

Albert-Ludwigs University Freiburg
Department of Computer Science
Bioinformatics Group

Master Thesis

Predicting Hi-C contact matrices using machine learning approaches

Author:
Ralf Krauth

Examiner:
Prof. Dr. Rolf Backofen

Second Examiner:
Prof. Dr. Ralf Gilsbach

Advisors:
Anup Kumar, Joachim Wolff

Submission date:
20.04.2021

Abstract

Harhar!

Zusammenfassung

Hohoho!

Erklärung

Hiermit erkläre ich, dass ich die vorliegende Abschlussarbeit selbständig verfasst habe, keine anderen als die angegebenen Quellen und Hilfsmittel verwendet habe und alle Stellen, die wörtlich oder sinngemäß aus anderen Werken entnommen wurden, als solche kenntlich gemacht habe. Darüber hinaus erkläre ich, dass die eingereichte Masterarbeit weder vollständig, noch in wesentlichen Teilen Gegenstand eines anderen Prüfungsverfahrens war oder ist.

Bad Krozingen, den 7. Februar 2021

Ralf Krauth

Contents

1	Introduction	7
1.1	Spatial structure of DNA	7
1.2	The Hi-C process for determining spatial DNA structure	8
1.3	The ChIP-seq process for determining protein binding sites	9
1.4	Motivation and goal of the thesis	9
2	Related work	11
2.1	Methods for predicting DNA-DNA interactions and contact matrices	11
2.2	Image synthesis techniques from computer vision	13
2.3	Discussion of existing work	13
3	Advancing predictions of Hi-C interaction matrices	15
3.1	Dense Neural Network approach	15
3.1.1	Basic network setup	15
3.1.2	Modifying the convolutional part of the network	16
3.1.3	Using a combination of MSE-, TV- and perceptual loss	17
3.1.4	Using a TAD-based loss function	18
3.1.5	Modifying binsize and window size	19
3.1.6	DNA sequence as an additional network input branch	19
3.2	Hi-cGAN approach	20
3.2.1	Pix2Pix as a generic generative network	20
3.2.2	Using a DNN for 1D - 2D conversion	20
3.2.3	Using a CNN for 1D - 2D conversion	20
4	Methods	21
4.1	Input data	21
4.1.1	Hi-C matrices	21
4.1.2	ChIP-seq data	21
4.1.3	Sample generation process for the dense neural network	22
4.1.4	Sample generation for the cGAN	24
4.1.5	Generalization of feature binning	24
4.2	Dense neural network approach	25
4.2.1	Basic setup	25
4.2.2	Modifying kernel size, number of filter layers and filters	26
4.2.3	Custom loss function based on TAD insulation score	26
4.2.4	Combination of mean squared error, perception loss and TV loss	26
4.2.5	Combination of mean squared error and TAD-score-based loss	27
4.3	HiC-GAN approach	27
4.3.1	Pix2Pix network layout	27
4.3.2	Using a DNN for 1D-2D conversion	27
4.3.3	Using a pre-trained DNN for 1D-2D conversion	27
4.3.4	Using a CNN for 1D-2D conversion	27

5	Results	29
5.1	Dense Neural Network approaches	29
5.1.1	Initial results for comparison	29
5.1.2	Results for variations of the convolutional part	29
5.1.3	Results for variations of the loss function	32
5.1.4	Results for different binsizes and windowsizes	32
5.2	Hi-cGAN approaches	32
5.2.1	Initial results	32
6	Discussion and Outlook	47
	References	49
	Acronyms	53

1 Introduction

In recent years, the three-dimensional organization of DNA has been shown to be a key factor for important processes in molecular biology. However, even with the most advanced experimental methods, it remains comparatively expensive to determine the spatial folding of DNA directly, so that current knowledge of three-dimensional DNA organization is still sketchy. In the last five years, several methods have thus been proposed to improve on this situation by determining DNA-DNA interactions *in-silico*. All of these are using existing experimental data which is easier to obtain than spatial data, but correlates with 3D chromatin structure in certain ways. However, most current in-silico approaches have disadvantages and shortcomings, and the current thesis attempts to improve on these.

1.1 Spatial structure of DNA

In the late 1970s, Watson and Crick discovered the now well-known double helix structure of DNA molecules [1]. However, this is not the only relevant spatial structure of genomes. At larger scales, DNA molecules can be wound around certain proteins, so-called histones, forming DNA-protein complexes named nucleosomes. Several of these can further be compacted into fibers, which in turn can be “supercoiled” into the also well-known chromosomes, fig. 1.



Figure 1: spatial chromatin structures (simplified, cf. [2])

While the spatial structure of chromatin outlined above brings along compaction (fig. 1 from left to right), e.g. to make chromosomes fit into eucaryotic cell nuclei, it also has other functional implications. One well known effect of spatial structure is the regulation of gene transcription by establishing or releasing contacts between gene enhancers and the relevant promoters [3, 4]. Since enhancers can be up to a million basepairs away from the promoters, the two can only interact by means of spatial DNA structure. Many effects of spatial structure are still under investigation, two recent studies have for example found dependencies between chromatin conformation and cell differentiation in *Drosophila Melanogaster* [5] and investigated spatial structure dynamics during phase transitions in murine cells [6].

While the driving mechanisms behind the formation of spatial chromatin structures are partially still under research, certain proteins like CTCF or modified histones are already well known to mediate or prevent DNA-DNA contacts [7, 8, 9].

In the last two decades, DNA-sequencing-based techniques have increasingly been utilized to capture the spatial structure of DNA experimentally. The method of choice within this thesis is called Hi-C and shall be explained in the following section.

1.2 The Hi-C process for determining spatial DNA structure

The Hi-C process is an elaborate biochemical procedure for investigating the spatial structure of DNA by detecting DNA-DNA interactions within and across chromosomes. The original Hi-C workflow has been developed by Lieberman-Aiden et al. in 2009 [10] and is depicted in simplified form in fig. 2.

