, ExplainNN is not suitable for detecting heterotypic cooperative interactions; however, the fully connected layers within each unit allow homotypic cooperativity to be learned.



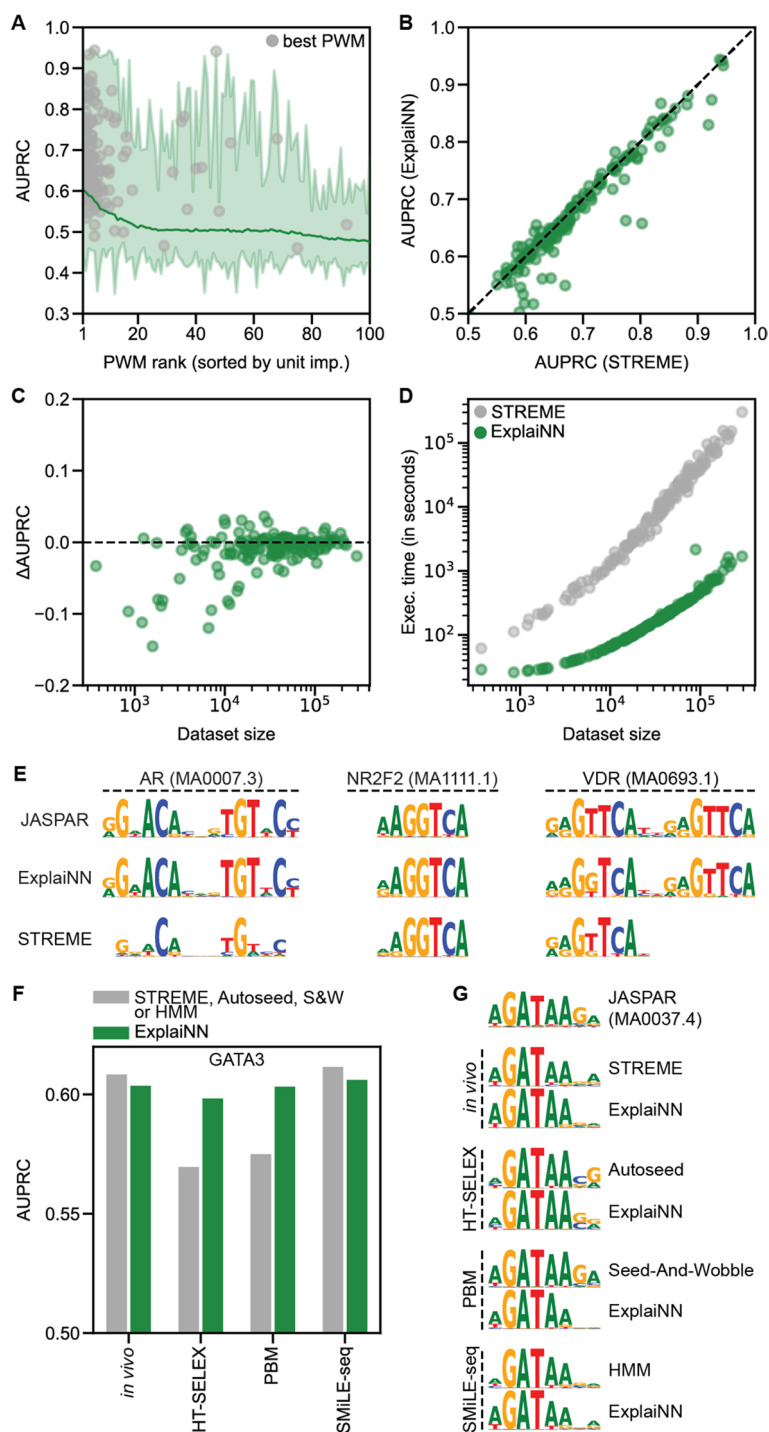
### ExplaiNN learns high-quality motifs comparable to de novo motif discovery tools

De novo motif discovery methods continue to emerge and improve [31–36]. With the dramatic escalation in the size of datasets, the execution time of these methods is increasingly a limitation. Furthermore, many de novo motif discovery methods are assay-specific, as exemplified by the DREAM5 challenge evaluation on protein binding microarray (PBM) data [37], requiring an extensive adjustment of parameters for their application to different assays. We sought to explore the capacity of ExplaiNN for efficient de novo motif discovery within a unified platform. For each of the 163 TFs from the previous dataset, we trained a model with 100 units and then visualized the filters and importance scores of each unit, resulting in 100 PWMs for each TF (see the “Methods” section). As expected, PWMs derived from visualizing filters associated with important units performed better: for 139 TFs (85.3%), the best-performing PWM was derived from the filter of one of the 10 most important units (Fig. 3(A)). Next, we applied STREME [34], a state-of-the-art method for de novo motif discovery, on the same sequences used for training the ExplaiNN models to discover 100 motifs for each TF (see the “Methods” section). Pairwise comparison of the performances of the best PWMs obtained by each method revealed subtle differences (Fig. 3(B)), although PWMs discovered by STREME were superior for TFs with small training datasets (Fig. 3(C)). Notably, the execution times exhibited by the two methods differed greatly, with ExplaiNN being > 100 times faster for TFs with large training datasets of  $\geq 50,000$  sequences (Fig. 3(D)). These differences were consistent with recent reports related to the benefits of GPU-enabled de novo motif discovery [33, 36].

The convolutional filters of CNNs are not guaranteed to learn full motif representations. Instead, CNNs are designed to learn motifs in a hierarchical manner, with the filters learning partial motifs that are assembled into full motifs in deeper layers, such as through higher-order interactions facilitated by the fully connected layers [19]. To explore the ability of ExplaiNN to learn full motif representations, we considered the case of nuclear receptors, which can bind to DNA sequences in the form of palindromes, direct repeats, or extended monomeric sites [42]. Focusing on 12 nuclear receptors from the previous analysis with at least 10,000 training sequences, ExplaiNN was able to learn motif representations faithful to the profiles of those TFs from the JASPAR database,

(See figure on next page.)

**Fig. 3** De novo motif discovery with ExplaiNN. **A** Average performances (AUPRC;  $y$ -axis) by rank ( $x$ -axis) of PWMs derived from training ExplaiNN models with 100 units on in vivo datasets of 163 TFs and then visualizing the filter of each unit (i.e., 100 PWMs per TF). The rank of each PWM is given by the importance of its unit. The gray dots indicate the rank of the best performing PWM for each TF. **B** Pairwise comparison of the individual performances (AUPRC) of the best PWMs derived for each TF from the previous dataset using ExplaiNN ( $y$ -axis) or STREME [34] ( $x$ -axis). **C** Performance difference (i.e.,  $\Delta$ AUPRC) of the previous PWMs ( $x$ -axis) is plotted with respect to the dataset size of the corresponding TF ( $x$ -axis). **D** Execution time (in seconds;  $y$ -axis) of the de novo motif discovery application of ExplaiNN (green dots) and STREME (gray dots) is plotted with respect to the dataset size of the corresponding TF ( $x$ -axis). **E** Logos derived using ExplaiNN or STREME for the nuclear receptors AR, NR2F2, and VDR from the previous dataset. For comparison, the JASPAR [18] logos for these TF profiles are shown: MA0007.3 (AR), MA1111.1 (NR2F2), and MA0693.1 (VDR). **F** Performances (AUPRC;  $y$ -axis) of PWMs derived from different experimental assay datasets related to the TF GATA3 by different methods ( $x$ -axis), including ExplaiNN (green bars) and four assay-specific methods [34, 38–40] (gray bars). **G** GATA3 logos derived from the dataset of each experimental assay using ExplaiNN or the assay-specific method. The JASPAR logo for this TF profile (MA0037.4), derived by applying RSAT [35] on the mouse Gata3 ChIP-seq data from ReMap [41], is shown at the top. AUPRC, area under the precision-recall curve; HMM, hidden Markov model; PBM, protein binding microarray; PWM, position weight matrix; S&W, Seed-And-Wobble; TF, transcription factor



**Fig. 3** (See legend on previous page.)

while STREME often reported partial motifs consistent with monomeric sites (Fig. 3(E) and Additional file 2: Fig. S3). These results are consistent with two recent reports in which we applied ExplainNN to uncover the underlying mechanism of action of two different mutations in the TF IRF4 associated with immunodeficiency in children [43] and lymphoma [44]. In both cases, ExplainNN was able to discover DNA motifs that were

