

# Genomics enters the deep learning era.

**Etienne Routhier** <sup>Corresp., 1</sup>, **Julien Mozziconacci** <sup>Corresp., 2</sup>

<sup>1</sup> LPTMC, Sorbonne Université, Paris, France

<sup>2</sup> StrInG Lab, Museum national d'Histoire naturelle, Paris, France

Corresponding Authors: Etienne Routhier, Julien Mozziconacci  
Email address: etiennerouthier@gmail.com, julien.mozziconacci@mnhn.fr

The tremendous amount of biological sequence data available combined with the recent methodological breakthrough in deep learning in domains such as computer vision or natural language processing, is leading today to the transformation of bioinformatics through the emergence of deep genomics, the application of deep learning to genomic sequences. We review here the new applications that the use of deep learning enables in the field, focusing on three aspects: the functional annotation of genomes, the sequence determinants of the genome functions and the possibility to write synthetic genomic sequences.

# Genomics enters the deep learning era

Etienne Routhier<sup>1</sup> and Julien Mozziconacci<sup>2</sup>

<sup>1</sup>LPTMC, Sorbonne Université, Paris, France ,

[etienne.routhier@sorbonne-universite.fr](mailto:etienne.routhier@sorbonne-universite.fr) ,

<sup>2</sup>StrInG Lab, MNHN, Paris, France , [julien.mozziconacci@mnhn.fr](mailto:julien.mozziconacci@mnhn.fr) ,

## ABSTRACT

The tremendous amount of biological sequence data available combined with the recent methodological breakthrough in deep learning in domains such as computer vision or natural language processing, is leading today to the transformation of bioinformatics through the emergence of deep genomics, the application of deep learning to genomic sequences. We review here the new applications that the use of deep learning enables in the field, focusing on three aspects: the functional annotation of genomes, the sequence determinants of the genome functions and the possibility to write synthetic genomic sequences.

## INTRODUCTION

Genomics is the field in life science focusing on genomic sequences (Fig.1a) and attempting to link the DNA sequence of a living organism with its physical and molecular characteristics. High-throughput sequencing techniques provide huge amounts of data to reconstruct this link. These techniques can now provide both the linear genome sequence and a lot of other information such as the genome 3D structure in cells (Hi-C), the nucleosome and other proteins bindings sites found along the molecule (MNase-seq, ChIP-seq), the local accessibility of the DNA sequence (DNase-seq), the epigenetic marks found on nucleosomes (ChIP-seq) and the activity of genes (RNA-seq, CAGE). Machine learning has long played an important role in the processing of these huge amounts of data (Libbrecht and Noble (2015)) and deep learning has recently emerged as a promising methodology to renew these machine learning approaches. This trend is shared by all bio-medical fields (Fig.1a) for which the number of publications regarding the application of deep learning is exploding (Holder et al. (2017); Ho et al. (2019); Zitnik et al. (2019); Zemouri et al. (2019)).

Schematically, machine or deep learning has been applied to genomics for two main tasks (Fig.1b). First, it has been used to predict second order (i.e. functional) annotation using the first order annotation (i.e. the experimental measures such as ChIP-seq, RNA-seq, ...). This process consists in labeling each DNA segment along the genome with a function (e.g. promoter, protein binding site, enhancer, ...). We will call here this general task *genome annotation*, going beyond the mere annotation of genes. Secondly, machine learning can also be used to annotate (first and/or second order) the genome directly from the DNA sequence. This review focuses on the application of deep learning for this second task (Fig1a-b, red dotted line boxes).

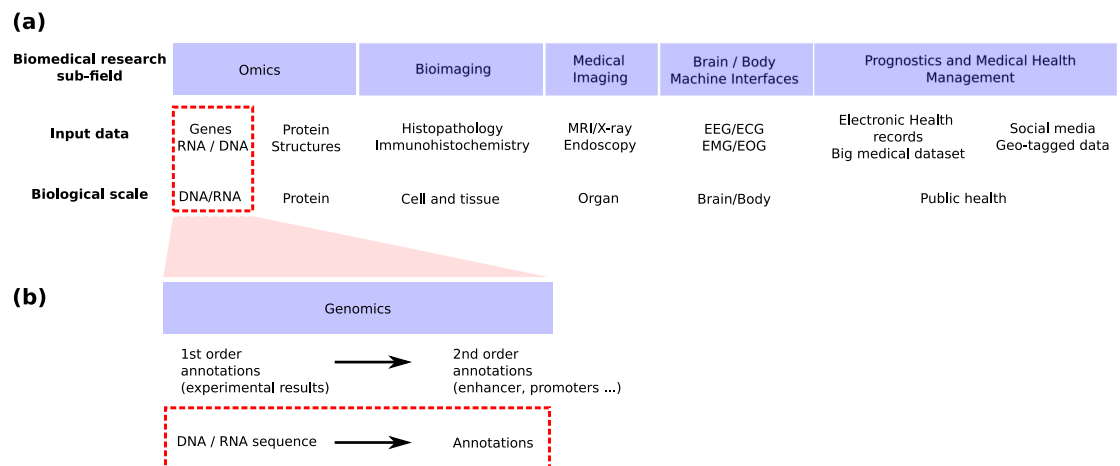
Other reviews focus on the application of deep learning to genomics and proteomics (Zou et al. (2019); Eraslan et al. (2019); Zhang et al. (2019c); Yue and Wang (2018)), often with an introduction to technical aspects and a rather broad domain focus. They present the different neural network architectures used in various types of applications as well as the potential pitfalls Koumakis (2020). Our goal here is to provide a complementary view, focusing on the practical benefits of the application of deep learning to the task of genome annotation from the DNA sequence. This review is intended to biologists and bioinformaticians who are curious to know what new questions can be efficiently tackled using deep learning, how deep learning may help them in their own studies and maybe change their perspectives on their field.

Amongst the first and most emblematic methodologies that were proposed, DeepSEA (Zhou and Troyanskaya (2015)), DeepBind (Alipanahi et al. (2015)) and Basset (Kelley et al. (2016)), are similar in both technical set up and goals. They all use a Convolutional Neural Network (CNN) architecture, which was originally used in computer vision Fukushima (1980); LeCun et al. (1989), to predict whether an input sequence is accessible (i.e. harbors a DNase peak), contains transcription factors binding sites (TFBS), or specific histone modifications (as assessed by ChIP-seq). Since these pioneering approaches were

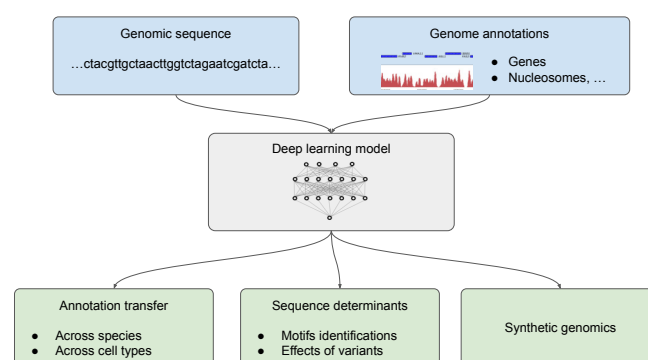
developed, the number of different methodologies used in the field has rapidly grown and the diversity in both domain of application and technical methods has exploded (see Table 1).

We divide the application of deep learning to genome annotation into three different goals (Fig.2): (1) transferring a known annotation for a given species or a given cell type to another species or a different cell type, practically enabling the automatic annotation of genomes from the DNA sequence; (2) getting a deeper understanding of sequence determinants of the genome function by predicting the effect of non-coding mutations and determining sequence motifs that are recognized by the cell machinery; (3) designing synthetic DNA sequences with a tailored annotations (Fig. 2). Two related applications for which classification methodologies from machine learning are inherently useful are the tasks of species attribution for a given, usually short, sequence and the task of assessing sequencing errors are covered in the insert "Sequence classification".

For the sake of completeness and for those who are interested in more technical aspects, we also compiled a short description of the publications that focus on methodological advances (Table 2).



**Figure 1. Positioning of this review within the field of deep learning for biological/biomedical application** (a) We adapted the segmentation of the field proposed by Zemouri et al. (Zemouri et al. (2019)) to position this review. (b) Zoom into the genomics field. This review focus on the application of deep learning to annotate the genome directly from the DNA sequence (red dashed line boxes).



**Figure 2. Different possible uses of deep learning in genomics** Deep learning models trained with genome annotations together with the underlying genomic sequence (in light blue) can be used for three different applications (in light green) (1) to automatically annotate the genome of a given species and for a given cell type, (2) to determine the sequence determinants of the genome functions by identifying sequence motifs (such as position weight matrix, PWM) and the effect of sequence variants, or (3) to design artificial sequences.

# **SURVEY METHODOLOGY**

Deep learning for genomics is a rapidly evolving field. We did our best to gather many of the studies published so far focusing on the use of deep neural networks that take as input genomic sequences. To ensure an unbiased analysis of literature, a comprehensive analysis of published articles was carried using the following online databases: Medline (PubMed), Science Direct (<http://sciencedirect.com>) database, and Google Scholar database. We used the following keywords: deep learning, neural networks, genomics, DNA sequence and then gathered articles together with articles that were cited within the recovered articles.

## **1 ANNOTATION TRANSFER**

Annotation transfer relies on what the machine learning field calls “out-of-distribution generalization”. The general idea behind is that when a neural network learns the link between a given annotation and a sequence, it may be able to generalize, that is to infer the annotation for another, albeit related sequence. While neural networks have been found to generalize well in this way (for a contemporary review, see Shen et al. (2021)), it is worth noting that it is always possible to design specific examples for which this generalization will fail Wiyatno et al. (2019). For this reason, it is always a plus to have in hand some experimental data of the annotation you seek to predict. Even if the data has a low coverage, it can be used to validate or even fine tune the model using this new data Iman et al. (2022).

The use of deep learning methodologies to annotate genomes from the sequence by transferring annotations learnt in a different context has been reported so far for three different applications. The first is the transfer from a species on which the network was trained to another species on which the predictions are made. The second one is the transfer of annotations from one cell type (or environmental condition) on which the network was trained to another cell type (or another environmental condition). This application relies on the use of secondary annotation that has to be used on top the DNA sequence to train the model.