The typical input (In) to Hi-C consists of several millions of cells, which are treated chemically to fix existing DNA-DNA contacts, commonly using formaldehyde, before they are lysed. Next, the DNA is extracted and cut into fragments by certain restriction enzymes (1), usually HindIII or DpnII, and the cut ends are repaired with nucleotides, some of which are marked by biotin (2). The free ends are then joined (3) under conditions which prefer ligations among open ends over ligations between different fragments. Originally, such conditions were achieved by high dilution of the fragments in solvents, but especially this part of the protocol has been replaced by more efficient methods in later works [11, 12]. The ligated fragments are then purified and cut into shorter sequences, some of which contain biotinylated nucleotides and some not (4). The fragments of interest – the ones containing biotinylated nucleotides – are then selected by pulling down biotin (5), for example using magnetic tags, and subjected to paired-end DNA-sequencing (6). In the end, the output of the Hi-C lab process is a large number of short genomic “reads”, which are subsequently processed in the bioinformatics part of the Hi-C protocol outlined in the following section.

On the software side of the protocol, the reads first need to be mapped to the corresponding reference genome. Here, only reads are kept where the “left” sequence (6)(a) uniquely maps to a different region of the reference genome than the “right” sequence (6)(b). These so-called chimeric reads are subjected to quality control, and those passing are counted as an interaction between the two genomic positions (1)(a) and (1)(b) to which the two ends belong. However, at reasonable sequencing depths, reads cannot be counted per base pair. Instead, the reference genome is split into equally sized bins (or regions), and the reads are counted for those regions where they belong. The final outcome of a Hi-C experiment is then a sparse matrix, henceforth referred to as “Hi-C matrix”, which records the interaction count for all possible pairs of regions in the reference genome.

For the thesis at hand, Hi-C matrices can abstractly be considered as sets of sparse, square matrices; one matrix for each chromosome. The individual elements $m_{i,j}$ of the matrices are counts of interactions between the bins with indices $i, j \in \{0, 1, \dots, l\}$, where the last index l is given by the chromosome size cs and the binsize b in basepairs as $l = \lceil \frac{cs}{b} \rceil - 1$. Because interaction counts do not have a “direction”, it holds that $m_{i,j} = m_{j,i}$, so Hi-C matrices are always symmetric.

A bin with index n then corresponds to the genomic region $p_n = [n \cdot b, (n+1) \cdot b)$, except the last bin l , where $p_l = [l \cdot b, cs]$, with l, b, cs as above. As an example, if $m_{0,10} = 22$ for a matrix with

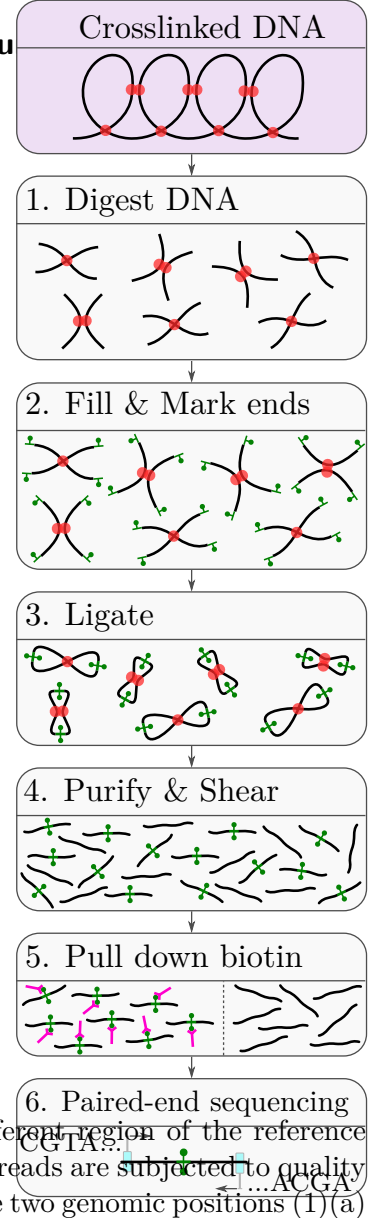


Figure 2: Hi-C lab process

binsize $b = 25$ kbp, this means that in the underlying experiment, 22 interactions have been counted between region p_0 , ranging from 0 to 24 999 bp and region p_{10} , ranging from 250 000 to 274 999 bp. Due to insufficient read coverage, $b = 1$ bp is impossible in practice. Instead, common values include $b \in \{1, 5, 10, 25, 50, 100, 1000\}$ kbp.

In the bioinformatics part of the Hi-C protocol, often just a small fraction of all reads fulfill the selection criteria outlined above, for example due to reads not being chimeric or uniquely mappable. This makes Hi-C a comparatively inefficient, slow and thus expensive process. For example, the well-known dataset by Rao et al. [11] with matrix bin sizes down to 1 kbp, required several *billions* of reads being made.

Parts of fig. 2 and the process description above have been adapted from the preceding study project [13].

XXXmod input steps in graph XXXmark regions a, b throughout the graph

1.3 The ChIP-seq process for determining protein binding sites

The ChIP-seq process is a combination of Chromatin-Immunoprecipitation (ChIP) and DNA-sequencing, designed for investigating DNA-protein interactions [14, 15]. As with Hi-C, the input consists of a sufficient number of cells, which are first fixed with formaldehyde. The cells are then lysed, the DNA-protein structure is extracted and cut into fragments, usually by sonication (2). Next, specific antibodies are added, designed to bind only to a certain protein of interest (3). These antibodies are additionally equipped with a tag, for example a magnetic one, so that the DNA-protein-antibody structures can be precipitated, while fragments without antibodies are discarded (4). The proteins and antibodies are then removed (5), the DNA is purified and finally sequenced (6). Typically, a control experiment is performed together with the ChIP-seq process, which comprises all steps described above, except the immunoprecipitation (3),(4).