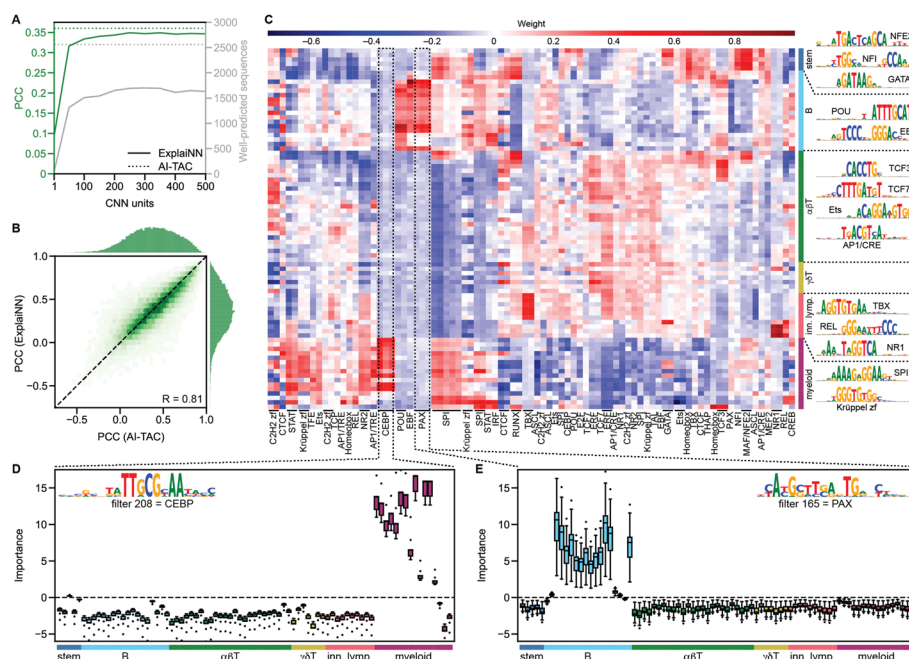


shown to impair the genome-wide binding landscape of the mutant IRF4s and the transcriptomes of the mutant cells.

Motivated by the success of ExplainNN in discovering motifs in *in vivo* data, and to demonstrate its capacity on data from various assays, we benchmarked ExplainNN against other *de novo* methods, but this time using *in vitro* data (see the “Methods” section). We downloaded publicly available HT-SELEX [45], PBM [46], and SMiLE-seq [38] data for the TF GATA3, as well as the PWMs discovered in these datasets by three assay-specific *de novo* motif discovery methods: Autoseed [39] (HT-SELEX), Seed-and-Wobble [40] (PBM), and a hidden Markov model (HMM) approach (SMiLE-seq). The performance of the best PWMs derived with ExplainNN was consistent regardless of the type of assay (Fig. 3(F)), as were their logos (Fig. 3(G)). Focusing on specific assays, the capacity of ExplainNN to discover *de novo* motifs in the *in vivo* and SMiLE-seq data, as measured by the performance of the best PWMs derived, was comparable to that of STREME and the HMM-based method, respectively, while outperforming Autoseed and Seed-and-Wobble in their corresponding assays (Fig. 3(F)). All GATA3 logos were very similar to each other and were also similar to the logo from JASPAR for this TF profile (MA0037.4), derived originally by applying RSAT [35] on mouse Gata3 ChIP-seq data from ReMap [41] (Fig. 3(G)). Taken together, this supports the potential of ExplainNN as a fast, universal method for *de novo* motif discovery.

#### **ExplainNN recapitulates the cis-regulatory lexicon of immune cell differentiation**

To further explore the capabilities of ExplainNN on distinct data types, we compared its performance against AI-TAC in predicting chromatin accessibility of sequences across 81 immune cell types from 6 different lineages [9]. We started with an exploratory analysis to determine the optimal number of units to train ExplainNN. Saturation in model performance by means of average PCC between predicted and actual ATAC signals, as well as in the number of well-predicted sequences, occurred at ~250 units (Fig. 4 (A); see the “Methods” section); however, we decided to use 300 units, which is the same number of convolutional filters used in the first layer of AI-TAC. The performance of ExplainNN by means of average PCC was comparable to that of AI-TAC (1–2% difference; Fig. 4 (A)), and the PCCs of individual sequences correlated well between the two models (Fig. 4 (B)); however, AI-TAC correctly predicted more sequences (Fig. 4 (A)). Next, we visualized both the filters and weights of each unit of the ExplainNN model, identifying the same lineage-specific TF motifs reported in AI-TAC without having to undergo the computationally intensive process of filter nullification (Fig. 4 (C)): NFE2, NFI, and GATA (in stem cells); POU, EBF, and PAX (in B cells); TCF3, TCF7, Ets, and AP1 (in T cells); NR1 (nuclear receptor type 1), TBX, and REL (in innate lymphoid cells); and SPI, Krüppel zinc fingers, and CEBP (in myeloid cells). Moreover, we visualized the importances of each unit to understand their influence on the model’s predictions. For instance, CEBP- and PAX-like units were important for predicting accessibility in most cell types of the myeloid and B lineages, respectively (Fig. 4 (D and E)). Taken together, using ExplainNN, we replicated the results of AI-TAC without the need to apply complex and time-consuming interpretation techniques by simply visualizing the weights and importances of each unit.



**Fig. 4** Application of ExplainNN in predicting chromatin accessibility in the mouse immune system. **A** Performances (average PCC; y-axis; green) and number of well-predicted sequences (secondary y-axis; gray) for ExplainNN models (solid lines) with increasing numbers of units (x-axis) and AI-TAC (dashed lines) trained on OCRs in 81 immune cell types from 6 different lineages [9]. **B** Pairwise comparison of the individual performances (PCC) of the OCRs from the previous dataset between the ExplainNN model trained using 300 units (y-axis) and AI-TAC (x-axis). The Pearson correlation coefficient (R) of the individual OCR performances between the two methods is shown at the lower right corner. **C** Heatmap of the final linear layer weights of the ExplainNN model trained using 300 units, with columns representing units with assigned biological annotations based on their Tomtom [27] similarity to known TF profiles from the JASPAR database [18] and rows representing the 81 immune cell types colored by lineage: stem cells (navy blue), B cells (turquoise), alpha/beta (forest green) and gamma/delta T cells (olive green), innate lymphoid cells (pink), and myeloid cells (purple). The logos derived from visualizing the filters of selected lineage-specific units are shown at the right. **D, E** Visualizations of importance scores colored by lineage for two units from the previous model annotated as CEBP and PAX, revealing their importance to the prediction of chromatin accessibilities in myeloid and B cell lineages, respectively. The logos of the filters of these units are shown at the top. OCR, open chromatin region; PCC, Pearson correlation coefficient; TF, transcription factor

**ExplainNN is suitable for the analysis of single-cell chromatin accessibility data**

Single-cell (sc) sequencing methods enable profiling of a wide range of genomic information in individual cells (reviewed in [47]), including chromatin accessibility. To explore the utility of ExplainNN for deciphering *cis*-regulatory properties from granular sc data, we reanalyzed a recent scATAC-seq dataset cataloging 228,873 OCRs across 15,298 human pancreatic islet cells that were grouped into 12 clusters based on their accessibility profiles [48]. We trained an ExplainNN model with 400 units (i.e., the number of units at which model performance plateaued) on the sc data to predict the activity of each OCR across the 12 clusters and visualized both the filters and weights of each unit. We observed that the weights of some units exhibited cell type-specific patterns that had also been found using chromVAR [49] in the original study (Fig. 5(A)). For example, PDX-like units had high positive weights for beta and delta cells, MAF-like units for alpha and beta cells, HNF1-like units for alpha and gamma but not for ductal cells, and FOX-like units for alpha, beta, and gamma cells. However, there were also differences:

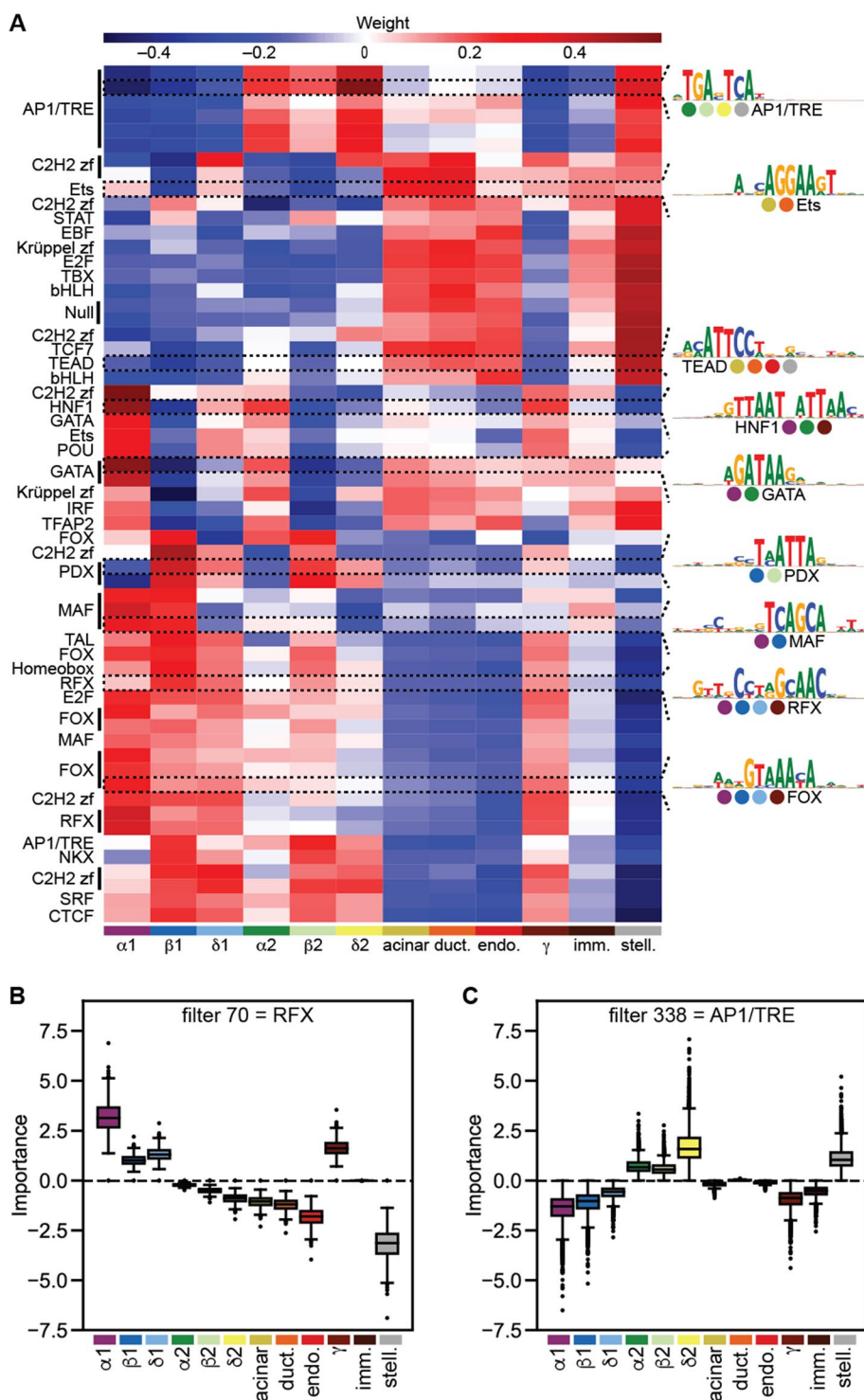
some units did not exhibit high positive weights in expected cell types (e.g., the highlighted Ets-like unit in endothelial cells), while others exhibited cell type-specific patterns not reported in the original study (Fig. 5(A)). Next, we visualized the importances of each unit for their contributions to the physiological stratification of the cells, finding that RFX-like units were important for predicting OCRs in hormone-high (i.e., alpha, beta, and delta type 1 cells) but not in hormone-low cells (i.e., alpha, beta, and delta type 2 cells) (Fig. 5(B)). It was the opposite for AP1-like units: they were important for hormone-low but not hormone-high cells (Fig. 5(C)). Both results were in agreement with the original study. Taken together, ExplainNN was able to reproduce and expand on the results from the original study while demonstrating its utility for the analysis and interpretation of scATAC-seq data.

### ExplainNN as a plug-and-play platform for TF motifs and deep learning models

Given that ExplainNN models can be conceptualized as a PWM scanning layer feeding into fully connected layers, we reasoned that initializing the weights of each unit filter with a JASPAR profile would facilitate the interpretation process because the biological annotations of the units would be known beforehand. To confirm this, we trained ExplainNN models with increasing numbers of units (from 300 to 1492) on the AI-TAC dataset in which the filter weights of each unit had been initialized with JASPAR profiles (Fig. 6(A); see the “Methods” section). During training, the filter weights were frozen to prevent them from being refined (i.e., the models were only allowed to learn the weights of the fully connected and final linear layers). These JASPAR-initialized, frozen models, even the largest one with 1492 units, performed much worse than both AI-TAC and an ExplainNN model with 300 units trained from scratch (Fig. 6(B)). Still, the importances of some units were informative. For example, a unit whose filter weights had been initialized with the profile of Lef1 (MA0768.2) was important for predicting accessibility in T cells (Additional file 2: Fig. S4), in agreement with the role of this TF in establishing T cell identity [50]. We attributed the overall poor performance of these frozen models to the fact that many JASPAR profiles used to initialize the filters might be from TFs irrelevant for immune cells and, therefore, a refinement process would be required to allow them to better resemble the motifs of relevant TFs. Indeed, unfreezing the filter weights during training (and thus allowing their refinement) improved the performance of the model, approaching that of AI-TAC (Fig. 6(B)). Moreover, as a consequence of this refinement, >35% of the filters underwent substantial changes; their visualization as PWMs revealed that they had become different from the original JASPAR profiles

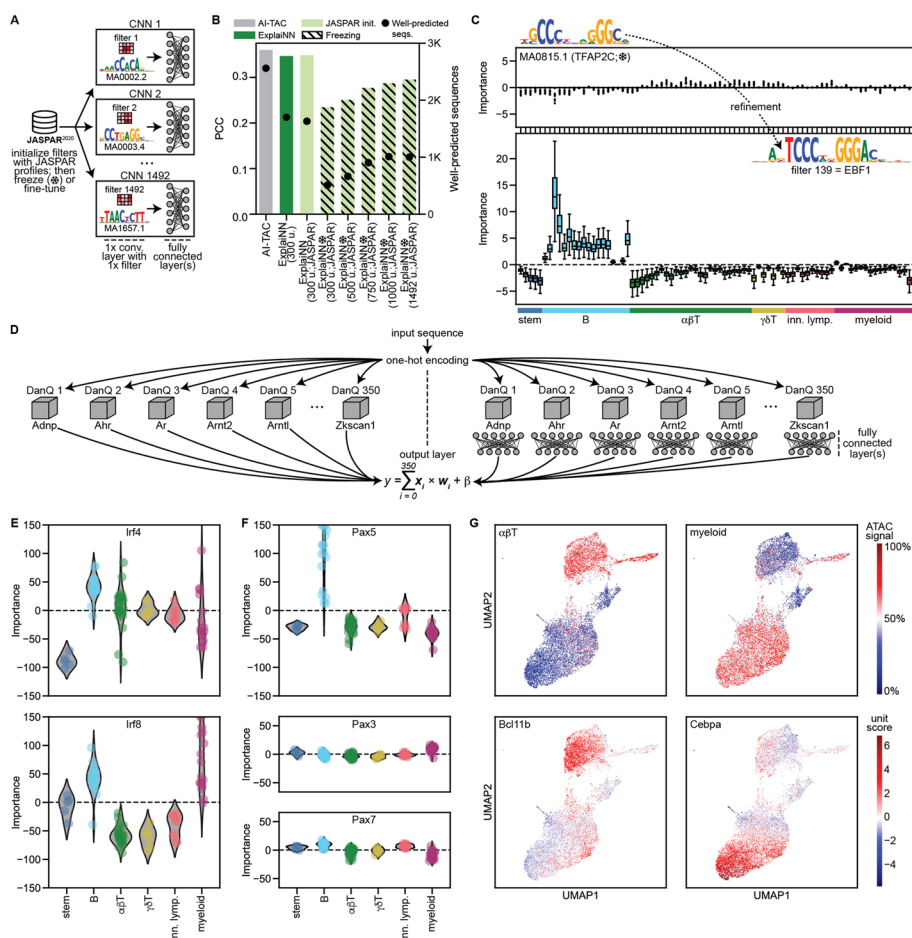
(See figure on next page.)