### **Annotation transfer across species**

The potential of using deep learning to transfer annotations is especially relevant for different species since experimentally annotating all the sequenced genomes is today impossible. Here again, deep learning methodologies may help to close the gap by automatically annotating new genomes after being trained on reference, well annotated genomes. Initial studies provide a proof of principle of this possibility and highlight some of the limitations. Khodabandelou et al. (Khodabandelou et al. (2020)) demonstrate that a model trained on a given species can be used to annotate another species. They developed a model trained on the human genome to detect transcription start sites (TSS) and faithfully predict TSS on the mouse and on the chicken genomes with this model. Nevertheless, their model failed to generalise to other species such as the zebrafish. This study highlights some limitations of cross species prediction as this possibility relies on the conservation of molecular mechanisms through evolution. The conservation of annotation logic has been further illustrated in Minnoye et al. (Minnoye et al. (2020)) and Chen et al. (Chen et al. (2018)) in the context of enhancer prediction. Minnoye et al. studied the gene expression level associated to enhancers in melanoma for six different species. They demonstrated that the melanoma chromatin accessibility landscape is conserved for homologous enhancers and that the associated TF motifs are also conserved between the six species. More practically, Chen et al. demonstrated that an enhancer predictor trained on one species among human, macaque, mouse, dog, cow and opossum performs correctly on all the other species. A model trained on species A and applied to species B has an AUROC score equal to 96% of the AUROC score of a model trained on B and applied to B (this figure is 85% for AUPRC). Taken together, these studies show the potential of training a model for a specific task on a given context and applying it to another, a process known as generalization. A promising idea in the field would be increasing the number of species on which the model is trained to increase its cross-species generalization. Kelley showed that deep learning models can be improved when trained on both the human and the mouse genome, especially in the context of predicting the effect of non-coding variant (Kelley (2020)). This study shows that the diversity in the training set provided by training on two species helps the model generalise in the context of predicting the effect of mutations. Pushing the idea further, Cochran et al. showed the benefits of training on multiple genomes to increase the cross species generalization of models in the context of TFBS predictions (Cochran et al. (2021)). They also improved their predictions by developing a model which penalises the learning of species specific features during the training.

We expect cross species learning and prediction to play a major role in the near future in automatic annotation of ill-annotated genomes. In order to be successful these approaches will need to overcome two potential pitfalls. The first is that the further away species are on the tree of life, the less conserved are molecular mechanisms. The second being the potential over-fitting of a species specific logic. Overfitting is a modeling error in statistics that occurs when a prediction is closely aligned to the training set of data but fails to predict other unseen data. To avoid these pitfalls, testing the predictions with experimental cues is needed to confirm or not the computational predictions. The computational predictions are just one line of evidence about the true annotation. The strength of that evidence depends on how well the model has predicted out-of-distribution instances in the past and how different the distribution of interest is to the training distribution and to other tested distributions. The best solution in an unknown case would be to get a small dataset of experimental annotations to be used to validate and/or fine tune the models.

### Annotation transfer across cell type

In bio-medical research, a few cell types are considered as reference models due to their availability or their potential to be grown in culture. These cell types are extensively studied while for the overwhelming majority of cells types in the human body there is a blatant lack of data. Developing methods that could extrapolate the annotation of the reference cell types to the others can help address this lack. With this goal in mind, specific neural network approaches are developed to annotate the genome in a cellular context that differs from the training context. Knowing that the DNA sequence is conserved between cell types, these methodologies use the DNA sequence and some complementary cell type specific annotation as input.

Nair et al. (Nair et al. (2019)) developed a CNN based model to predict the chromatin accessibility across human cellular context. The model uses the DNA sequence and the gene expression as input and predicts the local DNA accessibility. They evaluate the model genome wide on 10 cell types which were not used in the training set (113 cell types). The model achieves an average area under the precision recall curve (AUPRC) of 0.76 and an area under the receiver operating characteristic curve (AUROC) of 0.954 across five folds. Quang et al. investigated the possibility of cross species prediction in the context of TFBS prediction (Quang and Xie (2019)). They developed a model to predict cell type-specific transcription factor binding from the DNA sequence, the gene expression, the DNaseI cleavage profile and the mappability. The model was also evaluated on cell types that differ from the ones used for training. On 51 TF/cell pairs on which the model was evaluated, it typically achieves an AUROC above 0.97 for most of the TF/cell pairs with a more contrasted figure regarding the AUPRC (between 0.21 and 0.87 depending on the pair). These two studies show promising results for the task of extrapolating annotations from reference cell types to other cell types.

A related problem is the prediction of the annotation of a genome in a different environmental context. Liu et al. showed that deep learning methods can be used to predict the gene expression in yeast from the DNA sequence and from 472 TF and signal molecules binding contexts for different stress conditions (Liu et al. (2019)). Their evaluation was done on the same stress condition as their training as it is not possible to predict gene expression across different conditions. Their model was nevertheless able to perform in-silico TF knock-out experiments that were validated by micro-array experimental results.

## 2 UNVEILING SEQUENCE DETERMINANTS OF GENOMIC ANNOTATIONS

We have seen above that when trained appropriately on large datasets, neural networks are capable of predicting annotations in new contexts. These neural networks in some ways mimic the machinery of DNA binding: metaphorically, they both “read” the DNA sequence and “annotate” it, with computationally predicted biochemical labels or molecules. Such a neural network can be treated as an emulation of the DNA binding machinery and used for in silico experiments to generate hypotheses and theories, like any other model. The real benefit comes from the ease to perform new “experiments” with this neural network and from the possibility to reveal how the model makes its predictions. In other words, the advantage of studying a biological mechanism with a neural network is twofold. Firstly, by studying how the model makes its prediction, one can discover the motifs associated with the biological mechanism in the form of a PWM. Going further, a successful dissection of the model can give access to the grammar of motifs, *i.e.* understanding how motifs interact between themselves, forming motifs of motifs. Methods used to deal with these two tasks are described in the insert “Opening the black box”. Secondly, a model can be used to predict the effect of non-coding mutations, also called variants, on the annotation. Variants are changes

in the genomic sequence relative to the reference sequence. They arise from the natural variability of individual DNA sequences and can be either single nucleotide variations (also called SNPs for Single Nucleotide Polymorphisms), or insertions or deletions of small sequences. The interpretation of the effect of variants within coding sequences is another topic that relates to genetics rather than genomics. In genetics, that covers protein functions and how they are affected by mutations, deep learning is also a game changer but this whole field will not be discussed here. Interested readers may refer to this recent work and reference therein for gene function predictions Brandes et al. (2022) and to this comparative study for protein physical and chemical properties predictions Xu et al. (2020). Non-coding variants can statistically be associated with phenotypic traits or diseases but their mechanistic role cannot be immediately inferred. The statistical approach, also known as Genome Wide Association Study (GWAS), reveals variants that are significantly over-represented in people with a certain trait. This analysis has an important drawback: many variants are linked, i.e. their co-occurrence is significantly more frequent, but within these linked variants, some have no role in creating the phenotype. In other words, GWAS is prone to false positives. Here again, deep learning can be used to prioritize variants, i.e. trying to find the variants responsible for the trait. In order to uncover how the DNA sequence drives the local assembly of various chromatin context we review below the different experimental datasets that have been studied using a deep learning based approach.

# **Epigenomics: transcription factor binding site, histone modification and chromatin accessibility**

The DeepSEA model (Zhou and Troyanskaya (2015)) was specifically developed to study the effect of non-coding SNPs on a huge set of epigenomic data (from 690 TF ChIP-seq, 125 DNase-seq and 104 histone marker ChIP-seq experiments). Despite the fact that the network was trained using a unique human reference genome, it is able to predict the decrease in DNase-seq sensitivity for 57,407 experimentally identified SNPs and these changes were confirmed by experiments. As a paradigmatic example, the network is able to predict the deleterious effect of SNP rs4784227 for FoxA1 protein binding, a mechanism associated with breast cancer. These encouraging results led the authors to use DeepSEA in a general way to discriminate variants associated with functional modification from innocuous variants. The network obtained at the time the best results on this task. Other teams have, since then, improved the quality of predictions on the same experimental dataset and same testset (Tayara and Chong (2019); Quang and Xie (2016)).

The Basenji network (Kelley et al. (2018)) was also used to predict the effect of variants. This network predicts the outcome of 2307 ChIP-seq experiments on histone marks, 949 DNase-seq experiments, and 973 CAGE experiments in human. The predictions of the model in the presence of variants known to alter gene expression were compared to predictions obtained with non-variant sequence. Again, this network, although trained on a unique reference genome, is able to predict a change of the chromatin context at these loci. The model is for instance able to predict the effect of a variant (rs78461372) on the two surrounding genes, one of which is located 13 kbp away. Again, predictions were confirmed experimentally.

Liu et al. demonstrated that a neural network could be used to identify mechanistically disease-associated SNPs from SNPs that co-occur with them (Liu et al. (2018)). They developed a model to predict the chromatin accessibility given the DNA sequence. A test set was created consisting of 29 SNPs known to be related to breast cancer and 1057 harmless SNPs that co-occur with them. A score quantifying the variations in network predictions were found to be significantly higher on disease-associated SNPs than on co-occurring SNPs (one-sided Mann-Whitney U test,  $p\_value = 1.63 \times 10^{-3}$ ). The network training protocol can improve predictions associated with SNPs. Hoffman et al. used a CNN to predict the signals associated with one DNase-seq experiment and 3 ChIP-seq experiments on histone marks from the DNA sequence (Hoffman et al. (2019)). They used the genomes of the individuals on which the experiments were performed as a source of sequences and not the reference genome. They defined a score to evaluate the consequences of 438 million variants. They showed that SNPs with a link to a disease or modifying the expression level of a gene are often attributed with a higher score. Wesolowska et al. (Wesolowska-Andersen et al. (2020)) emphasize the importance of training the network on data obtained on a cell type which is relevant to the studied disease. In order to study the effect of variants associated with type II diabetes, they targeted islet cells of the pancreas. They developed a CNN predicting epigenomic data from the DNA sequence and show consistency between the network prediction-based

method and traditional methods for refining the detection of diabetes-associated variants. A part (roughly 10%) of the initial set of variants can be labeled as important by looking at the model predictions. The authors show that those variants are indeed significantly more likely to be evolutionary conserved than the original set (one sided Wilcoxon rank sum test,  $p\_value = 7.3 \times 10^{-4}$ ). Using this methodology, they were able to find 80% of expression quantitative loci (eQTL, loci associated with a quantitative trait) present in the variant set.