The outcome of the ChIP-seq lab process is again a bunch of short genomic sequences, which are then fed into the bioinformatics part of the pipeline.

On the software side of the process, the reads are filtered for quality and mapped to an appropriate reference genome. The number of mapped reads per genomic position can then be simply be counted. It is common to process reads from the control experiment in the same way as reads from the ChIP-seq experiment and then use special software to call “peaks”, i. e. protein binding sites, at those genomic positions where the read count from ChIP-seq is statistically significant compared to the read count from the control group.

For the thesis at hand, the reads from the ChIP-seq experiment have been used directly, because peak data was found to be too sparse for the machine learning approach used in the study project [13].

1.4 Motivation and goal of the thesis

Due to the high effort, Hi-C experiments have been conducted only by few labs and for comparatively few organisms and cell lines so far. For example, as of January 2021, Hi-C assays are available for just 26 of more than 150 human cell lines listed in the Encyclopedia of DNA elements [16, 17], and these have been produced by only five university labs.

However, investigations and experiments with the available Hi-C data have shown correlations between the spatial DNA structure and the binding sites of certain transcription factors and histone modifications [11, 18]. Since these binding sites can be detected with the less costly ChIP-seq method, a computational method able to predict contact matrices from DNA-protein bindings would be very helpful.

The goal of this master thesis is thus to predict missing Hi-C data from existing ChIP-seq data, making use of the known correlations between the binding sites of transcription factors and histone modifications of the one hand and Hi-C interaction counts on the other hand. Because this is a wide field and lead time for a master thesis is limited, the present work will focus on machine learning techniques, which have proven a good choice for exploiting complex patterns and relationships in various kinds of data. XXXcitations needed More specifically, the goal of this thesis is to develop an easy-to-use machine learning approach for predicting Hi-C interaction matrices from ChIP-seq data, using standard in- and output formats with minimal pre- and postprocessing.

To underscore the usefulness of such an approach, table 1 to the side shows the first 16 human cell lines in ENCODE, sorted by total number of selected assays available. It is obvious that even otherwise well-investigated cell lines like HEK293, MCF-7 and H1 lack Hi-C data, while ChIP-seq data is mostly available. XXXdata for H1 is available in

cell line	TF ChIP-seq	Histone ChIP-seq	Dnase-seq	Hi-C
K562	680	20	10	1
HepG2	668	15	3	1
A549	251	87	14	6
HEK293	231	6		
GM12878	211	15	3	7
MCF-7	155	18	8	
H1	88	55	3	
HeLa-S3	77	15	4	1
SK-N-SH	62	21	2	
IMR-90	16	34	2	2
HCT116	27	17	1	8
H9		35	7	
GM23338	15	13	1	
Ishikawa	25		4	
HEK293T	17		1	
GM23248	2	13	1	1

Table 1: availability of selected assays in ENCODE (extract)

2 Related work

In the last five years, several approaches have been presented to determine DNA-DNA interactions *in silico*, using existing data from various experiments. Section 2.1 gives an overview about these methods. Furthermore, some methods originally developed for image synthesis and similar tasks in computer vision might also be useful in the field of Hi-C matrix generation and are thus summarized in subsection 2.2. The section is then concluded by a short discussion of the existing work.

2.1 Methods for predicting DNA-DNA interactions and contact matrices

As of 2020, there is quite a body of existing work in the field of predicting DNA-DNA interactions, using various approaches and different types of input data.

Two conceptually similar methods have been proposed by Brackley et al. in 2016 and MacPherson et al. in 2018 [19, 20]. In both approaches, DNA is modeled as a “beads-on-a-string” polymer, and simulation techniques are employed to find energy-optimal spatial structures of these polymers. Apart from constraints derived from the molecule’s DNA sequence itself, the models also consider spatial contact constraints derived from ChIP-seq experiments of chromatin factors which are known to mediate such DNA-DNA contacts. The interaction matrices derived from the simulations look interesting, but the paper from Brackley et al. [19] is unfortunately lacking a comparison with “true” experimentally measured Hi-C matrices, and the results from MacPherson et al. [20] seem inferior to most other ones presented in this section.

Another simulation-based method has been developed by di Pierro et al. in 2017 [21] and later extended by Qi and Zhang [22]. In both cases, a convolutional neural network (CNN) is trained to learn different “open” and “closed” chromatin states from 11 chromatin factors, and the predicted chromatin states are then taken as constraints for beads-on-a-string models. The difference between [21] and [22] lies mainly in the number of states considered and the simulation methods applied; the results are mathematically convincing in both cases.

A further approach using chromatin states is due to Farrè and Emberly [23]. Here, the conditional probability of two genomic regions being in contact, given their distance and the chromatin state around them, is estimated using Bayes’ rule. In this case, the chromatin state – reduced to active or inactive – is derived from DNA adenine methyltransferase identification (DamID) signals of 53 chromatin factors using probabilistic methods [24]. The conditional probabilities on the right side of Bayes’ rule are either computed from training data or estimated with different probabilistic approaches, too. While the predicted contact matrices do not look like real Hi-C matrices with this approach, highly interacting regions are still often well identifiable.

Three further approaches in the field make use of random forests. 3DEpiLoop by Bkhetan and Plewczynski and Lollipop by Kai et al., both from 2018, use a random forest classifier to predict DNA loops, but differ in input data and preprocessing [25, 26]. While 3DEpiLoop is using only ChIP-seq data of histone modifications and transcription factors [25], Lollipop additionally takes ChIA-PET-, RNA-seq- and DNase-seq-data, CTCF motif orientation and loop length as inputs [26]. Both approaches show good coincidence of predicted loops with experimental data, but their output is binary and rather sparse. Contrary to these two, the third random-forest-based approach, HiC-Reg by Zhang et al. from 2019, allows predicting real-valued Hi-C interaction

matrices directly [27]. To this end, it employs random forest regression to predict interactions between two genomic “windows”, using ChIP-seq data of 13 chromatin factors and the genomic distance of the two windows. The published results for five human cell lines look interesting.