**Fig. 5** ExplainNN analysis and interpretation of scATAC-seq data of human pancreatic islets. **A** Heatmap of the final linear layer weights of an ExplainNN model with 400 units trained on OCRs from human pancreatic islet sc data [48], with rows representing units with assigned biological annotations based on their Tomtom [27] similarity to known TF profiles from the JASPAR database [18] and columns representing the 12 clusters of cells based on their accessibility profiles and colored by cell type: alpha type 1 cells (purple), beta type 1 cells (navy blue), delta type 1 cells (turquoise), alpha type 2 cells (forest green), beta type 2 cells (light green), delta type 2 cells (yellow), acinar cells (olive green), ductal cell (orange), endothelial cells (red), gamma cells (brown), immune cells (dark brown), and stellate cells (gray). The logos derived from visualizing the filters of selected units are shown at the right. **B, C** Visualizations of importance scores colored by cell cluster for two units from the previous model that were annotated as RFX and AP1 (TRE site), revealing their respective importance to the prediction of chromatin accessibilities in hormone-high and hormone-low cells. OCR, open chromatin region; scATAC-seq, single-cell ATAC-seq; TF, transcription factor



**Fig. 5** (See legend on previous page.)

used for their initialization (Tomtom [27]  $q$ -value > 0.05). For instance, a unit whose filter weights had been initialized with the profile of TFAP2C (MA0815.1), and whose importance across the different immune lineages when freezing the filter weights during



**Fig. 6** Initializing ExplainINN with JASP profiles and DanQ models. **A** Schematic representation of the proposed ExplainINN model in which the convolutional filters of each unit are initialized with JASP profiles [18]. **B** Performance (average PCC; y-axis) and number of well-predicted sequences (secondary y-axis; black dots) for different models trained on OCRs from the AI-TAC dataset [9]: one ExplainINN model with 300 units that was trained normally (green), as well as different ExplainINN models with increasing numbers of units (300, 500, 750, 1000, and 1492) that were initialized with JASP profiles (light green), with (pinstripes) and without (no pinstripes) freezing the weights of the filters during training. The performance of AI-TAC is provided as baseline (gray). **C** Refinement example for a filter whose weights had been initialized with the JASP profile of TFAP2C (MA0815.1). Freezing the filter weights during training was detrimental (top): the importance scores of that filter’s unit across all immune lineages were null. However, if filter weight refinement was allowed during training (i.e., no freezing; bottom), that filter was modified until it resembled the motif of EBF1, an important TF for maintaining B cell identity [51], and that same unit became important for predicting accessibility in B cells. **D** Schematic representation of the proposed transfer learning strategy in which ExplainINN units are replaced with pretrained DanQ [7] models that can be followed (right) or not (left) by fully connected layers. **E, F** We used the transfer learning strategy on the left to train one ExplainINN model with 350 units on the AI-TAC dataset in which the units had been replaced with 350 different pretrained DanQ models, each predicting the binding of a single TF to the mouse genome. During the training process of the ExplainINN model, the DanQ models were frozen (i.e., their weights were not modified). Unit importance scores of DanQ models belonging to different members of the IRF and PAX TF families are shown. **G** We repeated the same process but using the transfer learning strategy on the right: each unit was replaced with a pretrained DanQ model but adding two fully connected layers after each model. Then, we applied UMAP [52] to cluster the OCRs based on their unit outputs. (top) UMAP clusters display cell-type specificity (e.g., alpha/beta T and myeloid cells). (bottom) The outputs of the Bcl11b and Cebpa DanQ model units strongly agree with their biologically relevant clusters. OCR, open chromatin region; PCC, Pearson correlation coefficient; TF, transcription factor; UMAP, uniform manifold approximation and projection



training was null, became important for predicting accessibility in B cells; its filter was refined to such an extent that when visualized as a PWM it resembled the motif of EBF1, an important TF for maintaining B cell identity [51] (Fig. 6(C)).

Another limitation that arises during model interpretation is the inability to distinguish between TFs from the same structural family (i.e., TFs sharing the same class of DNA-binding domain), as they often have highly similar DNA-binding specificities [53]. In order to generate individual units with the resolution to highlight the importance of different paralogs, we implemented a transfer learning strategy (Fig. 6(D); see the “Methods” section): We pretrained 350 single-task DanQ models, each predicting the binding of a single TF to the mouse genome. The AUPRCs of the pretrained DanQ models ranged from 0.51 (E4f1) to 0.96 (Snai2), with a median of 0.78 across all models (Additional file 3). Then, we initialized an ExplainNN model with 350 units in which we replaced the layers of each unit with one of the pretrained DanQ models. Similar to NAMs, we initialized a second ExplainNN model in which we added two fully connected layers after each DanQ model. Next, we fine-tuned both ExplainNN models on the AI-TAC dataset, freezing the pretrained DanQ models (i.e., their weights were not modified). The performance of the fine-tuned models nearly reached that of AI-TAC (PCC = 0.335 and 0.343 vs. 0.360 for AI-TAC). Next, we visualized the importances of each DanQ model unit in the ExplainNN model without fully connected layers, allowing us to disambiguate the contribution of individual TF family members (see the “Methods” section). For instance, the Irf4 and Irf8 units were important for predicting accessibility in different immune lineages (Fig. 6(E)), in agreement with the distinct roles of these TFs: Irf4 regulates B, T, myeloid, and dendritic cell differentiation [54], while Irf8 regulates B and myeloid cell lineages [55]. Similarly, the Pax5 unit was important for predicting accessibility in B cells, consistent with the role of this TF in establishing B lineage identity and function [56]; in contrast, the importances of the Pax3 and Pax7 units across immune lineages were negligible, in agreement with their role in regulating myogenesis [57] (Fig. 6(F)). Finally, we applied UMAP [52] to cluster the sequences based on their unit outputs in the second ExplainNN model, resulting in three main clusters that were associated with ATAC signals in alpha/beta T, myeloid, and B cells (Fig. 6(G) and Additional file 2: Fig. S5). The outputs of some units were in strong agreement with a biologically relevant cluster. For example, the Bcl11b unit outputs were confined within the boundaries of the alpha/beta T cell cluster, consistent with its role in the differentiation and survival of these lymphocytes [58], the Cebpa unit outputs within the myeloid cluster, in agreement with its role in myeloid differentiation [59], and the Ebf1 unit outputs within the B cluster (Fig. 6(G) and Additional file 2: Fig. S5). Taken together, replacing the units of ExplainNN with pretrained high-resolution TF binding models achieved performance levels comparable to state-of-the-art deep learning models while providing the additional value of gaining insights into the roles of individual TFs. This flexibility of incorporating different components into ExplainNN offers the potential for increased performance while retaining interpretability.



## Discussion

ExplaiNN is an explainable deep learning approach designed for straightforward discovery of genomic features contributing to model performance for a broad range of predictive tasks related to DNA sequence data. Inspired by the recently introduced NAMs, the architecture of ExplaiNN is based on multiple, simple, independent CNN units that recognize sequence patterns. The outputs of these units are subsequently combined in a similar manner to classical regression analysis techniques. We benchmarked ExplaiNN on diverse tasks, including sequence-based prediction of TF binding, chromatin accessibility, both from bulk and single-cell data, enhancer activity, and de novo motif discovery. Through this array of diverse applications, ExplaiNN performed comparably to leading methods specialized for the independent tasks.

While deep learning models have been successful at revealing novel biological insights from high-throughput sequencing data, outperforming traditional PWM-based approaches such as motif enrichment analysis, their adoption in the genomics field remains largely in the hands of deep learning experts. This is in part a consequence of their perceived opaque nature and lack of consensus on best practices for their interpretation, making it difficult for the non-specialist to apply these methods routinely. There is a growing need for deep learning models that are more interpretable and easy to use [14, 15]. ExplaiNN is one such approach. Architecturally, ExplaiNN models are constrained to only capturing homotypic cooperativity, excluding heterotypic interactions between pairs of motifs. It may therefore be surprising to some that in practice ExplaiNN performs so well across a wide range of predictive tasks. One reason behind this could be that excluding higher-order interactions acts as a regularizer leading to more accurate predictions than multi-layer CNN models on simple tasks such as TF binding prediction. Another reason could be that the presence of higher-order interactions in the genomics datasets explored could be marginal, in agreement with recent studies suggesting that *cis*-regulatory logic may be less complex than some believe [60, 61]. Alternatively, these kinds of complex interactions may be highly specific to individual enhancers, and thus not detectable by models generalized for analysis across the genome. As exemplified by the DeepSTARR analysis in this report, comparisons between ExplaiNN and non-linear models can enable researchers to assess the benefit from allowing non-linear interactions. In such cases where there is not substantial benefit, the use of the simpler, readily interpretable method may be preferred.

The emergence of new methods for regulatory sequence analysis is a recurrent process. An important early approach to identifying recurring subsequences associated with TFBSs used the expectation–maximization algorithm to learn a PWM from a set of unaligned regulatory sequences [62]. These individually learned PWMs could then be used in a linear model to make predictions regarding various properties of a given sequence [63]. ExplaiNN could be perceived as a direct extension of this approach, wherein its convolutional filters learn a large amount of PWMs simultaneously; however, unless the fully connected layers of each unit were replaced with a global maximum pooling layer, their activation values would not be equivalent to those of independent PWM scans. The fully connected layers might transform the activation values into more expressive and informative scores leading to better performance. As we have shown, this expressiveness is likely achieved by learning relevant motif positional information.