A number of biomedical studies have demonstrated practical application of deep learning for variant analysis and SNP interpretation. Illustrating this every day increasing number, a large study of congenital heart disease was recently performed (Richter et al. (2020)). Another practical application of variant analysis using neural networks is provided by Zhou et al. (Zhou et al. (2019)) to study variants related to autism. The study uses 7097 genomes, 1790 of which are from siblings, and overall covering 127,140 SNPs. These siblings' groups are formed by one member diagnosed with autism while the other is not. On average, the alleles of individuals with autism have a higher effect on transcription than the alleles of their siblings. The effect is larger when SNPs are close to genes and in particular close to loss-of-function intolerant genes. Finally, 34 SNPs considered particularly important were experimentally tested. For 32 of them, a significant modification of the expression of associated genes was observed and among these genes many are active during brain development.

### RNA binding

Deep learning methodologies have also provided new insight into RNA binding processes. For example, by interpreting the first-layer filters (see insert "Opening the black box") of a CNN designed to predict RNA binding sites on the genome, Wang et al. were able to identify new patterns of triple bond formation (between the two DNA strands and RNA). The authors also validated their finding with experiments (Wang et al. (2018a)). Koo et al. (Koo et al. (2018)) developed a CNN model to predict the binding sites of RNA-binding proteins. Analysis of their network by *in silico mutagenesis* (see insert "Opening the black box") showed that it was sensitive not only to consensus sites but also to their number and spacing. This analysis also revealed that the network learned to take into account the RNA secondary structure.

### DNA and RNA methylation

Angermueller et al. developed a neural network capable of predicting methylation sites from the DNA sequence (Angermueller et al. (2017)). An analysis of the filters in the first layer associated with the observation of the predictions finds that GC-rich motifs tend to decrease the methylation of nearby CpG islands in contrast to AT-rich motifs. The motifs associated with the filters (see insert "opening the black box") were then compared with many TF consensus motifs. This analysis showed that 17 filters out of 128 correspond to TFs involved in methylation while 13 others are close to motifs of enzymes involved in methylation.

The miCLIP-seq protocol can be used to measure N6-methyladenosine (m6A) methylation on RNA. Zhang et al. (Zhang and Hamada (2018)) used these data to train a network to detect methylation positions on mRNA sequences. By analyzing the filters in the first layer, they were able to find patterns associated with known m6A readers. Interestingly, they were also able to detect a reader of these methylations, *FMRI*, which was discovered via traditional methods in a paper published the previous year Edupuganti et al. (2017).

### Gene expression

Vaishnav et al. trained a deep transformer network to predict the gene expression level associated with 20 million randomly sampled 80-bp long DNA sequences introduced in a *Saccharomyces cerevisiae* promoter region (Vaishnav et al. (2021)). They assessed the effect of all single mutations in promoter regions and discovered four evolvability archetypes: robust promoters on which mutations have little effect, plastic promoters on which every mutations have a small effect and minimal or maximal promoters on which only some mutations can dramatically decrease or increase the associated expression level. Using this framework and analysing the promoter sequences in 1001 yeast strains, the authors were able to demonstrate that evolution tends to select robust promoters. Earlier, Liu et al. trained a network to predict the expression level of yeast genes under different stresses (Liu et al. (2019)). Analysis of the first convolution layer filters revealed that the network primarily searched for well-documented stress regulatory sites. Transcription factor silencing experiments *in silico* achieved results similar to real microarray experiments. Zrimec et al. also used a CNN to predict mRNA abundance directly from

mRNA sequence in *S. cerevisiae* (Zrimec et al. (2020)). They demonstrated that the entire sequence is useful for determining the level of gene expression. Four elements (promoter, 5'UTR, 3'UTR and termination sequence) are used by the model to make the prediction. By interpreting the model with *in silico mutagenesis*, the authors recovered typical motifs of the four regions: TF binding motifs for the promoter or 5'UTR, the so-called Kozak sequence Kozak (1989) in the 5'UTR, poly-A and T-rich sites for the 3'UTR, and termination sites. More importantly, they demonstrated that mRNA abundance cannot be predicted by the presence or absence of the motifs alone, but can be predicted by the combination of motifs.

Movva et al. trained a network to predict the expression level of genes subjected to artificial regulatory sequences in humans. Interpretation of the network with DeepLIFT (see insert "Opening the black box") reveals that the sequences used by the network to make the prediction correspond to transcription factor binding sites (Movva et al. (2019)). Agarwal et al. predicted gene activity from 10 kbp DNA sequences surrounding the TSS (Agarwal and Shendure (2020)). The authors could not find motifs used by the network but an analysis of the over-represented k-mers in the promoters of highly active genes (according to the network) reveals the importance of CpG islands in predicting gene activity.

### Splicing, translation and polyadenylation of RNA

Cheng et al. developed a neural network to predict gene splice sites from the RNA sequence (Cheng et al. (2019, 2021)). Analysis of the effect of variants using this network shows its utility in understanding the genomic causes of autism. They used the dataset provided by Zhou et al. (Zhou et al. (2019)) that we presented earlier and targeted 3884 mutations that are near exons. They demonstrated that the disruption score of mutations as provided by their model is significantly higher in the affected group than in their unaffected siblings (Wilcoxon rank sum test,  $p\_value = 0.0035$ ). Once again the effect is larger in brain tissues. Jaganathan et al. confirmed the relevance of the use of neural networks for the study of gene splicing in the context of intellectual disability and autism (Jaganathan et al. (2019)). They used data coming from 4293 individuals with intellectual disabilities, 3953 individuals with autism spectrum disorder and 2073 unaffected siblings. *De novo* mutations that are predicted to disrupt splicing are enriched 1.51-fold in intellectual disability and 1.30-fold in autism spectrum disorder compared to healthy controls.

Translation initiation of mRNAs does not always occur at the canonical AUG codon, as shown by the recent QTI-seq method which precisely maps translation initiation sites (Gao et al. (2015)). These data have paved the way for the use of deep learning to predict these initiation sites. Zhang et al. developed a network capable of predicting initiation sites from mRNA sequences (Zhang et al. (2017)). By interpreting their network using input optimization methods, they highlighted the importance of Kozak sequences around AUG codons, confirming the previously established role of these sequences Kozak (1989).

Understanding the mechanisms controlling polyadenylation sites within mRNA sequences is another area that benefited from the contribution of deep learning methods. By interpreting their network, which is able to predict the probability of polyadenylation site usage in human mRNAs, Leung et al. (Leung et al. (2018)) showed that poly(A) sites, the cutting factor UGUA, and GU-rich sequences tend to increase the probability of being a polyadenylation site, whereas the presence of CA-rich sequences decreases this probability. Gao et al. also demonstrated the importance of poly(A) sites in the polyadenylation code in the plant *Arabidopsis thaliana* (Gao et al. (2018)) using a gradient-based method to interpret their network (see insert "Opening the black box").

### 3D genome structure

Fudenberg et al. extended the Basenji method to predict the 3D structure of the genome directly from the DNA sequence (Fudenberg et al. (2020)). Targeted analysis of different areas of the human genome by *in silico mutagenesis* reveals that the CCCTC-binding factor (CTCF) binding sites are the most important elements for structure establishment. By testing the other TFs, the authors reveal that these have no influence apart from their possible interactions with CTCF. By performing CTCF site inversion experiments *in silico*, the authors show that the network is able to learn the role of CTCF motif direction in the 3D structure establishment. Finally, the attribution maps (see insert "Opening the black box") reveal the importance of cohesin ChIP-seq peak and, to a lesser extent, of promoters and enhancers.



### Nucleosome positioning

Routhier et al. used a CNN to predict nucleosome positioning in *S.cerevisiae* directly from the DNA sequence and evaluated the effect of every single mutation on the genome by *in silico mutagenesis*. They demonstrated the core role of the nucleosome depleted region (NDR) in nucleosome positioning and identified nucleosome repulsive motifs that were previously described in the literature. On the other hand, they did not find any motifs that would position nucleosomes by attracting them, suggesting that nucleosome repulsion is the main positioning mechanism (Routhier et al. (2021)). Cakiroglu et al. predicted nucleosome positioning as well as TFBS from the DNA sequence based on results obtained with Micrococcal Nuclease digestion treatment (MNase-seq, Cakiroglu et al. (2021)). The model was able to reproduce the competition between nucleosomes and TFs for binding on the DNA. The analysis of the first layer of the CNN shows that the model identifies TF consensus motifs as important for the prediction and, by removing the filters corresponding to these motifs, the authors also demonstrated that TFs tend to exclude nucleosomes.

## 3 DEEP LEARNING ASSISTED GENOME WRITING

### Anticipation of experimental results and sequence fine tuning

Many cell or developmental biology experiments require the introduction of an artificial DNA fragment into the genome or modifying in some other ways the genome sequence. Having a neural network able to anticipate the consequences of these modifications on many genomic annotations allows to fine tune the experimental protocols not only by refining the sequence to introduce but also its position within the genome. The Kipoi repository gathers many independently developed networks. Its objective is to standardize and simplify the use of trained networks in concrete situations such as experiment support. For example, this repository makes available the DeepMEL network (Minnoye et al. (2020)) developed to predict the accessibility of enhancers in melanoma in several different vertebrate species. The effectiveness of this model to anticipate the expression of enhancer-associated genes has been demonstrated using the CAGI5 challenge data. DeepMEL can be used to predict the activity of artificially introduced enhancers and to optimize their sequence. In a related work in yeast, Zrimec et al. (Zrimec et al. (2020)) used their model to anticipate the genes expression level for various gene constructs, especially changing the terminator (5'UTR + termination sequence) with the promoters left intact. Their predictions were experimentally validated for 6 different genes and show great promise for the experimental control of gene expression by the sequence of surrounding regulatory elements.

The development of synthetic genomics is today largely due to the combination of the CRISPR-Cas9 protocol Jinek et al. (2012), which allows to introduce tailored modifications in the DNA sequence in many organisms, with the industrialization of DNA synthesis Ostrov et al. (2019). The CRISPR-Cas9 protocol uses small RNAs (sgRNA, single guide RNA) to guide the Cas9 protein to its target by sequence complementarity. However, sgRNAs usually target both the desired site and other sites on the DNA. Methodologies have been proposed to anticipate the binding strength between sgRNAs, the desired position and the spurious positions from their sequences (Chuai et al. (2018); Xue et al. (2018)). These networks can be used to design sgRNA sequences that maximize interaction with the desired target and minimize interaction with spurious targets. In order to address the challenges of security and intellectual property raised by the development of synthetic genomics, Nielsen and Voigt have developed a deep learning model to predict the laboratory of origin of artificial plasmids from their DNA sequences (Nielsen and Voigt (2018)). For this specific question, however, deep learning methods do not necessarily deliver the best results (Wang et al. (2021)).