Another recent method which investigated decision-tree-based algorithms is due to Martens et al. (2020) [28]. Here, gradient boosted decision trees, logistic regression and neural networks were used to predict highly interacting chromatin regions and TAD domain boundaries from histone modifications and CTCF ChIP-seq data. However, in this setup, the neural network approach yielded the best, overall acceptable results of the three approaches, but again in form of a binary classifications. A neural-network approach with comparable input data, but without the restriction to binary classifications has been presented by Farrè et al. in 2018 [29]. Here, a one-dimensional convolutional filter is used to convert ChIP-seq data from 50 chromatin factors into a one-dimensional chromatin vector, which is then processed by a dense neural network (DNN). This allows predicting real-valued Hi-C interaction matrices, which resemble the general structure of experimentally derived matrices quite well.

While the approaches discussed so far have either modeled DNA as a “beads-on-a-string” polymer or not used it explicitly at all, there are also several machine-learning approaches which directly consider DNA sequence without a need for polymer modeling. In 2019, Singh et al. presented SPEID [30], an approach to predict promoter-enhancer interactions from DNA sequence, using a combination of CNNs, a recurrent network (LSTM) and a DNN. The results match well with experimental data, but are limited to promoter and enhancer loci by design, disallowing predictions of complete Hi-C contact matrices. Other researchers have tried to design similar methods without such limitations. The work by Peng from 2017 [31] is an extension of SPEID, based on a 2016 preprint [32], additionally taking into account a “middle sequence” between enhancer- and promoter sequences, CTCF motif counts within the sequences and genomic distance between two sequence snippets. However, the network lacks generalization, i. e. the results are only good in training regions [31, figs. 4, 5]. A conceptually similar method to the one by Peng [31], but with a different neural network design has been presented by Schreiber et al., also in 2017, named Rambutan [33]. It accepts DNA-sequence, DNase-seq data and distance between two genomic loci as inputs and then uses a combination of CNNs and a DNN to predict whether the given two loci interact or not. Unfortunately, it is difficult to decide whether the results of Schreiber et al. are useful for the task at hand, since the evaluation is done only by statistical means and no actual Hi-C matrices have been published. The original paper [33] also contains a known error and seems not to have appeared in a peer-reviewed journal in improved form yet. A probably more promising method working on DNA sequence, Akita, has been published by Fudenberg et al. in 2020 [34]. It is based on two rather involved convolutional neural networks. While the first one, “trunk”, processes one-dimensional, one-hot encoded DNA sequence input through convolutional filters, the second one, “head”, converts one-dimensional representations to 2D, further processes the data with convolutional filters and enforces symmetry. Although Fudenberg et al. initially seemed to focus on determining the influence of DNA modifications on spatial structure [35], predicting complete Hi-C matrices is an integral part of their work, and a large number of images of Hi-C matrices from the test set has been published alongside the article. The predicted matrices often hardly look like experimental Hi-C matrices, but mostly still indicate highly interacting regions quite well.

A further method by Schwessinger et al. [36] also makes use of DNA sequence and additional epigenetic data for its predictions, but is conceptually different from the ones presented so far.

Here, ChIP-seq tracks are initially used to train a CNN on the relationship between sequence and the corresponding chromatin factors. The weights of this first network are then used to seed another convolutional neural network, which is responsible for predicting DNA-DNA interactions from DNA sequence. In this case, the results were blurry, but generally in good accordance with experimental Hi-C data.

Yet another machine-learning concept based on sequence data, Samarth, has been proposed by Nikumbh and Pfeifer in 2017 [37]. Here, a support vector machine is trained on 5C data, using a specific representation for DNA sequence called oligomer distance histograms. The results showed acceptable consensus with experimental Hi-C data and allowed some interesting conclusions about the importance of certain k-mers for DNA folding. However, the applicability for the task at hand is hard to assess, because none of the predicted matrices used for statistical evaluation has been published alongside the paper.

2.2 Image synthesis techniques from computer vision

With the advent of deep learning, both the number of opportunities and the demand for using machine learning techniques in image synthesis tasks have risen, as recently summarized by Tsirikoglou et al. [38]. Since Hi-C contact matrices can be seen as square grayscale images, such techniques can also be relevant for this thesis.

Probably one of the first applications of computer vision methods for Hi-C matrices was presented by Liu and Wang in 2019 [39]. Here, deep convolutional neural networks – modified from established image super-resolution networks ~~XXX~~ – have been used to enhance low-resolution Hi-C matrices.

Another technique from computer vision that has been transferred to Hi-C matrices are conditional generative adversarial networks (GANs), invented by Goodfellow, Mirza and colleagues in 2014 [40, 41]. Again, several works have employed this comparatively new and involved method to increase the resolution of given Hi-C matrices, including the ones by Liu and colleagues [42], Hong et al. [43] and Dimmick et al. [44]. In general, all three works feature the typical GAN setup with two competitive networks – a generator, which is trained to produce realistic high-resolution Hi-C matrices from its low-resolution inputs, and a discriminator, which tries to distinguish real Hi-C matrices from generated ones and thus serves as a “critique”, an additional loss function, for the generator. While the convolutional building blocks for the discriminators and the residual building blocks for the generators are conceptually similar in all cases, the three works differ in the number of building blocks used and the activation functions applied within the blocks. Furthermore, Dimmick et al. and Hong et al. include additional loss functions beyond standard generator- and adversarial losses. This includes perceptual losses derived from other (pre-trained) neural networks to obtain “visually good” Hi-C matrices and total variation loss to suppress noise while maintaining edges. The method by Dimmick et al. outperformed the others for most test cases, but it is also the most elaborated.