As shown throughout this work, ExplaiNN has the potential to be applied to multiple problems. This flexibility may motivate users to become proficient with the system as they will be able to work across or within topics using the same tool. As presented here for de novo motif discovery, the flexibility to apply ExplaiNN to both ChIP-seq data and multiple special classes of in vitro TF binding data has the promise to reduce the number of methods that must be mastered.

There remains an ample variety of potential directions for further development of ExplaiNN, split between enhanced technical capabilities and user empowerment. It should be possible to integrate a more complex model architecture within the system to account for the “missing complexity” (e.g., non-linearity). As implemented, ExplaiNN cannot address analyses requiring large receptive fields (e.g., Akita [5], Xpresso [11], Basenji [64], or Enformer [12]), but it should be possible to implement such capacity within the system. Incorporating more advanced visualization techniques and data processing tools would ease the adoption of ExplaiNN by scientists working on applied problems who currently rely on heritage methods for regulatory sequence analysis.

The software for using ExplaiNN is provided open-source in a well-documented repository, and it is our hope that it will lead to or inspire widespread use of interpretable deep learning methods.

## Conclusions

Sequence-based deep learning models perform well on genomics-related tasks but can be challenging to interpret. ExplaiNN bridges this gap by using an interpretable neural network design, providing both global (cell state level) and local (individual sequence level) biological insights into the data without compromising on performance, thus allowing researchers to accelerate their genomic analyses.

## Methods

### Model architectures

Unless otherwise specified, each ExplaiNN unit consisted of:

- 1st convolutional layer with 1 filter ( $19 \times 4$ ), batch normalization, exponential activation, and max pooling (kernel size = 7, stride = 7)
- 1st fully connected layer with 100 nodes, batch normalization, ReLU activation, and 30% dropout
- 2nd fully connected layer with 1 node, batch normalization, and ReLU activation

In the 1st convolutional layer, exponential activation was used (as opposed to ReLU), as it has been shown to significantly improve the recovery of biologically meaningful motifs from the filters [26]. The final layer of ExplaiNN is linear: one fully connected output layer with  $n$  outputs (e.g., 50 in the initial TF binding prediction task or 81 when predicting chromatin accessibility states across immune cells).

For the  $CNN_1$  and  $CNN_1Exp$  architectures, we used the following specifications:

- 1st convolutional layer with 100 filters ( $19 \times 4$ ), batch normalization, ReLU (for  $\text{CNN}_1$ ) or exponential (for  $\text{CNN}_1\text{Exp}$ ) activation, and max pooling (kernel size = 7, stride = 7)
- 1st fully connected layer with 1000 nodes, batch normalization, ReLU activation, and 30% dropout
- 2nd fully connected layer with 1000 nodes, batch normalization, ReLU activation, and 30% dropout
- Fully connected output layer with 50 outputs

For the DeepCNN (adapted from Basset [8]):

- 1st convolutional layer with 100 filters ( $19 \times 4$ ), batch normalization, ReLU activation, and max pooling (kernel size = 3, stride = 3)
- 2nd convolutional layer with 200 filters ( $7 \times 1$ ), batch normalization, ReLU activation, and max pooling (kernel size = 3, stride = 3)
- 3rd convolutional layer with 200 filters ( $4 \times 1$ ), batch normalization, ReLU activation, and max pooling (kernel size = 3, stride = 3)
- 1st fully connected layer with 1000 nodes, batch normalization, ReLU activation and 30% dropout
- 2nd fully connected layer with 1000 nodes, batch normalization, ReLU activation, and 30% dropout
- Fully connected output layer with 50 outputs

For DanQ [7]:

- 1st convolutional layer with 320 filters ( $26 \times 4$ ), ReLU activation, 20% dropout, and max pooling (kernel size = 13, stride = 13)
- 2 bi-directional LSTM layers with a hidden state size of 320 and 50% dropout
- 1st fully connected layer with 925 nodes and ReLU activation
- Fully connected output layer with 1 or 50 outputs (depending on the task)

For DeepSTARR [30]:

- 1st convolutional layer with 256 filters ( $7 \times 4$ ), padding of size 3, batch normalization, ReLU activation function, and max pooling (kernel size = 2, stride = 2)
- 2nd convolutional layer with 60 filters ( $3 \times 1$ ), padding of size 1, batch normalization, ReLU activation function, and max pooling (kernel size = 2, stride = 2)
- 3rd convolutional layer with 60 filters ( $5 \times 1$ ), padding of size 2, batch normalization, ReLU activation function, and max pooling (kernel size = 2, stride = 2)
- 4th convolutional layer with 120 filters ( $3 \times 1$ ), padding of size 1, batch normalization, ReLU activation function, and max pooling (kernel size = 2, stride = 2)
- 1st fully connected layer with 256 nodes, batch normalization, ReLU activation function, and 40% dropout
- 2nd fully connected layer with 256 nodes, batch normalization, ReLU activation function, and 40% dropout
- Fully connected output layer with 2 outputs

The AI-TAC architecture was provided by Maslova and colleagues [9]. All architectures were implemented using the PyTorch framework [65] (version 1.9.0).

### Interpretability with ExplainNN

First, we converted the filter of each unit into a PWM by following the specifications from the AI-TAC manuscript: For each filter, we built a position frequency matrix (PFM) by aligning all 19-mers (i.e., 19 bp-long DNA sequences) activating that filter's unit by  $\geq 50\%$  of its maximum activation value in correctly predicted sequences. The resulting PFM was then transformed into a PWM by setting the background uniform nucleotide frequency to 0.25. Next, the PWM derived from each filter was mapped to one or more profiles from the vertebrate collection of the JASPAR 2020 database [18] using Tomtom [27] (version 5.3.0;  $q$ -value  $\leq 0.05$ ), which, in turn, were used to annotate that filter's unit. For instance, a unit whose filter's PWM were similar to the profiles of members from the SOX TF family would be annotated as "SOX-like."

For global interpretability, we visualized the weights of each unit in the final linear layer of the model (e.g., using a heatmap). Usually, these weights were associated with the importance of the unit for each task. In some cases, however, they were highly positive or negative for a task in which the unit was not activated by the input sequences. To overcome this limitation and obtain more fine-grained unit importances, for each unit and for each task, we computed the product of the unit's activation for each sequence with the final layer weight for that task and visualized them (e.g., using a boxplot). For visualization, for each unit, we only included the products from correctly predicted sequences activating that unit's filter by  $\geq 50\%$  of its maximum activation value. In addition, for each unit and for each task, the median of these products was used to assess the importance of that unit for that task.

### Interpretability with DeepLIFT and TF-MoDISco

For each correctly predicted sequence in the test set, we generated DeepLIFT [22] importance scores with 10 reference sequences using the Captum library [66] (version 0.4.0). We used TF-MoDISco [23] (version 0.5.14.0) with default settings to obtain motifs from DeepLIFT importance scores.

### Training, validation, and test datasets

To obtain human in vivo TF binding data, we repurposed a previously described data matrix aggregating the binding of 163 TFs to 1,817,918 200-bp-long DNase I hypersensitive sites (DHSs) in 52 cell and tissue types [28]. For each TF-DHS pair, a "1" was used to indicate that the DHS was accessible and bound by the TF (i.e., the DHS and at least one ChIP-seq peak summit of that TF from ReMap 2018 [41] overlapped), a "0" that the DHS was accessible but not bound by the TF, and a null sign (" $\emptyset$ ") that it was not accessible to the TF for binding (i.e., unresolved). For the CNN<sub>1</sub>, CNN<sub>1</sub>Exp, DeepCNN, DanQ, and ExplainNN models predicting the binding of 50 TFs, we extracted a slice of the matrix including the row vectors of the 50 TFs, removing any column vectors with unresolved elements. The resulting resolved regions for all 50 TFs were randomly split into training (80%), validation (10%), and test (10%) sets using the "train\_test\_split" function from

scikit-learn [67] (version 0.24.2; datasets were always randomly split in this way). For de novo motif discovery, since the number of bound versus unbound regions of all TFs was imbalanced, we subsampled the set of unbound regions to a 50:50 ratio for each TF while accounting for their %GC content distributions. The resulting datasets were randomly split into training (80%), validation (10%), and test (10%) sets while maintaining an equal proportion of bound and unbound regions.