### Synthetic sequence design

Possibly the most exciting prospective of the application of deep genomics is the computer assisted writing of genomes. Indeed, the use of neural networks to predict genomic functions from the sequence opens the possibility of optimizing sequences to control their function. This new research field has seen its first promising results in the recent years.

The study of transcription factor binding sites plays a key role in the application of deep learning to genomics, both for the development of architectures and interpretation methods. This problem has therefore naturally been approached from the perspective of sequence design. Lanchantin et al. (Lanchantin et al. (2016)) optimized the input sequence to maximize the network predictions for specific TF binding. Schreiber et al. developed Ledidi, a methodology to minimally modify the input sequence

of the network in order to modify its predictions (Schreiber et al. (2020)). Using this methodology, the authors are able to induce or destroy sites of CTCF binding or suppress JUND protein binding. Gupta et al. used a method inspired by variational autoencoders to induce SPI1 protein binding sites (Gupta and Kundaje (2019)).

On a different topic, Bogard et al. developed a sequence optimization methodology to design RNA sequences with controlled polyadenylation sites (Bogard et al. (2019)). Linder et al. improved this technique to use it for various problems such as controlling the level of gene transcription, RNA splicing, or RNA 3' cleavage (Linder et al. (2020)). More recently, Linder et al. used masks on the sequence to both determine whether each part of the input sequence was sufficient to explain the network predictions and use this information to generate new sequences with similar properties (Linder et al. (2021)). Other applications include Cuperus et al. who used their trained CNN to predict the translation level of mRNAs from their 5' untranslated sequence (Cuperus et al. (2017)). This network was used to design 5' untranslated regions that induce maximal translation level.

Vaishnav et al. designed promoter regions that induced unusually low or high level of gene expression in yeast *S.cerevisiae* (Vaishnav et al. (2021)). They used a genetic algorithm to write an 80-bp long sequences that produce the desired output. The predictions were made with a deep transformer network trained to predict the gene expression level associated to 20 million randomly sampled promoters. Experimental validation on 500 sequences demonstrated that the sequences actually led to unusual level of expression. On average, designed sequences led to an expression level higher (or lower) than 99% of natural sequences. About 20% of the designed sequences led to a higher (or lower) expression level than any natural sequences.

## CONCLUSION

We highlighted in this review the high potential that deep learning holds to transform classical bioinformatics and open the deep genomics era. We started our tour by listing the first applications in the transfer of genomic annotations between species or between cell types. Due to the tremendous number of genomes that are sequenced everyday, we posit that deep learning will be a game changer for the task of genome annotation. We have also reviewed demonstrations of the potential of these techniques to uncover the complex regulatory grammar of motifs, which go beyond simple motifs represented by PWM that are of common use in the field of functional genomics. We finally presented perhaps the most transformative application of deep learning: the generation of new sequences using sequence optimisation based on predictions or using deep generative models.

Having reviewed these new potential avenues at the intersection between deep learning and genomics, we wish also to mention the risks that comes with using such techniques, as the output of an algorithm should always be taken with caution, especially in cases for which human health is at stake. In clinical medicine, initial enthusiasm for deep learning driven by over-stated results has given way to broad cynicism as the rubber has hit the road. See Wynants et al. (2020); Roberts et al. (2021) for critical reviews of machine learning models for COVID-19 diagnosis.

A typical reason why predictions made by deep learning models may fail is out of distribution sampling. For instance, if the sequences on which a network is trained have an average GC content that is peaked around one value, there is no guarantee that predictions made on sequences harboring a different GC content will turn out to be correct. With that in mind, predictions that lead to annotation transfer, mechanisms discovery or sequence generation should be framed as a hypothesis generation tool to speed up research by suggesting targeted experiments.

We hope that our work will help colleagues in better understanding the profound impact that deep learning will have in the field of bioinformatics. We, as a community, now stand at the beginning of an exciting time: the deep genomics era.

## Sequence classification

### ***Species classification***

Bioinformatics tools used to determine a short sequence's species of origin, such as BLAST, align and compare this sequence with sequences from different reference species. These tools are therefore increasingly slow as the number of reference species increases. Deep learning approaches do not suffer from the same problem. Indeed, once a network is trained to attribute a species to a DNA sequence, the prediction will always take the same amount of time. This advantage also comes at a cost: adding a new reference species to the database would require retraining the network. That retraining process is much more complicated, error-prone and demanding in terms of computational resources than adding a new sequence for BLAST. Networks have nevertheless already been developed for this specific application. The use of k-mer preprocessing allows the prediction of the species with a simple dense network (Vervier et al. (2018)). The results get even better when improving the sequence embedding strategy or increasing the length  $k$  of the k-mers (Menegaux and Vert (2019, 2020)). Other methodologies use CNNs to predict the species of short ribosomal DNA fragments (Busia et al. (2019)) or to identify viruses and microbes from metagenomic data (Liang et al. (2020); Ren et al. (2020)). Another application that today shows great promise is the identification of viral DNA within metagenomic samples (Tampuu et al. (2019)).

### ***Classifying and correcting sequencing errors***

The variations of DNA sequences obtained from sequencers can be due either to the intrinsic diversity of the DNA sequence in the sample or to sequencing errors. In order to obtain the precise pool of sequences in a sample, it is necessary to differentiate the "true" variations from the sequencing errors. Several deep learning based methods have emerged for this purpose (Poplin et al. (2018); Ravasio et al. (2018)). Zhang et al. proposed an improved method that leverages the internal states of a RNN to model the distributions of biased and unbiased RNA-seq reads (Zhang et al. (2019a)). Luo et al. developed similar strategies for long read sequencing (Luo et al. (2018)) and Torracinta et al. developed a methodology that allows the correction of sequence errors for RNA-seq (Torracinta et al. (2016)).

424

## Opening the black box

Neural networks are often referred to as "black boxes", which are trained to give the best answer but from which it is not possible to extract comprehensive rules. In the context of genomics, these rules would allow to understand which sequence motifs are associated with a given annotation. This chapter is intended at reviewing the methods used to "open the black box" in genomics studies and access the DNA motifs and their combination that are the most important for the prediction. Its content is more technical than the rest of the review and it can be skipped by the readers who do not wish to get into these details.

We separate the different methods for interpreting deep networks into two categories. First-order methods allow to determine which DNA motifs play an important role in the network decision. Second-order methods allow to understand the grammar of motifs, i.e. a set of rules that is used to interpret motifs.

### First order – Motif discovery

Convolutional networks include filters (corresponding to short DNA motifs) that are optimized as they are trained. Within the first convolution layer, the network scans the sequence for the occurrence of motifs corresponding to the filters. Studying the motifs corresponding to the filters in the first layer is thus a way to see the DNA patterns deemed important by the network to make its predictions. Kelley et al. (Kelley et al. (2016)) and Alipanahi et al. (Alipanahi et al. (2015)) introduced this method to study the Basset and DeepBind networks that they developed. Computing the motif associated with a filter requires three steps. First, all sequences in the test set are scanned using the filter. Second, for all sequences, at each position where the sequence matches the filter, a subsequence of the filter length is extracted. Matching here means that the norm of the elementwise multiplication between the subsequence and the filter exceeds a threshold. Third, the frequency of the nucleotides A, C, G, T within the extracted subsequences is computed to give a position weight matrix (PWM) of the motif searched by the filter. This method has a major disadvantage: there is no guarantee that the network searches for biologically important patterns with its first layer. This information can be dispersed within all layers. To overcome this limitation, constraints can be applied on the first layer of the network to make its weights directly interpretable in terms of a frequency matrix (Ploenzke and Irizarry (2018)). The network architecture can also be adapted to force information to be contained in the first layer (Koo and Eddy (2019); Koo and Ploenzke (2020a)). An alternative option is to adjust the networks' training procedure to penalize the use of patterns that are too small and therefore not likely to be of biological interest (Tseng et al. (2020)). These methods work well for CNNs but cannot be used directly for RNNs. However, recurrent networks make intermediate predictions during their reading of the DNA sequence. Studying the positions that make these intermediate predictions vary the most can point at the nucleotides that are important in the establishment of the final prediction (Lanchantin et al. (2017)).

A second class of method assigns to each nucleotide of a sequence a score reflecting its importance in the prediction made by the network. If the network predicts several classes, a score will be calculated for each class. There are two ways to compute this score. The first method was introduced for the study of DeepBind and DeepSEA. For each nucleotide, the difference between the predictions obtained with the natural sequence and with a mutated sequence is computed. By summing up the contributions of the three possible mutations a mutation score is obtained (Alipanahi et al. (2015); Nair et al. (2020)). This method is called *in silico mutagenesis*. The second method is based on the estimation of the change in the prediction  $P_c(X_0)$ , obtained for the class  $c$ , that is obtained when the input, one-hot encoded, sequence  $X_0$  is changed to another sequence  $X$ . When  $X$  close to  $X_0$  the Taylor expansion to first order gives:

$$P_c(X) - P_c(X_0) \approx \left. \frac{\partial P_c}{\partial X} \right|_{X_0} (X - X_0) \quad (1)$$

The quantity  $\left. \frac{\partial P_c}{\partial X} \right|_{X_0}$  allows to estimate in which proportion the change of an input element will affect the prediction. As a bonus, this quantity can be easily and cheaply computed with the same methods used during network training. Multiplying this quantity term by term with the input  $X_0$  produces an importance score for each nucleotide of the input sequence (Lanchantin et al. (2017)). Based on similar principles, the DeepLIFT (Shrikumar et al. (2017a); Ancona et al. (2017)) and DeepSHAP (Lundberg and Lee (2017)) methods also use back-propagation to compute how much a change in the input changes the prediction. Both are inspired from methods used in the image recognition field. Some of their potential limitations have been put forward in this field Sturmfels et al. (2020) and it remains to be seen how these limitations will impact the interpretation of deep learning models in the context of genomics.

The importance scores assigned to nucleotides can be used to determine important motifs. Avsec et al. (Avsec et al. (2021b)) developed a methodology, called TF-Modisco, to determine globally important motifs from the assignment scores. This methodology works in three steps. In the first step, an importance score is associated to all positions of the test set sequences with the DeepLIFT model. In a second step, all sub-sequences with high scores are extracted. In order to define these high scores, the scores of the real sequences are compared to the scores obtained for random sequences having the same di-nucleotide distribution. Finally, the sub-sequences are grouped into motifs using hierarchical clustering. Peaks in the importance scores can also be interpreted as peak of ChIP-seq data and standard bioinformatic tools such as MEME (Bailey et al. (2015, 2009)) can be used to extract important motifs (Routhier et al. (2021)).