2.3 Discussion of existing work

Independent of chosen techniques, several of the methods shown above only allow predictions restricted to certain loci (e.g. promoters and enhancers) or yield binary predictions in the sense

of “interaction” or “no interaction” between certain loci [25, 26, 28, 30]. These methods are thus not directly suitable for the task at hand, but would require further development.

The chromatin-modeling based approaches [19, 20, 21, 22] allow indirect derivation of Hi-C matrices by performing sufficient simulation passes and estimating contacts from the resulting model ensembles. Depending on the chosen chromatin model – which seems not straightforward [45, 46] – and the number of constraints involved, the required computations can be expensive. However, the results still seem slightly inferior compared to other methods mentioned in section 2.1, maybe because not all constraints for chromatin modeling are known or cannot be considered in the models yet.

Three of the DNA-sequence based methods mentioned above also allow direct prediction of Hi-C matrices [33, 34, 36]. At first look, it is surprising that learning spatial structure from DNA-sequence actually works, because there seems no obvious correlation between sequence and structure – instead, 3D conformation can vary considerably for different cell lines of the same organism, which all share the same DNA. On the other hand, the chosen artificial networks might be able to figure out binding sites of relevant proteins from DNA sequence, which *do* have a correlation with spatial structure. This is especially true for the method by Schwessinger [36], where the network is seeded by training on exactly such binding sites. While the methods by Schreiber and Schwessinger use secondary inputs and therefore can – at least partially – adapt to cell lines sharing the same DNA, the method by Fudenberg focuses on DNA and is thus expected to produce the same outputs for all cell lines of a given organism. All three methods require comparatively deep networks which are expensive to train. The sequence-based method by Nikumbh et al. [37] is using a support vector machine, which is generally easier to train, but the adaptability to different cell lines is expected to remain problematic due to the chosen concept.

The random-forest-based method by Zhang et al. [27] is directly targeted at the goal of this thesis, since it directly predicts Hi-C matrices, is not limited to certain loci and can adapt to different cell lines using corresponding ChIP-seq data. However, this approach has extensively been investigated in two previous study projects at the university of Freiburg, and the results could not be reproduced for unknown reasons [13, 47].

The dense neural network approach by Farrè et al. [29] is also compatible with the goals of this thesis. The published results are visually and statistically convincing, and both the presented network and the training process offer opportunities for amendments, which will be explored in section 3.1.

Since all of the image synthesis methods presented above in section 2.2 require existing Hi-C matrices for training *and* prediction, none of them is directly suitable for the task at hand. However, some of the concepts can still be used to develop novel approaches, which will be discussed in section 3.2.

3 Advancing predictions of Hi-C interaction matrices

In the following subsections, two conceptually different approaches towards the goal of the thesis, predicting Hi-C matrices from ChIP-seq data, will be explored. While the first approach is a dense neural network based on work by Farré et al. [29], the second is a novel method based on conditional generative adversarial networks.

3.1 Dense Neural Network approach

In their 2008 paper [29], Pau Farré, Alexandre Heurtau, Olivier Cuvier and Eldon Emberly propose a combination of a 1D convolutional filter with a three-layer dense neural network which already fulfills most goals of this thesis with some exceptions regarding data formats and preprocessing. This thesis tries to build on the success of their method by extending the comparatively simple neural network in various ways, modifying the binsizes of the Hi-C matrices and changing the learning process. As a start, the basic network has been rebuilt and used on well-known data from human cell lines GM12878 and K562.

3.1.1 Basic network setup

The basic network setup taken over from Farré et al. [29] is shown in simplified form in fig. 3, see section 4.2.1 for technical details.