Mouse in vivo TF binding data was obtained as follows: First, we downloaded non-redundant mouse ChIP-seq peaks for 350 DNA-binding TFs [68] from ReMap. We resized each ChIP-seq peak to 201 bp by extending its summit 100 bp in both directions. Furthermore, we retrieved an atlas of 1,802,604 DHS regions in the mouse genome [69], which we also resized to 201 bp around the center of each DHS. In both cases, we applied BEDTools slop [70] (version 2.30.0). We then created a set of negative sequences for each TF by subsampling non-overlapping regions from the DHS atlas while matching the %GC content distribution of its ChIP-seq peaks. The resulting sequences for each TF were randomly split into training (80%), validation (10%), and test (10%) sets while keeping the same ratio between positive and negative sequences.

PBM data was downloaded from UniPROBE [46] (Gata3; UP00032; clone ID pTH1024). Probe signal intensities were quantile normalized (QN) using the “quantile\_transform” function from scikit-learn. Probes were 60 bp long, including both the de Bruijn and linker sequences. The arrays AMADID #015681 and #016060 were used for training and validation, respectively.

HT-SELEX [45] and SMiLE-seq [38] data were retrieved from the Sequence Repository Archive (run ids: ERR1003435, ERR1003437, ERR1003439, ERR1003441, and SRR3405148) using parallel-fastq-dump (version 0.6.7). For HT-SELEX data, we treated each cycle as an independent class as in Asif and Orenstein [71], thereby removing the need for negative sequences. Reads were randomly split into training (80%) and validation (20%) sets while preserving the proportions between reads from each cycle. For SMiLE-seq data, reads were left- and right-clipped 7 and 64 bp, respectively, for a final length of 30 bp corresponding to the randomized DNA. A set of negative sequences were obtained by dinucleotide shuffling using BiasAway [72] (version 3.3.0). Sequences were randomly split into training (80%) and validation (20%) while maintaining an equal proportion between positives and negatives.

Binarised human pancreatic islet scATAC-seq data was obtained from [48]. Data was denoised by training a PeakVI [73] model with 21 latent dimensions using scvi-tools [74] (version 0.14.6). The denoised profiles were used to compute the mean accessibility of each peak in the 12 clusters identified in the original study. Peaks were resized to 600 bp around the center of each peak using BEDTools slop. Peaks from chromosome 1 were used for validation, peaks from chromosomes 10, 11, and 12 for testing, and peaks from the remaining chromosomes for training.

The AI-TAC and DeepSTARR datasets were obtained from the original publications, and we used the same data splits.

### Model training

Unless otherwise specified, we trained all models using the Adam optimizer [75] and binary cross entropy as loss function, except for models trained on PBM data, wherein

we used the mean squared error. We applied one-hot encoding to convert nucleotides into 4-element vectors (i.e., A, C, G, and T), set the learning rate to 0.003 and batch size to 100, and used an early stopping criteria to prevent overfitting when the model performance on the validation set did not improve. The training and validation sets included the reverse-complement of each sequence. Sequences with Ns were discarded.

### De novoPWMs

For each of the 163 TFs from the data matrix, we trained one ExplainNN model with 100 units using its training and validation sets and then visualized the filters and importance scores of each unit, resulting in 100 PWMs per TF. As baseline, we applied STREME [34] (version 5.3.0) on the set of training sequences of each TF to also generate 100 PWMs: We set the fraction of sequences held-out for  $p$ -value computation to 0 (option `-hofract`), the maximum length of the PWMs to 19 bp (i.e., the filter size in ExplainNN; option `-maxw`), and the number of output PWMs to 100 (option `-nmotifs`). Next, for each TF and for each of its PWMs, we evaluated the PWM performance by keeping the score of the best hit from scanning that PWM along both strands of each sequence in the test set and computing the AUPRC.

Like above, for each GATA3 *in vitro* assay (i.e., HT-SELEX, PBM, and SMiLE-seq), we trained one ExplainNN model with 100 units using the training and validation sets and visualized the filter of each unit, resulting in 100 PWMs for each assay. However, the task of the ExplainNN model trained on HT-SELEX data was multiclass classification: the model tried to predict the cycle of origin of each input read, with the expectation that the last cycles would be enriched for bound sequences of the TF. In addition, for PBM data, the task of the ExplainNN model was to infer QN intensity signals of each probe.

For each TF and for each assay, the best PWM was selected based on its performance on the corresponding validation set.

### Plug-and-play

First, we resized all 746 profiles from the JASPAR 2020 vertebrate collection to 19 bp. Then, we applied a farthest point sampling procedure to remove profile redundancy based on Tomtom similarities starting from the most dissimilar profile. This resulted in four different sets of non-redundant profiles of sizes 150, 250, 375, and 500. Next, we computed the reverse complement of each profile, doubling the number of profiles in each non-redundant set to 300, 500, 750, and 1000, and raising the number of profiles in the JASPAR 2020 vertebrate collection to 1492. We trained ExplainNN models with increasing numbers of units (300, 500, 750, 1000, and 1492) on the AI-TAC dataset in which the filter weights of each unit were initialized with those JASPAR profiles. To initialize filter weights with profiles from JASPAR, we followed the specifications from DanQ: We reformatted JASPAR profiles as PWMs using Biopython [76] (version 1.79) and then converted them to filter weights by subtracting 0.25 from the probability of each nucleotide at each PWM position. During training, nullification of gradients could be applied to freeze the filter weights.

For each of the 350 TFs from the mouse *in vivo* TF binding data, we trained a DanQ model using the training and validation sets of that TF and assessed its performance by computing the AUPRC on the test set. We then scored the sequences from the



AI-TAC dataset keeping the outputs from each model (pre- or post-sigmoid). Next, we combined the outputs of each sequence in one fully connected layer with 81 output nodes. Alternatively, each DanQ output could be embedded in its own fully connected neural network with one layer and 200 nodes, resulting in 350 processed outputs ultimately combined in one fully connected layer with 81 outputs.

UMAP clusters were obtained and plotted using the UMAP Python library [52] (version 0.5.2).

## Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-023-02985-y>.

**Additional file 1.** TF binding modes used in this study.

**Additional file 2.** Supplementary figures.

**Additional file 3.** Individual performances of 350 DanQ models trained on mouse in vivo binding data in predicting DNA binding of their respective TFs.

**Additional file 4.** Datasets used in this study.

## Acknowledgements

We thank the members of the Mostafavi and Wasserman labs for providing useful feedback. We thank Dora Pak and Jonathan Chang for administrative and IT support, respectively. We thank Bernardo P. de Almeida and Dr. Alexander Stark for providing us with the Dref, Ohler1, and Ohler6 motifs for the distance dependence analysis. We thank UBC's Advanced Research Computing (ARC) (<https://arc.ubc.ca/>) for enabling HPC.

## Review history

This article has been reviewed at a different journal, and no review history is available.

## Peer review information

Anahita Bishop was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

## Authors' contributions

MS and GN conceived the idea; OF, GN, and MS designed and performed the analyses; GN, MS and OF wrote the code; GN developed the Python library; OF and GN created the figures with input from MS; WW and SM provided advice on all aspects of the work; all authors participated in the writing and editing of the manuscript.

## Authors' Twitter handles

Gherman Novakovsky, @NovakovskyG; Oriol Fornes, @ofornes; Manu Saraswat, @manusaraswat10; Sara Mostafavi, @sara\_mostafavi; Wyeth W. Wasserman, @WyWyWa.

## Funding

GN was supported by an International Doctoral Fellowship from the University of British Columbia. OF, MS, and WWW were supported by grants from the Canadian Institutes of Health Research (PJT-162120), Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grant (RGPIN-2017-06824), and the BC Children's Hospital Foundation and Research Institute. SM acknowledges support from the Canadian Institute for Advanced Research (CIFAR). Equipment enabled by NSERC Research Tools and Instruments Grant (RTI-2020-00778) to SM and WWW.

## Availability of data and materials

ExplaiNN is open-source software distributed under the MIT license. It is available on GitHub (<https://github.com/wasserlab/ExplaiNN> [77, 78]) accompanied by the IPython notebooks that we used for the different experiments, as well as a collection of Python scripts to allow users to deploy their own ExplaiNN models. The TF binding matrices are available on GitHub as numpy arrays [79] (<https://github.com/wasserlab/TF-Binding-Matrix> [80, 81]). The filter weights from JASPAR profiles, and the code for generating them, are available on GitHub (<https://github.com/wasserlab/PWM-to-filter-weights> [78, 82]). The ChIP-seq data for the 350 mouse TFs were downloaded from ReMap [46], the PBM data for the TF Gata3 from UniPROBE [42] (UP00032; clone ID pTH1024), and the GATA3 HT-SELEX (run ids ERR1003435 [83], ERR1003437 [84], ERR1003439 [85], ERR1003441 [86]) and SMiLE-seq (run id SRR3405148 [87]) data from the Sequence Repository Archive (Additional file 4; see the "Methods" section).