### **Second order – Grammar of motifs**

The methods used to explore the grammar of motifs can also be divided in two categories: methodologies that exploit the model architecture to compute the interactions between motifs and methodologies that take benefit of the attribution map to visualize the effect of varying the motifs organisation.

CNNs typically have multiple convolution layers. The filters of the second layer can be analyzed in the same way as the filters of the first layer and provide interactions between filters (i.e. motifs) of the first layer. Networks using an attention mechanism to exploit the patterns in the sequence are directly interpretable. The first convolution layer transforms the one-hot encoded sequence in which letters are replaced by vectors of ones and zeros into a 2D matrix, one dimension representing the 1D sequence and the second dimension representing the different filters of the first layer. The attention mechanism assigns a weight to each point of the sequence, a weight learned during training. The "encoded" sequence is then averaged along the spatial axis, weighted by the attention weights. Thus, spatial interactions between filters are readable in the weights assigned by the mechanism (Hu et al. (2019)).

Another approach to decipher the grammar of motif is to perform *in silico* experiments consisting in introducing, moving or destroying motifs and assessing the impact on the predictions. Greenside et al. propose to study the evolution of a motif importance score upon mutation of another motif to discover possible interactions between them (Greenside et al. (2018)). Koo et al. developed a methodology to quantify the global importance of any motif in a general context by evaluating the difference between the average of the predictions obtained for sequences randomly drawn from the natural distribution and the average of the predictions obtained for these same sequences in which the motif has been artificially included (Koo and Ploenzke (2020b)). This method can be used to analyze the interactions between motifs by adding two motifs within the random sequences. Avsec et al. also performed *in silico* experiments to understand the grammar of motifs by changing the genomic distance between the motifs and assessing the evolution of the prediction (Avsec et al. (2021b)).

426

## **ACKNOWLEDGMENTS**

427

428 We would like to thank Lou Duron and Alex Westbrook for their comments on the manuscript. We are  
429 also grateful to the reviewers for their invaluable work.

# REFERENCES

- Agarwal, V. and Shendure, J. (2020). Predicting mrna abundance directly from genomic sequence using deep convolutional neural networks. *Cell reports*, 31(7):107663.
- Al-Ajlan, A. and El Allali, A. (2019). Cnn-mgp: Convolutional neural networks for metagenomics gene prediction. *Interdisciplinary Sciences: Computational Life Sciences*, 11(4):628–635.
- Alipanahi, B., Delong, A., Weirauch, M. T., and Frey, B. J. (2015). Predicting the sequence specificities of dna-and rna-binding proteins by deep learning. *Nature biotechnology*, 33(8):831.
- Ancona, M., Ceolini, E., Öztireli, C., and Gross, M. (2017). Towards better understanding of gradient-based attribution methods for deep neural networks. *arXiv preprint arXiv:1711.06104*.
- Angermueller, C., Lee, H. J., Reik, W., and Stegle, O. (2017). Deepcpg: accurate prediction of single-cell dna methylation states using deep learning. *Genome biology*, 18(1):1–13.
- Arefeen, A., Xiao, X., and Jiang, T. (2019). Deeppasta: deep neural network based polyadenylation site analysis. *Bioinformatics*, 35(22):4577–4585.
- Avsec, Z., Agarwal, V., Visentin, D., Ledsam, J. R., Grabska-Barwinska, A., Taylor, K. R., Assael, Y., Jumper, J., Kohli, P., and Kelley, D. R. (2021a). Effective gene expression prediction from sequence by integrating long-range interactions. *bioRxiv*.
- Avsec, Z., Barekatin, M., Cheng, J., and Gagneur, J. (2018). Modeling positional effects of regulatory sequences with spline transformations increases prediction accuracy of deep neural networks. *Bioinformatics*, 34(8):1261–1269.
- Avsec, Z., Weilert, M., Shrikumar, A., Krueger, S., Alexandari, A., Dalal, K., Fropf, R., McAnany, C., Gagneur, J., Kundaje, A., et al. (2021b). Base-resolution models of transcription-factor binding reveal soft motif syntax. *Nature Genetics*, 53(3):354–366.
- Bailey, T. L., Boden, M., Buske, F. A., Frith, M., Grant, C. E., Clementi, L., Ren, J., Li, W. W., and Noble, W. S. (2009). Meme suite: tools for motif discovery and searching. *Nucleic acids research*, 37(suppl\_2):W202–W208.
- Bailey, T. L., Johnson, J., Grant, C. E., and Noble, W. S. (2015). The meme suite. *Nucleic acids research*, 43(W1):W39–W49.
- Bartoszewicz, J. M., Seidel, A., Rentzsch, R., and Renard, B. Y. (2020). Deepac: predicting pathogenic potential of novel dna with reverse-complement neural networks. *Bioinformatics*, 36(1):81–89.
- Bogard, N., Linder, J., Rosenberg, A. B., and Seelig, G. (2019). A deep neural network for predicting and engineering alternative polyadenylation. *Cell*.
- Brandes, N., Ofer, D., Peleg, Y., Rappoport, N., and Linial, M. (2022). Proteinbert: A universal deep-learning model of protein sequence and function. *Bioinformatics*, 38(8):2102–2110.
- Brown, R. C. and Lunter, G. (2019). An equivariant bayesian convolutional network predicts recombination hotspots and accurately resolves binding motifs. *Bioinformatics*, 35(13):2177–2184.
- Busia, A., Dahl, G. E., Fannjiang, C., Alexander, D. H., Dorfman, E., Poplin, R., McLean, C. Y., Chang, P.-C., and DePristo, M. (2019). A deep learning approach to pattern recognition for short dna sequences. *BioRxiv*, page 353474.
- Cakiroglu, S. A., Steinhauser, S., Smith, J., Xing, W., and Luscombe, N. M. (2021). Chromwave: Deciphering the dna-encoded competition between transcription factors and nucleosomes with deep neural networks. *Available at SSRN 3816949*.
- Cao, Z. and Zhang, S. (2019). Simple tricks of convolutional neural network architectures improve dna–protein binding prediction. *Bioinformatics*, 35(11):1837–1843.
- Chen, D., Jacob, L., and Mairal, J. (2019). Biological sequence modeling with convolutional kernel networks. *Bioinformatics*, 35(18):3294–3302.
- Chen, L., Fish, A. E., and Capra, J. A. (2018). Prediction of gene regulatory enhancers across species reveals evolutionarily conserved sequence properties. *PLoS computational biology*, 14(10):e1006484.
- Cheng, J., Çelik, M. H., Kundaje, A., and Gagneur, J. (2021). Mtsplce predicts effects of genetic variants on tissue-specific splicing. *Genome biology*, 22(1):1–19.
- Cheng, J., Nguyen, T. Y. D., Cygan, K. J., Çelik, M. H., Fairbrother, W. G., Gagneur, J., et al. (2019). Mmsplce: modular modeling improves the predictions of genetic variant effects on splicing. *Genome biology*, 20(1):1–15.
- Chuai, G., Ma, H., Yan, J., Chen, M., Hong, N., Xue, D., Zhou, C., Zhu, C., Chen, K., Duan, B., et al. (2018). Deepcrispr: optimized crispr guide rna design by deep learning. *Genome biology*, 19(1):1–18.
- Cochran, K., Srivastava, D., Shrikumar, A., Balasubramani, A., Kundaje, A., and Mahony, S. (2021).

- 485 Domain adaptive neural networks improve cross-species prediction of transcription factor binding.  
486 *bioRxiv*.
- 487 Cuperus, J. T., Groves, B., Kuchina, A., Rosenberg, A. B., Jojic, N., Fields, S., and Seelig, G. (2017).  
488 Deep learning of the regulatory grammar of yeast 5 prime untranslated regions from 500,000 random  
489 sequences. *Genome research*, 27(12):2015–2024.
- 490 Di Gangi, M., Bosco, G. L., and Rizzo, R. (2018). Deep learning architectures for prediction of nucleosome  
491 positioning from sequences data. *BMC bioinformatics*, 19(14):418.
- 492 Du, X., Yao, Y., Diao, Y., Zhu, H., Zhang, Y., and Li, S. (2018). Deepss: Exploring splice site motif  
493 through convolutional neural network directly from dna sequence. *IEEE Access*, 6:32958–32978.
- 494 Edupuganti, R. R., Geiger, S., Lindeboom, R. G., Shi, H., Hsu, P. J., Lu, Z., Wang, S.-Y., Baltissen, M.,  
495 Jansen, P. W., Rossa, M., et al. (2017). N6-methyladenosine (m6a) recruits and repels proteins to  
496 regulate mrna homeostasis. *Nature structural & molecular biology*, 24(10):870–878.
- 497 Eraslan, G. et al. (2019). Deep learning: new computational modelling techniques for genomics. *Nature*  
498 *Reviews Genetics*, page 1.
- 499 Fudenberg, G., Kelley, D. R., and Pollard, K. S. (2020). Predicting 3d genome folding from dna sequence  
500 with akita. *Nature methods*, 17(11):1111–1117.
- 501 Fukushima, K. (1980). A self-organizing neural network model for a mechanism of pattern recognition  
502 unaffected by shift in position. *Biol. Cybern.*, 36:193–202.
- 503 Gao, X., Wan, J., Liu, B., Ma, M., Shen, B., and Qian, S.-B. (2015). Quantitative profiling of initiating  
504 ribosomes in vivo. *Nature methods*, 12(2):147–153.
- 505 Gao, X., Zhang, J., Wei, Z., and Hakonarson, H. (2018). Deepppolya: a convolutional neural network  
506 approach for polyadenylation site prediction. *IEEE Access*, 6:24340–24349.
- 507 Greenside, P., Shimko, T., Fordyce, P., and Kundaje, A. (2018). Discovering epistatic feature interactions  
508 from neural network models of regulatory dna sequences. *Bioinformatics*, 34(17):i629–i637.
- 509 Guo, Y., Zhou, D., Li, W., Cao, J., Nie, R., Xiong, L., and Ruan, X. (2021). Identifying polyadenylation  
510 signals with biological embedding via self-attentive gated convolutional highway networks. *Applied*  
511 *Soft Computing*, 103:107133.
- 512 Gupta, A. and Kundaje, A. (2019). Targeted optimization of regulatory dna sequences with neural editing  
513 architectures. *bioRxiv*, page 714402.
- 514 Gupta, A. and Rush, A. M. (2017). Dilated convolutions for modeling long-distance genomic dependencies.  
515 *arXiv preprint arXiv:1710.01278*.
- 516 Ho, D. S. W., Schierding, W., Wake, M., Saffery, R., and O’Sullivan, J. (2019). Machine learning snp  
517 based prediction for precision medicine. *Frontiers in genetics*, 10:267.
- 518 Hoffman, G. E., Bendl, J., Girdhar, K., Schadt, E. E., and Roussos, P. (2019). Functional interpretation of  
519 genetic variants using deep learning predicts impact on chromatin accessibility and histone modification.  
520 *Nucleic acids research*, 47(20):10597–10611.
- 521 Holder, L. B., Haque, M. M., and Skinner, M. K. (2017). Machine learning for epigenetics and future  
522 medical applications. *Epigenetics*, 12(7):505–514.
- 523 Hu, H., Xiao, A., Zhang, S., Li, Y., Shi, X., Jiang, T., Zhang, L., Zhang, L., and Zeng, J. (2019). Deephint:  
524 understanding hiv-1 integration via deep learning with attention. *Bioinformatics*, 35(10):1660–1667.
- 525 Iman, M., Rasheed, K., and Arabnia, H. R. (2022). A review of deep transfer learning and recent  
526 advancements. *arXiv preprint arXiv:2201.09679*.
- 527 Jaganathan, K., Panagiotopoulou, S. K., McRae, J. F., Darbandi, S. F., Knowles, D., Li, Y. I., Kosmicki,  
528 J. A., Arbelaez, J., Cui, W., Schwartz, G. B., et al. (2019). Predicting splicing from primary sequence  
529 with deep learning. *Cell*, 176(3):535–548.
- 530 Jinek, M., Chylinski, K., Fonfara, I., Hauer, M., Doudna, J. A., and Charpentier, E. (2012). A pro-  
531 grammable dual-rna-guided dna endonuclease in adaptive bacterial immunity. *science*, 337(6096):816–  
532 821.
- 533 Jing, F., Zhang, S., Cao, Z., and Zhang, S. (2019). An integrative framework for combining sequence  
534 and epigenomic data to predict transcription factor binding sites using deep learning. *IEEE/ACM*  
535 *transactions on computational biology and bioinformatics*.
- 536 Kalkatawi, M., Magana-Mora, A., Jankovic, B., and Bajic, V. B. (2019). Deepgsr: an optimized deep-  
537 learning structure for the recognition of genomic signals and regions. *Bioinformatics*, 35(7):1125–1132.
- 538 Kelley, D. R. (2020). Cross-species regulatory sequence activity prediction. *PLoS computational biology*,  
539 16(7):e1008050.