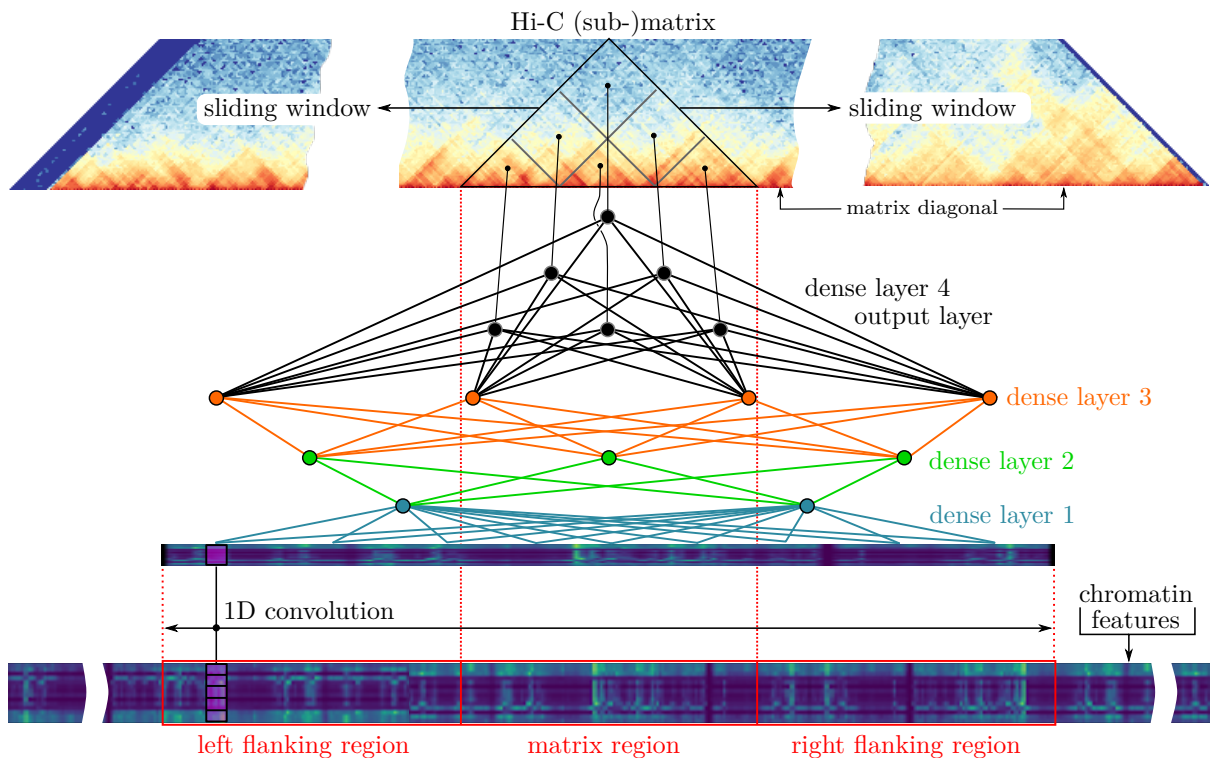


Figure 3: Principle of basic dense neural network

Since the network implements a supervised learning technique, it requires two kinds of inputs for training – ChIP-seq data of the chosen chromatin factors and target Hi-C matrices for each training chromosome. However, it is generally infeasible to learn Hi-C matrices for complete chromosomes at once with machine learning approaches, because the latter heavily depend on availability of training data.

The target matrices are thus taken as submatrices of size $w \times w$ with fixed window size w , centered at the diagonal of the original Hi-C matrices, fig. 3 top. The ChIP-seq data is taken as $3w \times n$ subarray of the original array, whereby the middle w bins are aligned with the position of the submatrix, the first w bins correspond to the left flanking region and the last w bins correspond to the right flanking region of the current submatrix region, fig. 3 bottom. Training samples are then obtained from the input data by sliding the input windows along the diagonal of the target Hi-C matrix. For the technical details of the sample generation process, see section 4.1.3.

Within the network, a 1D convolutional filter compresses the $3w \times n$ input arrays to 1D vectors of size $3w \times 1$, and four dense layers further process the compressed input, fig. 3 middle. The number of neurons in the last dense layer corresponds to the number of bins in the upper triangular part of the target submatrix, i.e. it consists of $(w \cdot (w + 1))/2$ neurons, exploiting the symmetry of Hi-C matrices. For implementation details, please refer to section 4.2.1.

Training of the network is performed by minimizing the mean squared error of the predicted matrix versus the target Hi-C matrix using stochastic gradient descent. The technical details are also given in section 4.2.1

Farré et al. propose a window size of $w = 80$ at $b_{feat} = b_{mat} = 10$ kbp, but larger binsizes of $b_{feat} = b_{mat} = 25$ kbp were found beneficial for the data used throughout this master thesis. The results are shown in section 5.1.1.

The network shown above is quite simple, and immediately offers some opportunities for expansion, partially already proposed in the original paper [29]. These will be explored below.

3.1.2 Modifying the convolutional part of the network

One starting point for modifying the neural network is its convolutional part.

With only a single 1D convolutional filter in one layer, the network might have difficulties capturing complex relationships between Hi-C interaction counts and more than one of the chromatin features. For this reason, an extended “longer” network was created, comprising three 1D convolutional filter layers with 16, 8 and 4 filters, respectively. This is still a comparatively low number of layers and filters, but the choice seemed justified in order to avoid overfitting to the low-dimensional input.

Next, a “wider” network was created, featuring the same setup as the basic network except the width of the filter kernel, which was set to 4 instead of 1. The idea here is to allow the network to capture correlations between Hi-C interaction counts and chromatin features which span more than one bin. The actual number has been kept low, since at binsize $b = 25$ kbp, 4 bins already correspond to 100 kbp. Of course, increasing filter width and using more filters can also be combined, hopefully allowing the “wider-longer” network to capture both correlations spanning more than one bin and more than one chromatin feature.

Another approach to potentially improve the predictions that goes somewhat into the direction of the “wider” network has been proposed, but not implemented by Farré et al. in their paper [29]. ChIP-seq experiments can usually be binned at smaller binsizes than Hi-C data due to the nature of the process. This can be exploited to capture finer details in the ChIP-seq data without a need for higher (training-)matrix resolutions. To this end, the initial network can be generalized by binning the ChIP-seq data at k times the bin size of the matrices, whereby $k \in \mathbb{N}^{\geq 1}$, cf. section 4.1.5 for the technical details. This yields an input data size of $k \cdot 3w \times n$, which is then again compressed to a $3w \times 1$ vector by a 1D convolutional filter with kernel size n and strides k . For practical reasons, $k = 5$ was chosen for the thesis at hand, and the results for binsizes $b_{mat} = 25$ kbp, $b_{feat} = 5$ kbp are shown in XXX.

3.1.3 Using a combination of MSE-, TV- and perceptual loss

In image regression tasks, optimizing for mean squared error is known to produce blurry images, because it is computed independently for each image pixel, ignoring spatial proximity [48, 49]. Another approach to improve the predictions of the basic neural network is thus using a different loss function.

It has been shown that loss functions based on the (multiscale-)structural similarity index (SSIM) [50]. can outperform mean squared error (L2) and mean absolute error (L1) in image regression tasks. While Zhao et al. used a combination of L1- and multiscale SSIM loss [51], Lu proposed a custom level-weighted SSIM loss [49]. The results were better than with L1- or L2 loss alone, but sometimes not much – depending on the machine learning model in use.