## Declarations

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

**Competing interests**

The authors declare that they have no competing interests.

Received: 20 November 2022 Accepted: 12 June 2023

Published online: 27 June 2023

**References**

- Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat Methods*. 2013;10:1213–8.
- Johnson DS, Mortazavi A, Myers RM, Wold B. Genome-wide mapping of in vivo protein-DNA interactions. *Science*. 2007;316:1497–502.
- Libbrecht MW, Noble WS. Machine learning applications in genetics and genomics. *Nat Rev Genet*. 2015;16:321–32.
- Eraslan G, Avsec Ž, Gagneur J, Theis FJ. Deep learning: new computational modelling techniques for genomics. *Nat Rev Genet*. 2019;20:389–403.
- Fudenberg G, Kelley DR, Pollard KS. Predicting 3D genome folding from DNA sequence with Akita. *Nat Methods*. 2020;17:1111–7.
- Zhou J, Troyanskaya OG. Predicting effects of noncoding variants with deep learning-based sequence model. *Nat Methods*. 2015;12:931–4.
- Quang D, Xie X. DanQ: a hybrid convolutional and recurrent deep neural network for quantifying the function of DNA sequences. *Nucleic Acids Res*. 2016;44:e107.
- Kelley DR, Snoek J, Rinn JL. Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res*. 2016;26:990–9.
- Maslova A, Ramirez RN, Ma K, Schmutz H, Wang C, Fox C, et al. Deep learning of immune cell differentiation. *Proc Natl Acad Sci*. 2020;117:25655–66.
- Zhou J, Theesfeld CL, Yao K, Chen KM, Wong AK, Troyanskaya OG. Deep learning sequence-based ab initio prediction of variant effects on expression and disease risk. *Nat Genet*. 2018;50:1171–9.
- Agarwal V, Shendure J. Predicting mRNA abundance directly from genomic sequence using deep convolutional neural networks. *Cell Rep*. 2020;31:107663.
- Avsec Ž, Agarwal V, Visentin D, Ledsam JR, Grabska-Barwinska A, Taylor KR, et al. Effective gene expression prediction from sequence by integrating long-range interactions. *Nat Methods*. 2021;18:1196–203.
- Koo PK, Ploenzke M. Deep learning for inferring transcription factor binding sites. *Curr Opin Syst Biol*. 2020;19:16–23.
- Rudin C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat Mach Intell*. 2019;1:206–15.
- Novakovsky G, Dexter N, Libbrecht MW, Wasserman WW, Mostafavi S. Obtaining genetics insights from deep learning via explainable artificial intelligence. *Nat Rev Genet*. 2023;24:125–37.
- Stormo GD. DNA binding sites: representation and discovery. *Bioinformatics*. 2000;16:16–23.
- Alipanahi B, DeLong A, Weirauch MT, Frey BJ. Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning. *Nat Biotechnol*. 2015;33:831–8.
- Fornes O, Castro-Mondragon JA, Khan A, van der Lee R, Zhang X, Richmond PA, et al. JASPAR 2020: update of the open-access database of transcription factor binding profiles. *Nucleic Acids Res*. 2020;48:D87–92.
- Koo PK, Eddy SR. Representation learning of genomic sequence motifs with convolutional neural networks. *PLOS Comput Biol*. 2019;15:e1007560.
- Nair S, Shrikumar A, Schreiber J, Kundaje A. fastISM: performant in silico saturation mutagenesis for convolutional neural networks. *Bioinformatics*. 2022;38:2397–403.
- Sundararajan M, Taly A, Yan Q. Axiomatic attribution for deep networks. *arXiv*. 2017;1703.01365. Available from: <http://arxiv.org/abs/1703.01365>.
- Shrikumar A, Greenside P, Kundaje A. Learning important features through propagating activation differences. *arXiv*. 2019;1704.02685. Available from: <http://arxiv.org/abs/1704.02685>.
- Shrikumar A, Tian K, Avsec Ž, Shcherbina A, Banerjee A, Sharmin M, et al. Technical note on transcription factor motif discovery from importance scores (TF-MoDISco) version 0.5.6.5. *arXiv*. 2020;1811.00416. Available from: <http://arxiv.org/abs/1811.00416>.
- Koo PK, Ploenzke M. Interpreting deep neural networks beyond attribution methods: quantifying global importance of features. *bioRxiv*. 2020;2020.02.19.956896. Available from: <https://doi.org/10.1101/2020.02.19.956896>.
- Agarwal R, Melnick L, Frosst N, Zhang X, Lengerich B, Caruana R, et al. Neural additive models: interpretable machine learning with neural nets. *arXiv*. 2021;2004.13912. Available from: <http://arxiv.org/abs/2004.13912>.
- Koo PK, Ploenzke M. Improving representations of genomic sequence motifs in convolutional networks with exponential activations. *Nat Mach Intell*. 2021;3:258–66.
- Gupta S, Stamatoyannopoulos JA, Bailey TL, Noble WS. Quantifying similarity between motifs. *Genome Biol*. 2007;8:R24.
- Novakovsky G, Saraswat M, Fornes O, Mostafavi S, Wasserman WW. Biologically relevant transfer learning improves transcription factor binding prediction. *Genome Biol*. 2021;22:280.
- Teng M, Zhou S, Cai C, Lupien M, He HH. Pioneer of prostate cancer: past, present and the future of FOXA1. *Protein Cell*. 2021;12:29–38.
- de Almeida BP, Reiter F, Pagani M, Stark A. DeepSTARR predicts enhancer activity from DNA sequence and enables the de novo design of synthetic enhancers. *Nat Genet*. 2022;54:613–24.
- Kulakovskiy IV, Boeva VA, Favorov AV, Makeev VJ. Deep and wide digging for binding motifs in ChIP-seq data. *Bioinformatics*. 2010;26:2622–3.