Kelley, D. R., Reshef, Y. A., Bileschi, M., Belanger, D., McLean, C. Y., and Snoek, J. (2018). Sequential regulatory activity prediction across chromosomes with convolutional neural networks. *Genome research*, 28(5):739–750.

Kelley, D. R., Snoek, J., and Rinn, J. L. (2016). Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome research*, 26(7):990–999.

Khodabandelou, G., Routhier, E., and Mozziconacci, J. (2020). Genome annotation across species using deep convolutional neural networks. *PeerJ Computer Science*, 6:e278.

Koo, P. K., Anand, P., Paul, S. B., and Eddy, S. R. (2018). Inferring sequence-structure preferences of rna-binding proteins with convolutional residual networks. *BioRxiv*, page 418459.

Koo, P. K. and Eddy, S. R. (2019). Representation learning of genomic sequence motifs with convolutional neural networks. *PLoS computational biology*, 15(12):e1007560.

Koo, P. K. and Ploenzke, M. (2020a). Improving representations of genomic sequence motifs in convolutional networks with exponential activations. *bioRxiv*.

Koo, P. K. and Ploenzke, M. (2020b). Interpreting deep neural networks beyond attribution methods: Quantifying global importance of genomic features. *bioRxiv*.

Koumakis, L. (2020). Deep learning models in genomics; are we there yet? *Computational and Structural Biotechnology Journal*, 18:1466–1473.

Kozak, M. (1989). The scanning model for translation: an update. *The Journal of cell biology*, 108(2):229–241.

Lanchantin, J., Singh, R., Lin, Z., and Qi, Y. (2016). Deep motif: Visualizing genomic sequence classifications. *arXiv preprint arXiv:1605.01133*.

Lanchantin, J., Singh, R., Wang, B., and Qi, Y. (2017). Deep motif dashboard: Visualizing and understanding genomic sequences using deep neural networks. In *Pacific Symposium on Biocomputing 2017*, pages 254–265. World Scientific.

LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., and Jackel, L. D. (1989). Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551.

Leung, M. K., Delong, A., and Frey, B. J. (2018). Inference of the human polyadenylation code. *Bioinformatics*, 34(17):2889–2898.

Liang, Q., Bible, P. W., Liu, Y., Zou, B., and Wei, L. (2020). Deepmicrobes: taxonomic classification for metagenomics with deep learning. *NAR Genomics and Bioinformatics*, 2(1):lqaa009.

Libbrecht, M. W. and Noble, W. S. (2015). Machine learning applications in genetics and genomics. *Nature Reviews Genetics*, 16(6):321–332.

Linder, J., Bogard, N., Rosenberg, A. B., and Seelig, G. (2020). A generative neural network for maximizing fitness and diversity of synthetic dna and protein sequences. *Cell Systems*, 11(1):49–62.

Linder, J., La Fleur, A., Chen, Z., Ljubetić, A., Baker, D., Kannan, S., and Seelig, G. (2021). Interpreting neural networks for biological sequences by learning stochastic masks. *bioRxiv*.

Liu, B., Hussami, N., Shrikumar, A., Shimko, T., Bhate, S., Longwell, S., Montgomery, S., and Kundaje, A. (2019). A multi-modal neural network for learning cis and trans regulation of stress response in yeast. *arXiv preprint arXiv:1908.09426*.

Liu, Q., Xia, F., Yin, Q., and Jiang, R. (2018). Chromatin accessibility prediction via a hybrid deep convolutional neural network. *Bioinformatics*, 34(5):732–738.

Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. In *Proceedings of the 31st international conference on neural information processing systems*, pages 4768–4777.

Luo, R., Sedlazeck, F. J., Lam, T.-W., and Schatz, M. C. (2018). Clairvoyante: a multi-task convolutional deep neural network for variant calling in single molecule sequencing. *bioRxiv*, page 310458.

Luo, X., Tu, X., Ding, Y., Gao, G., and Deng, M. (2020). Expectation pooling: an effective and interpretable pooling method for predicting dna–protein binding. *Bioinformatics*, 36(5):1405–1412.

Menegaux, R. and Vert, J.-P. (2019). Continuous embeddings of dna sequencing reads and application to metagenomics. *Journal of Computational Biology*, 26(6):509–518.

Menegaux, R. and Vert, J.-P. (2020). Embedding the de bruijn graph, and applications to metagenomics. *bioRxiv*.

Min, X., Chen, N., Chen, T., and Jiang, R. (2016). Deepenhancer: Predicting enhancers by convolutional neural networks. In *2016 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 637–644. IEEE.



- Min, X., Zeng, W., Chen, N., Chen, T., and Jiang, R. (2017a). Chromatin accessibility prediction via convolutional long short-term memory networks with k-mer embedding. *Bioinformatics*, 33(14):i92–i101.
- Min, X., Zeng, W., Chen, S., Chen, N., Chen, T., and Jiang, R. (2017b). Predicting enhancers with deep convolutional neural networks. *BMC bioinformatics*, 18(13):35–46.
- Minnoye, L., Taskiran, I. I., Mauduit, D., Fazio, M., Van Aerschot, L., Hulselmans, G., Christiaens, V., Makhzami, S., Seltenhammer, M., Karras, P., et al. (2020). Cross-species analysis of enhancer logic using deep learning. *Genome research*, 30(12):1815–1834.
- Mostavi, M., Salekin, S., and Huang, Y. (2018). Deep-2'-o-me: predicting 2'-o-methylation sites by convolutional neural networks. In *2018 40th annual international conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 2394–2397. IEEE.
- Movva, R., Greenside, P., Marinov, G. K., Nair, S., Shrikumar, A., and Kundaje, A. (2019). Deciphering regulatory dna sequences and noncoding genetic variants using neural network models of massively parallel reporter assays. *PLoS One*, 14(6):e0218073.
- Nair, S., Kim, D. S., Perricone, J., and Kundaje, A. (2019). Integrating regulatory dna sequence and gene expression to predict genome-wide chromatin accessibility across cellular contexts. *Bioinformatics*, 35(14):i108–i116.
- Nair, S., Shrikumar, A., and Kundaje, A. (2020). fastism: Performant in-silico saturation mutagenesis for convolutional neural networks. *bioRxiv*.
- Nielsen, A. A. and Voigt, C. A. (2018). Deep learning to predict the lab-of-origin of engineered dna. *Nature communications*, 9(1):1–10.
- Ostrov, N., Beal, J., Ellis, T., Gordon, D. B., Karas, B. J., Lee, H. H., Lenaghan, S. C., Schloss, J. A., Stracquadanio, G., Trefzer, A., et al. (2019). Technological challenges and milestones for writing genomes. *Science*, 366(6463):310–312.
- Pachganov, S., Murtazaliev, K., Zarubin, A., Sokolov, D., Chartier, D. R., and Tatarinova, T. V. (2019). Transprise: a novel machine learning approach for eukaryotic promoter prediction. *PeerJ*, 7:e7990.
- Ploenzke, M. S. and Irizarry, R. A. (2018). Interpretable convolution methods for learning genomic sequence motifs. *bioRxiv*, page 411934.
- Poplin, R., Chang, P.-C., Alexander, D., Schwartz, S., Colthurst, T., Ku, A., Newburger, D., Dijamco, J., Nguyen, N., Afshar, P. T., et al. (2018). A universal snp and small-indel variant caller using deep neural networks. *Nature biotechnology*, 36(10):983–987.
- Quang, D. and Xie, X. (2016). Danq: a hybrid convolutional and recurrent deep neural network for quantifying the function of dna sequences. *Nucleic acids research*, 44(11):e107–e107.
- Quang, D. and Xie, X. (2019). FactorNet: a deep learning framework for predicting cell type specific transcription factor binding from nucleotide-resolution sequential data. *Methods*, 166:40–47.
- Ravasio, V., Ritelli, M., Legati, A., and Giacomuzzi, E. (2018). Garfield-ngs: Genomic variants filtering by deep learning models in ngs. *Bioinformatics*, 34(17):3038–3040.
- Ren, J., Song, K., Deng, C., Ahlgren, N. A., Fuhrman, J. A., Li, Y., Xie, X., Poplin, R., and Sun, F. (2020). Identifying viruses from metagenomic data using deep learning. *Quantitative Biology*, pages 1–14.
- Richter, F., Morton, S. U., Kim, S. W., Kitaygorodsky, A., Wasson, L. K., Chen, K. M., Zhou, J., Qi, H., Patel, N., DePalma, S. R., et al. (2020). Genomic analyses implicate noncoding de novo variants in congenital heart disease. *Nature genetics*, 52(8):769–777.
- Roberts, M., Driggs, D., Thorpe, M., Gilbey, J., Yeung, M., Ursprung, S., Aviles-Rivero, A. I., Etmann, C., McCague, C., Beer, L., et al. (2021). Common pitfalls and recommendations for using machine learning to detect and prognosticate for covid-19 using chest radiographs and ct scans. *Nature Machine Intelligence*, 3(3):199–217.
- Routhier, E., Pierre, E., Khodabandelou, G., and Mozziconacci, J. (2021). Genome-wide prediction of dna mutation effect on nucleosome positions for yeast synthetic genomics. *Genome Research*, 31(2):317–326.
- Schreiber, J., Libbrecht, M., Bilmes, J., and Noble, W. S. (2017). Nucleotide sequence and dnasei sensitivity are predictive of 3d chromatin architecture. *bioRxiv*, page 103614.
- Schreiber, J., Lu, Y. Y., and Noble, W. S. (2020). Ledidi: Designing genomic edits that induce functional activity. *bioRxiv*.
- Schwessinger, R., Gosden, M., Downes, D., Brown, R. C., Oudelaar, A. M., Telenius, J., Teh, Y. W., Lunter, G., and Hughes, J. R. (2020). Deepc: predicting 3d genome folding using megabase-scale