Another type of loss function used in image processing is the so-called perceptual- or perception loss. The idea here is not to compute L1 or L2 loss directly from the output of the network to be trained, but instead use a pre-trained loss network to determine structures in “predicted” and “real” images and then compute e.g. L1 or L2 loss on these structures, fig. 4. The optimization

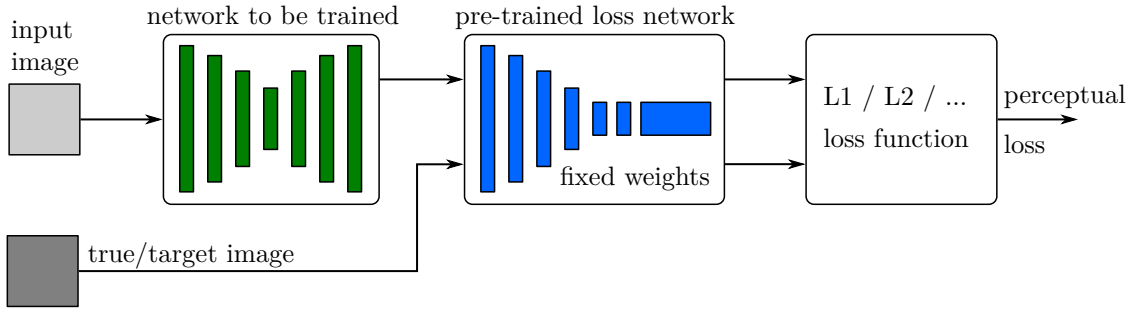


Figure 4: perceptual loss

of the trainable weights can be performed as usual, for example with gradient descent and backpropagation, simply keeping the weights of the loss networks constant. Often, complex image classification networks like VGG-16 [52] are taken as loss network, e.g. in the well-known style-transfer network by Johnson et al. [53], because these are known to be good at detecting relevant structures in images.

To check whether the given learning task benefits from using a perceptual loss, a custom combined loss function was generated, consisting of mean squared error l_{MSE} between true- and

predicted matrices, perceptual loss l_{VGG} based on VGG-16 and total variation loss l_{TV} to reduce noise in the output while preserving edges [54]. This choice was inspired by the custom loss function used by Hong et al. in their successful Hi-C super-resolution network DeepHiC [43], which is otherwise not similar to the network used here. In short, the combined loss function L is shown in eq. (1), see section 4.2.4 for details,

$$L_{combined} = \lambda_{MSE} l_{MSE} + \lambda_{VGG} l_{VGG} + \lambda_{TV} l_{TV} \quad (1)$$

with individual loss weights λ .

Unfortunately, there is no straightforward way for determining the optimal parameters λ , and an exhaustive parameter search was infeasible due to the computation time requirements of about 4:30 min per epoch. Therefore, only few runs were conducted with different sets of parameters, and the results for $\lambda_{MSE} = 0.8999$, $\lambda_{VGG} = 0.1$, $\lambda_{TV} = 0.0001$, which should not be considered optimal, are shown in section 5.1.3.

3.1.4 Using a TAD-based loss function

Looking at the results obtained from the networks so far, see [XXX](#), it seemed that highly interacting regions, especially topologically interacting domains (TADs), were not well predicted and either absent in the matrix plots or blurred. Assuming availability of a TAD scoring function $tad(z)$, where z is a predicted submatrix, this might be improved by directly optimizing a loss function as shown in eq. (2).

$$L_{combined} = \lambda_{MSE} + \lambda_{TAD} tad(z) \quad (2)$$

However, this approach suffers from several restrictions. First and foremost, there seems no consensus on the exact definition of TADs, and no less than 22 algorithms for TAD detection existed as of 2018 [55, 56]. Additionally, many of these algorithms have several tuning parameters, are notoriously parameter-dependent and may not even yield any results if parametrized in an unfavorable way [56]. A further restriction results from the context – since the loss function needs to be optimized, one needs to be able to compute gradients of it with respect to the network’s weights. Due to their complexity, this is generally very difficult to implement in a computationally efficient way for all known TAD calling algorithms.

To overcome the limitations, a novel loss function based on TAD scores [57] was developed. Figure 5 exemplarily shows its basic idea for a 16×16 submatrix with window size 4. First, the mean is taken from diamond-shaped (or rhombus-shaped) matrix cutouts along the diagonal, fig. 5 (A, C), and stored in a score vector, fig. 5 (D / E). The size of the diamonds, 2 in the figure, is configurable, but a reasonable balance with the submatrix size, i.e. the size of the sliding window w , and the expected size of TADs in the matrix must be maintained. Next, loss is computed by taking the mean squared error between the score vectors of the “real” submatrices and the “predicted” submatrices.

The idea behind the score-based approach is visualized in fig. 5 (E). Local minima, i.e. dents in the line plot of the score values, fig. 5 (F), often correspond with highly interacting regions in the matrix, since the mean of diamonds from *inside* TADs is normally significantly higher than the mean of diamonds *outside* TADs. Indeed, some TAD calling algorithms like TopDom [58] and hicFindTADs [59, W12f.] do compute insulation scores – usually for more than one