32. Grau J, Posch S, Grosse I, Keilwagen J. A general approach for discriminative de novo motif discovery from high-throughput data. *Nucleic Acids Res.* 2013;41:e197.
33. Quang D, Guan Y, Parker SCJ. YAMDA: thousandfold speedup of EM-based motif discovery using deep learning libraries and GPU. *Bioinformatics.* 2018;34:3578–80.
34. Bailey TL. STREME: accurate and versatile sequence motif discovery. *Bioinformatics.* 2021;37:2834–40.
35. Santana-Garcia W, Castro-Mondragon JA, Padilla-Gálvez M, Nguyen NTT, Elizondo-Salas A, Ksouri N, et al. RSAT 2022: regulatory sequence analysis tools. *Nucleic Acids Res.* 2022;50:W670–6.
36. Chu S, Stormo G. Deep unfolded convolutional dictionary learning for motif discovery. *bioRxiv.* 2022;2022.11.06.515322. Available from: <https://doi.org/10.1101/2022.11.06.515322>.
37. Weirauch MT, Cote A, Norel R, Annala M, Zhao Y, Riley TR, et al. Evaluation of methods for modeling transcription factor sequence specificity. *Nat Biotechnol.* 2013;31:126–34.
38. Isakova A, Groux R, Imbeault M, Rainer P, Alpern D, Dainese R, et al. SMiLE-seq identifies binding motifs of single and dimeric transcription factors. *Nat Methods.* 2017;14:316–22.
39. Nitta KR, Jolma A, Yin Y, Morgunova E, Kivioja T, Akhtar J, et al. Conservation of transcription factor binding specificities across 600 million years of bilateria evolution. *eLife.* 2015;4:e04837.
40. Berger MF, Philippakis AA, Qureshi AM, He FS, Estep PW, Bulky ML. Compact, universal DNA microarrays to comprehensively determine transcription-factor binding site specificities. *Nat Biotechnol.* 2006;24:1429–35.
41. Chèneby J, Gheorghe M, Artufel M, Mathelier A, Ballester B. ReMap 2018: an updated atlas of regulatory regions from an integrative analysis of DNA-binding ChIP-seq experiments. *Nucleic Acids Res.* 2018;46:D267–75.
42. Weikum ER, Liu X, Ortlund EA. The nuclear receptor superfamily: a structural perspective. *Protein Sci.* 2018;27:1876–92.
43. IRF4 International Consortium. A multimorphic mutation in IRF4 causes human autosomal dominant combined immunodeficiency. *Sci Immunol.* 2023;8:eade7953.
44. Schlessner N, Cauchy P, Franke V, Giefing M, Fornes O, Vankadari N, et al. A new type of transcriptional reprogramming by an IRF4 mutation in lymphoma. *bioRxiv.* 2022;2022.12.29.522203. Available from: <https://doi.org/10.1101/2022.12.29.522203>.
45. Yin Y, Morgunova E, Jolma A, Kaasinen E, Sahu B, Khund-Sayeed S, et al. Impact of cytosine methylation on DNA binding specificities of human transcription factors. *Science.* 2017;356:eaaj2239.
46. Hume MA, Barrera LA, Gisselbrecht SS, Bulky ML. UniPROBE, update 2015: new tools and content for the online database of protein-binding microarray data on protein-DNA interactions. *Nucleic Acids Res.* 2015;43:D117–22.
47. Stuart T, Satija R. Integrative single-cell analysis. *Nat Rev Genet.* 2019;20:257–72.
48. Chiou J, Zeng C, Cheng Z, Han JY, Schlichting M, Miller M, et al. Single-cell chromatin accessibility identifies pancreatic islet cell type- and state-specific regulatory programs of diabetes risk. *Nat Genet.* 2021;53:455–66.
49. Schep AN, Wu B, Buenrostro JD, Greenleaf WJ. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat Methods.* 2017;14:975–8.
50. Shan Q, Li X, Chen X, Zeng Z, Zhu S, Gai K, et al. Tcf1 and Lef1 provide constant supervision to mature CD8+ T cell identity and function by organizing genomic architecture. *Nat Commun.* 2021;12:5863.
51. Nechanitzky R, Akbas D, Scherer S, Györy I, Hoyle T, Ramamoorthy S, et al. Transcription factor EBF1 is essential for the maintenance of B cell identity and prevention of alternative fates in committed cells. *Nat Immunol.* 2013;14:867–75.
52. McInnes L, Healy J, Melville J. UMAP: uniform manifold approximation and projection for dimension reduction. *arXiv.* 2018;1802.03426. Available from: <http://arxiv.org/abs/1802.03426>.
53. Ambrosini G, Vorontsov I, Penzar D, Groux R, Fornes O, Nikolaeva DD, et al. Insights gained from a comprehensive all-against-all transcription factor binding motif benchmarking study. *Genome Biol.* 2020;21:114.
54. Yanai H, Negishi H, Taniguchi T. The IRF family of transcription factors. *Oncol Immunology.* 2012;1:1376–86.
55. Wang H, Morse HC. IRF8 regulates myeloid and B lymphoid lineage diversification. *Immunol Res.* 2008;43:109.
56. Cobaleda C, Schebesta A, Delogu A, Busslinger M. Pax5: the guardian of B cell identity and function. *Nat Immunol.* 2007;8:463–70.
57. Buckingham M, Relaix F. PAX3 and PAX7 as upstream regulators of myogenesis. *Semin Cell Dev Biol.* 2015;44:115–25.
58. Wakabayashi Y, Watanabe H, Inoue J, Takeda N, Sakata J, Mishima Y, et al. Bcl11b is required for differentiation and survival of  $\alpha\beta$  T lymphocytes. *Nat Immunol.* 2003;4:533–9.
59. Pabst T, Mueller BU, Zhang P, Radoska HS, Narravula S, Schnittger S, et al. Dominant-negative mutations of CEBPA, encoding CCAAT/enhancer binding protein- $\alpha$  (C/EBP $\alpha$ ), in acute myeloid leukemia. *Nat Genet.* 2001;27:263–70.
60. Wei B, Jolma A, Sahu B, Orre LM, Zhong F, Zhu F, et al. A protein activity assay to measure global transcription factor activity reveals determinants of chromatin accessibility. *Nat Biotechnol.* 2018;36:521–9.
61. Patel ZM, Hughes TR. Global properties of regulatory sequences are predicted by transcription factor recognition mechanisms. *Genome Biol.* 2021;22:285.
62. Bailey TL, Elkan C. Unsupervised learning of multiple motifs in biopolymers using expectation maximization. *Mach Learn.* 1995;21:51–80.
63. Wasserman WW, Fickett JW. Identification of regulatory regions which confer muscle-specific gene expression. Edited by G. Von Heijne. *J Mol Biol.* 1998;278:167–81.
64. Kelley DR, Reshef YA, Bileschi M, Belanger D, McLean CY, Snoek J. Sequential regulatory activity prediction across chromosomes with convolutional neural networks. *Genome Res.* 2018;28:739–50.
65. Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: an imperative style, high-performance deep learning library. *Adv Neural Inf Process Syst.* 2019;32:8026–37.
66. Kohlikiyan N, Miglani V, Martin M, Wang E, Alsallakh B, Reynolds J, et al. Captum: a unified and generic model interpretability library for PyTorch. *arXiv.* 2020;2009.07896. Available from: <http://arxiv.org/abs/2009.07896>.
67. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: machine learning in python. *J Mach Learn Res.* 2011;12:2825–30.
68. Lovering RC, Gaudet P, Acencio ML, Ignatchenko A, Jolma A, Fornes O, et al. A GO catalogue of human DNA-binding transcription factors. *Biochim Biophys Acta BBA - Gene Regul Mech.* 2021;1864:194765.

69. Breeze CE, Lazar J, Mercer T, Halow J, Washington I, Lee K, et al. Atlas and developmental dynamics of mouse DNase I hypersensitive sites. *bioRxiv*. 2020;2020.06.26.172718. Available from: <https://doi.org/10.1101/2020.06.26.172718>.
70. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*. 2010;26:841–2.
71. Asif M, Orenstein Y. DeepSELEX: inferring DNA-binding preferences from HT-SELEX data using multi-class CNNs. *Bioinformatics*. 2020;36:i634–42.
72. Khan A, Riudavets Puig R, Boddie P, Mathelier A. BiasAway: command-line and web server to generate nucleotide composition-matched DNA background sequences. *Bioinformatics*. 2021;37:1607–9.
73. Ashuach T, Reidenbach DA, Gayoso A, Yosef N. PeakVL: a deep generative model for single-cell chromatin accessibility analysis. *Cell Rep Methods*. 2022;2:100182.
74. Gayoso A, Lopez R, Xing G, Boyeau P, Valiollah Pour Amiri V, Hong J, et al. A Python library for probabilistic analysis of single-cell omics data. *Nat Biotechnol*. 2022;40:163–6.
75. Kingma DP, Ba J. Adam: a method for stochastic optimization. *arXiv*. 2017;1412.6980. Available from: <http://arxiv.org/abs/1412.6980>.
76. Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*. 2009;25:1422–3.
77. Novakovsky G, Fornes O, Saraswat M, Mostafavi S, Wasserman WW. ExplainNN: interpretable and transparent neural networks for genomics: code. GitHub. 2023. Available from: <https://github.com/wassermanlab/ExplainNN>.
78. Novakovsky G, Fornes O, Saraswat M, Mostafavi S, Wasserman WW. ExplainNN: interpretable and transparent neural networks for genomics: code and data. Zenodo. 2023. Available from: <https://doi.org/10.5281/zenodo.7995928>.
79. Harris CR, Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, et al. Array programming with NumPy. *Nature*. 2020;585:357–62.
80. Novakovsky G, Saraswat M, Fornes O, Mostafavi S, Wasserman WW. Biologically-relevant transfer learning improves transcription factor binding prediction: TF binding matrices. Zenodo. 2021. <https://doi.org/10.5281/zenodo.5283416>.
81. Novakovsky G, Saraswat M, Fornes O, Mostafavi S, Wasserman WW. Biologically-relevant transfer learning improves transcription factor binding prediction: TF binding matrices. GitHub. 2021. Available from: <https://github.com/wassermanlab/TF-Binding-Matrix>.
82. Novakovsky G, Fornes O, Saraswat M, Mostafavi S, Wasserman WW. ExplainNN: interpretable and transparent neural networks for genomics: data. GitHub. 2023. Available from: <https://github.com/wassermanlab/PWM-to-filter-weights>.
83. HT-SELEX data: GATA3: cycle 1. SRA. 2015. Available from: <https://www.ncbi.nlm.nih.gov/sra/?term=ERR1003435>.
84. HT-SELEX data: GATA3: cycle 2. SRA. 2015. Available from: <https://www.ncbi.nlm.nih.gov/sra/?term=ERR1003437>.
85. HT-SELEX data: GATA3: cycle 3. SRA. 2015. Available from: <https://www.ncbi.nlm.nih.gov/sra/?term=ERR1003439>.
86. HT-SELEX data: GATA3: cycle 4. SRA. 2015. Available from: <https://www.ncbi.nlm.nih.gov/sra/?term=ERR1003441>.
87. SMiLE-seq data: GATA3. SRA. 2016. Available from: <https://www.ncbi.nlm.nih.gov/sra/?term=SRR3405148>.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