- transfer learning. *Nature methods*, 17(11):1118–1124.
- Shen, Z., Liu, J., He, Y., Zhang, X., Xu, R., Yu, H., and Cui, P. (2021). Towards out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*.
- Shrikumar, A., Greenside, P., and Kundaje, A. (2017a). Learning important features through propagating activation differences. In *International Conference on Machine Learning*, pages 3145–3153. PMLR.
- Shrikumar, A., Greenside, P., and Kundaje, A. (2017b). Reverse-complement parameter sharing improves deep learning models for genomics. *bioRxiv*, page 103663.
- Singh, S., Yang, Y., Póczos, B., and Ma, J. (2019). Predicting enhancer-promoter interaction from genomic sequence with deep neural networks. *Quantitative Biology*, 7(2):122–137.
- Sturmfels, P., Lundberg, S., and Lee, S.-I. (2020). Visualizing the impact of feature attribution baselines. *Distill*. <https://distill.pub/2020/attribution-baselines>.
- Tampuu, A., Bzhalava, Z., Dillner, J., and Vicente, R. (2019). Viraminer: Deep learning on raw dna sequences for identifying viral genomes in human samples. *PloS one*, 14(9):e0222271.
- Tayara, H. and Chong, K. T. (2019). Improving the quantification of dna sequences using evolutionary information based on deep learning. *Cells*, 8(12):1635.
- Tian, R., Zhou, P., Li, M., Tan, J., Cui, Z., Xu, W., Wei, J., Zhu, J., Jin, Z., Cao, C., et al. (2021). Deephpv: a deep learning model to predict human papillomavirus integration sites. *Briefings in Bioinformatics*, 22(4):bbaa242.
- Torracinta, R., Mesnard, L., Levine, S., Shakhovich, R., and Hanson, M. (2016). Adaptive somatic mutations calls with deep learning and semi-simulated data. *bioRxiv*, page 079087.
- Trabelsi, A., Chaabane, M., and Ben-Hur, A. (2019). Comprehensive evaluation of deep learning architectures for prediction of dna/rna sequence binding specificities. *Bioinformatics*, 35(14):i269–i277.
- Tseng, A., Shrikumar, A., and Kundaje, A. (2020). Fourier-transform-based attribution priors improve the interpretability and stability of deep learning models for genomics. *Advances in Neural Information Processing Systems*, 33.
- Umarov, R., Kuwahara, H., Li, Y., Gao, X., and Solovyev, V. (2019). Promoter analysis and prediction in the human genome using sequence-based deep learning models. *Bioinformatics*, 35(16):2730–2737.
- Umarov, R. K. and Solovyev, V. V. (2017). Recognition of prokaryotic and eukaryotic promoters using convolutional deep learning neural networks. *PloS one*, 12(2):e0171410.
- Vaishnav, E. D., de Boer, C. G., Yassour, M., Molinet, J., Fan, L., Adiconis, X., Thompson, D. A., Cubillos, F. A., Levin, J. Z., and Regev, A. (2021). A comprehensive fitness landscape model reveals the evolutionary history and future evolvability of eukaryotic cis-regulatory dna sequences. *bioRxiv*.
- Vervier, K., Mahé, P., and Vert, J.-P. (2018). Metavw: Large-scale machine learning for metagenomics sequence classification. In *Data Mining for Systems Biology*, pages 9–20. Springer.
- Wang, F., Chainani, P., White, T., Yang, J., Liu, Y., and Soibam, B. (2018a). Deep learning identifies genome-wide dna binding sites of long noncoding rnas. *RNA biology*, 15(12):1468–1476.
- Wang, J. and Wang, L. (2019). Deep learning of the back-splicing code for circular rna formation. *Bioinformatics*, 35(24):5235–5242.
- Wang, M., Tai, C., E, W., and Wei, L. (2018b). Define: deep convolutional neural networks accurately quantify intensities of transcription factor-dna binding and facilitate evaluation of functional non-coding variants. *Nucleic acids research*, 46(11):e69–e69.
- Wang, Q., Kille, B., Liu, T. R., Elworth, R. L., and Treangen, T. J. (2021). Plasmidhawk improves lab of origin prediction of engineered plasmids using sequence alignment. *Nature communications*, 12(1):1–12.
- Wang, Z., Lei, X., and Wu, F.-X. (2019). Identifying cancer-specific circrna-rbp binding sites based on deep learning. *Molecules*, 24(22):4035.
- Wesolowska-Andersen, A., Yu, G. Z., Nylander, V., Abaitua, F., Thurner, M., Torres, J. M., Mahajan, A., Gloyn, A. L., and McCarthy, M. I. (2020). Deep learning models predict regulatory variants in pancreatic islets and refine type 2 diabetes association signals. *Elife*, 9:e51503.
- Wiyatno, R. R., Xu, A., Dia, O., and de Berker, A. (2019). Adversarial examples in modern machine learning: A review. *arXiv preprint arXiv:1911.05268*.
- Wynants, L., Van Calster, B., Collins, G. S., Riley, R. D., Heinze, G., Schuit, E., Bonten, M. M., Dahly, D. L., Damen, J. A., Debray, T. P., et al. (2020). Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal. *bmj*, 369.

- 705 Xu, W., Zhu, L., and Huang, D.-S. (2019). Dcde: an efficient deep convolutional divergence encoding  
706 method for human promoter recognition. *IEEE transactions on nanobioscience*, 18(2):136–145.
- 707 Xu, Y., Verma, D., Sheridan, R. P., Liaw, A., Ma, J., Marshall, N. M., McIntosh, J., Sherer, E. C., Svetnik,  
708 V., and Johnston, J. M. (2020). Deep dive into machine learning models for protein engineering.  
709 *Journal of chemical information and modeling*, 60(6):2773–2790.
- 710 Xue, L., Tang, B., Chen, W., and Luo, J. (2018). Prediction of crispr sgRNA activity using a deep  
711 convolutional neural network. *Journal of chemical information and modeling*, 59(1):615–624.
- 712 Yang, B., Liu, F., Ren, C., Ouyang, Z., Xie, Z., Bo, X., and Shu, W. (2017). Biren: predicting enhancers  
713 with a deep-learning-based model using the DNA sequence alone. *Bioinformatics*, 33(13):1930–1936.
- 714 Yin, Q., Wu, M., Liu, Q., Lv, H., and Jiang, R. (2019). DeepHistone: a deep learning approach to  
715 predicting histone modifications. *BMC genomics*, 20(2):11–23.
- 716 Yue, T. and Wang, H. (2018). Deep learning for genomics: A concise overview. *arXiv preprint*  
717 *arXiv:1802.00810*.
- 718 Zemouri, R., Zerhouni, N., and Racocceanu, D. (2019). Deep learning in the biomedical applications:  
719 Recent and future status. *Applied Sciences*, 9(8):1526.
- 720 Zeng, H., Edwards, M. D., Liu, G., and Gifford, D. K. (2016). Convolutional neural network architectures  
721 for predicting DNA–protein binding. *Bioinformatics*, 32(12):i121–i127.
- 722 Zeng, H. and Gifford, D. K. (2017). Predicting the impact of non-coding variants on DNA methylation.  
723 *Nucleic acids research*, 45(11):e99–e99.
- 724 Zeng, W., Wang, Y., and Jiang, R. (2020). Integrating distal and proximal information to predict gene  
725 expression via a densely connected convolutional neural network. *Bioinformatics*, 36(2):496–503.
- 726 Zhang, J., Peng, W., and Wang, L. (2018). Lenup: learning nucleosome positioning from DNA sequences  
727 with improved convolutional neural networks. *Bioinformatics*, 34(10):1705–1712.
- 728 Zhang, Q., Shen, Z., and Huang, D.-S. (2019a). Modeling in-vivo protein–DNA binding by combining  
729 multiple-instance learning with a hybrid deep neural network. *Scientific reports*, 9(1):1–12.
- 730 Zhang, Q., Shen, Z., and Huang, D.-S. (2019b). Predicting in-vitro transcription factor binding sites using  
731 DNA sequence+ shape. *IEEE/ACM transactions on computational biology and bioinformatics*.
- 732 Zhang, S., Hu, H., Jiang, T., Zhang, L., and Zeng, J. (2017). Titer: predicting translation initiation sites by  
733 deep learning. *Bioinformatics*, 33(14):i234–i242.
- 734 Zhang, Y. and Hamada, M. (2018). Deepm6aseq: prediction and characterization of m6A-containing  
735 sequences using deep learning. *BMC bioinformatics*, 19(19):1–11.
- 736 Zhang, Y., Qiao, S., Ji, S., and Li, Y. (2020). DeepSite: bidirectional LSTM and CNN models for predicting  
737 DNA–protein binding. *International Journal of Machine Learning and Cybernetics*, 11(4):841–851.
- 738 Zhang, Z., Zhao, Y., Liao, X., Shi, W., Li, K., Zou, Q., and Peng, S. (2019c). Deep learning in omics: a  
739 survey and guideline. *Briefings in functional genomics*, 18(1):41–57.
- 740 Zhou, J. (2021). Sequence-based modeling of genome 3D architecture from kilobase to chromosome-scale.  
741 *bioRxiv*.
- 742 Zhou, J., Park, C. Y., Theesfeld, C. L., Wong, A. K., Yuan, Y., Scheckel, C., Fak, J. J., Funk, J., Yao, K.,  
743 Tajima, Y., et al. (2019). Whole-genome deep-learning analysis identifies contribution of noncoding  
744 mutations to autism risk. *Nature genetics*, 51(6):973–980.
- 745 Zhou, J. and Troyanskaya, O. G. (2015). Predicting effects of noncoding variants with deep learning–based  
746 sequence model. *Nature methods*, 12(10):931.
- 747 Zitnik, M., Nguyen, F., Wang, B., Leskovec, J., Goldenberg, A., and Hoffman, M. M. (2019). Ma-  
748 chine learning for integrating data in biology and medicine: Principles, practice, and opportunities.  
749 *Information Fusion*, 50:71–91.
- 750 Zou, J., Huss, M., Abid, A., Mohammadi, P., Torkamani, A., and Telenti, A. (2019). A primer on deep  
751 learning in genomics. *Nature genetics*, 51(1):12–18.
- 752 Zrimec, J., Börlin, C. S., Buric, F., Muhammad, A. S., Chen, R., Siewers, V., Verendel, V., Nielsen, J.,  
753 Töpel, M., and Zelezniak, A. (2020). Deep learning suggests that gene expression is encoded in all  
754 parts of a co-evolving interacting gene regulatory structure. *Nature communications*, 11(1):1–16.

Annotation	Usage	Preprocessing	Data	Species	Architecture	Reference
TFBS	Transfer	one-hot-encoding	DNA + gene expression + DNaseI cleavage	human	CNN + RNN	Quang and Xie (2019)
			DNA sequence	human + mouse	CNN	Cochran et al. (2021)
	Bio. mechanism	one-hot-encoding	DNA sequence	human	CNN	Wang et al. (2018b)
				human + mouse + drosophila	CNN	Wang et al. (2018a)
			RNA sequence	human	CNN	Koo et al. (2018)
	Syn. genomics	one-hot-encoding	DNA sequence	human	RNN + Attention	Gupta and Kundaje (2019)
					CNN	Lanchantin et al. (2016)
TFBS + histone + chromatin accessibility	Transfer	one-hot-encoding	DNA sequence	human + mouse	CNN	Kelley (2020)
	Bio. mechanism	one-hot-encoding	DNA sequence	human	CNN	Kelley et al. (2018)
					CNN	Alipanahi et al. (2015)
						Zhou et al. (2019)
						Hoffman et al. (2019)
						Zhou and Troyanskaya (2015)
						Richter et al. (2020)
	Syn. genomics	one-hot-encoding	DNA sequence	human	CNN	Schreiber et al. (2020)
TFBS (circRNA)	Bio. mechanism	one-hot-encoding	RNA sequence	human	CNN	Wang et al. (2019)
chromatin accessibility	Transfer + Bio. mechanism	one-hot-encoding	DNA + gene expression	human	CNN	Nair et al. (2019)
	Bio. mechanism	one-hot-encoding + embedding	DNA sequence	human	CNN	Liu et al. (2018)
gene expression	Transfer + Bio. mechanism	one-hot-encoding	DNA + TF expression level	yeast	CNN	Liu et al. (2019)
	Bio. mechanism Syn. genomics	one-hot-encoding	RNA sequence	7 species	CNN	Zrimec et al. (2020)
				yeast	CNN	Cuperus et al. (2017)
			DNA sequence	Random promoters (yeast)	CNN + Attention + RNN	Vaishnav et al. (2021)
	Bio. mechanism	one-hot-encoding	DNA + mRNA half-life + CG content + ORF length	human	CNN	Agarwal and Shendure (2020)
			DNA + promoter-enhancer interaction	human	CNN	Zeng et al. (2020)
			DNA sequence	human	CNN	Movva et al. (2019)
gene expression + RNA splicing	Syn. genomics	one-hot-encoding	DNA sequence	human	CNN	Linder et al. (2020)

**Table 1. Overview of studies applying deep learning in genomics, segmented by their usage.** CNN: Convolutional Neural Network, RNN: Recurrent Neural Network. After the pioneering use of CNN in genomics in 2015, the methodologies have diversified according to four different aspects: the model inputs (that may include other annotations on top of the sole DNA sequence), the sequence encoding (mainly one-hot-encoding or k-mer embedding), the neural network architecture (CNN, RNN, Attention mechanism) and the output format, which can be either binary or continuous.

Annotation	Usage	Preprocessing	Data	Species	Architecture	Reference
epigenetic mark	Bio. mechanism	one-hot-encoding	DNA + chromatin accessibility	human	CNN	Yin et al. (2019)
			DNA + CpG neighborhood of cells	human	CNN + RNN	Angermueller et al. (2017)
			DNA sequence	human	CNN	Zeng and Gifford (2017)
			RNA sequence	human + mouse + zebrafish	CNN + RNN	Zhang and Hamada (2018)
polyadenylation	Bio. mechanism	one-hot-encoding	DNA sequence	Arabidopsis thaliana	CNN	Gao et al. (2018)
				human	CNN	Leung et al. (2018)
	Syn. genomics	one-hot-encoding	DNA sequence	human	CNN	Bogard et al. (2019)
polyadenylation + translation initiation site	transfer	one-hot-encoding	DNA sequence	human + mouse + bovine + drosophila	CNN	Kalkatawi et al. (2019)
splicing	Bio. mechanism	one-hot-encoding	DNA sequence	human	CNN	Cheng et al. (2019)
						Cheng et al. (2021)
			RNA sequence	human	CNN	Du et al. (2018)
					CNN	Jaganathan et al. (2019)
3D architecture	Bio. mechanism	one-hot-encoding	DNA sequence	human	CNN	Wang and Wang (2019)
					CNN	Zhou (2021)
					CNN	Fudenberg et al. (2020)
				human + mouse	CNN + RNN	Singh et al. (2019)
nucleosome	Bio. mechanism	one-hot-encoding	DNA sequence	yeast	CNN	Schwessinger et al. (2020)
nucleosome + TFBS	Bio. mechanism	one-hot-encoding	DNA sequence	yeast	CNN	Routhier et al. (2021)
enhancer	transfer	embedding	DNA sequence	yeast + human	CNN	Cakiroglu et al. (2021)
	Bio. mechanism + Syn. genomics	one-hot-encoding	DNA sequence	6 species	CNN	Chen et al. (2018)
	Bio. mechanism	one-hot-encoding	DNA sequence	6 pecies	CNN + RNN	Minnoye et al. (2020)
promoter	transfer	one-hot-encoding	DNA sequence	human	CNN	Min et al. (2017b)
promoter + enhancer + TFBS + chromatin accessibility	transfer	one-hot-encoding	DNA sequence	5 species	CNN	Khodabandelou et al. (2020)
translation initiation site	Bio. mechanism	one-hot-encoding	DNA sequence	human	CNN	Wesolowska-Andersen et al. (2020)
sgRNA binding site	Syn. genomics	one-hot-encoding	DNA sequence	human	CNN + RNN	Zhang et al. (2017)
Virus integration	Bio. mechanism	one-hot-encoding	DNA + TFBS + epigenetic + accessibility	human	CNN	Chuai et al. (2018)
			DNA sequence	human + mouse	CNN	Xue et al. (2018)
			DNA sequence	human	CNN + Attention	Tian et al. (2021)

Overview of studies applying deep learning in genomics, segmented by their usage.

Annotation	Preprocessing	Data	Species	Architecture	Reference
TFBS	benchmark	DNA sequence	human	benchmark	Trabelsi et al. (2019)
	embedding	DNA sequence	human	CNN + RNN	Zhang et al. (2019a)
	one-hot-encoding	DNA + distance to various annotations	human	CNN	Avsec et al. (2018)
		DNA + histone marks + accessibility	human	CNN	Jing et al. (2019)
		DNA + shape	human	CNN	Zhang et al. (2019b)
		DNA sequence	human	CNN	Brown and Lunter (2019)
				CNN + RNN	Zhang et al. (2020)
				CNN	Shrikumar et al. (2017b)
				CNN	Luo et al. (2020)
				CNN	Zeng et al. (2016)
				CNN	Chen et al. (2019)
TFBS + histone marks + accessibility	one-hot-encoding	DNA sequence	human	CNN	Cao and Zhang (2019)
				CNN	Kelley et al. (2016)
				CNN + RNN	Tayara and Chong (2019)
				CNN + Attention	Quang and Xie (2016)
				CNN	Avsec et al. (2021a)
chromatin accessibility	embedding	DNA sequence	human	CNN + RNN	Gupta and Rush (2017)
epigenetic marks	embedding	RNA sequence	human	CNN + RNN	Min et al. (2017a)
polyadenylation	embedding	DNA sequence	4 species	CNN	Mostavi et al. (2018)
	one-hot-encoding	RNA + secondary structure	human	CNN + Attention	Guo et al. (2021)
3D architecture	one-hot-encoding	DNA + DNaseI signal	human	CNN + RNN	Arefeen et al. (2019)
	one-hot-encoding	DNA sequence	human	CNN	Schreiber et al. (2017)
nucleosome	one-hot-encoding	DNA sequence	yeast	CNN	Zhang et al. (2018)
				CNN + RNN	Di Gangi et al. (2018)
enhancer	one-hot-encoding	DNA sequence	human	CNN	Yang et al. (2017)
				CNN	Min et al. (2016)
promoter	embedding	DNA sequence	human	CNN	Xu et al. (2019)
	one-hot-encoding	DNA sequence	human	CNN	Umarov and Solovyev (2017)
			human + <i>Oryza sativa</i>	CNN	Umarov et al. (2019)
gene	one-hot-encoding	DNA sequence	metagenomics	CNN	Pachganov et al. (2019)
pathogenicity	one-hot-encoding	DNA sequence	bacterias	CNN	Al-Ajlan and El Allali (2019)
species	one-hot-encoding	DNA sequence	metagenomics	CNN + RNN	Bartoszewicz et al. (2020)
	embedding	DNA sequence	metagenomics	RNN + Attention	Liang et al. (2020)
	one-hot-encoding	DNA sequence	13,838 species	CNN	Busia et al. (2019)
			viruses	CNN	Ren et al. (2020)
					Tampuu et al. (2019)

**Table 2. Overview of studies developing deep learning methodologies in genomics.**