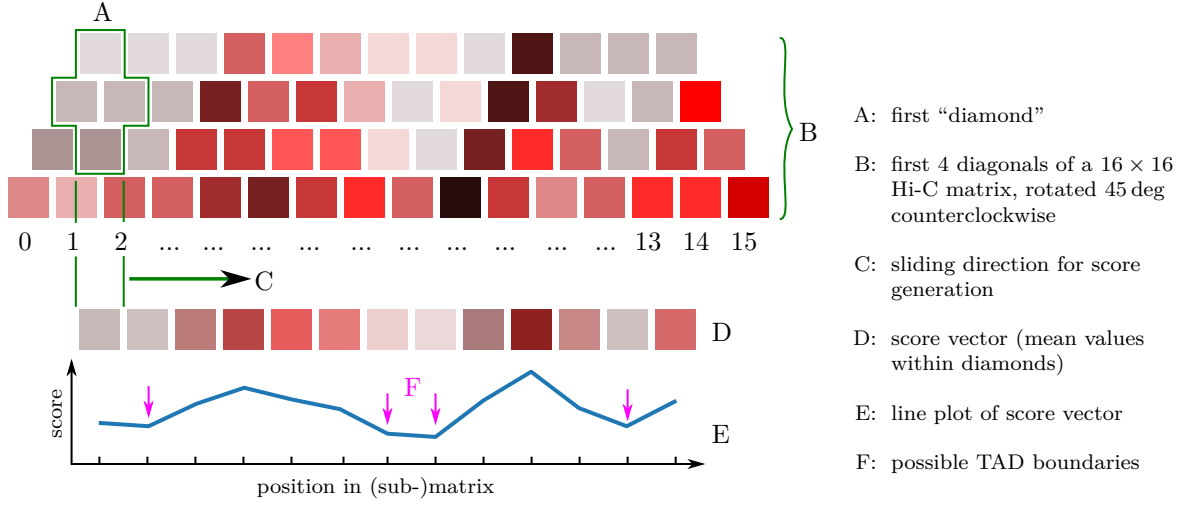


Figure 5: score vector generation for TAD-based loss

diamond size – and then use diverse techniques to detect meaningful local minima in the score values. However, finding meaningful local minima in the given context is still computationally involved, so it was left to the network to make sense out of the score vectors. This way, the loss function was well defined, efficiently computable and tensorflow standard functionality could be used to compute gradients with respect to weights.

For the thesis at hand, window size $w = 80$ and diamond size $ds = 12$ were used at binsizes of 25 kbp. The score-based loss was always used in combination with mean squared error, as shown above in eq. (2). For technical details please refer to section 4.2.5.

3.1.5 Modifying binsize and window size

- rationale: smaller binsizes - predict smaller structures, larger binsizes - predict larger structures
- then combine the predictions (e.g. sum them up, take the mean or even use a NN)
- use trapezoids, i.e. capped larger submatrices and flankingsize smaller than window size
- rationale: larger window size without increasing training time too much
- train on matrices/features with different binsizes in the same training run (e.g. two matrices with 25k, 50k and bin at 25k, 50k or 5k, 10k)

3.1.6 DNA sequence as an additional network input branch

- use DNA as an additional input
- rationale a): allow the network to figure out true binding sites in conjunction with cs data
- rationale b): given the success of pure DNA based methods, allow the network to find yet unknown sequence structure correlations

- probably not the most important subsection, leave it out in case of time problems

3.2 Hi-cGAN approach

3.2.1 Pix2Pix as a generic generative network

Describe the basic setup of pix2pix and what it is currently used for Pix2Pix using 2D input, but we have n times 1D. Bad. 2 ideas how to come around, using DNN like Farre et al, potentially pretrained, or using new CNN. Maybe mention restriction of 2^x for image size

3.2.2 Using a DNN for 1D - 2D conversion

- use network like Farre et al for converting inputs to 2D
- rationale: this nw has been shown to create comparatively good Hi-C matrices
- rationale: pix2pix could then improve them
- rationale: transfer learning can be employed

3.2.3 Using a CNN for 1D - 2D conversion

- use somewhat deep CNN for conversion
- rationale: the DNN approach didn't work
- rationale: downsampling like in image processing, just in 1D

4 Methods

4.1 Input data

For the thesis at hand, data from human cell lines GM12878, K562, HMEC, HUVEC and NHEK was used. The exact data sources and data preprocessing will be outlined in the following subsections 4.1.1 and 4.1.2.

4.1.1 Hi-C matrices

Hi-C data due to Rao et al. [11] was downloaded in .hic format from Gene Expression Omnibus under accession key GSE63525. Here, the quality-filtered “combined_30” matrices were taken, which contain only high-quality reads from both replicates.

Next, matrices at 5 kbp binsize were extracted and converted to cooler format using `hic2cooler` and subsequently coarsened to resolutions of 10, 25, 50 and 100 kbp using `cooler coarsen`.

Contrary to the work from Farré et al. [29], which is using a distance-based normalization, and many others in the field which are using ICE- or KR-normalization, these matrices have not been normalized for the thesis at hand because no benefit of doing so was found during the study project [13].

4.1.2 ChIP-seq data

For this thesis, ChIP-seq data for 13 chromatin features and DNaseI-seq data was used, cf. table 2. Here, ChIP-seq data for the 13 features was downloaded in .bam format from the ENCODE project XXX for both replicate one and two, and DNaseI-seq data was downloaded in .bam format from XXX; the download links are also given in table 2.

feature name	download link
CTCF	
DNaseI	
H3k27ac	
H3k27me3	
H3k36me3	
H3k4me1	
H3k4me2	
H3k4me3	
H3k79me2	
H3k9ac	
H3k9me3	
H4k20me1	
Rad21	
Smc3	

Table 2: chromatin features used for the thesis

The data were then converted to bigwig format, which is more convenient to handle, and the replicates were merged into one bigwig file. Pseudocode for the full conversion process is given in [XXX](#).

4.1.3 Sample generation process for the dense neural network

This thesis follows the sliding window approach proposed by Farré et al. [29].

First, all chromatin features were binned to binsize b_{feat} by splitting each chromosome of size cs into $l_{feat} = \left\lceil \frac{cs}{b_{feat}} \right\rceil$ non-overlapping bins of size b_{feat} and taking the mean feature value within the genomic region belonging to each bin. All n chromatin factors were processed in this way and then stacked into a $l_{feat} \times n$ array.

Hi-C matrices were taken as provided by the cooler format, cf. section 4.1.1, i. e. as $(l_{mat} \times l_{mat})$ -matrices, $l_{mat} = \left\lceil \frac{cs}{b_{mat}} \right\rceil$ being the number of bins in the given chromosome. Initially, $b_{feat} = b_{mat}$ was used, which leads to $l_{feat} = l_{mat}$, because cs is a constant for a given chromosome.

A sliding window approach was now applied to generate training samples for the neural networks G described below. Here, subarrays of size $(3w_{mat} \times n)$ were cut out from the feature array such that the i -th training sample corresponded to the subarray containing the columns $i, i + 1, i + 2, \dots, i + 3w_{mat}$ of the full array. Sliding the window along the array with stepsize one obviously yields $N = l - 3w_{mat} + 1$ training samples. The corresponding Hi-C matrices were then cut out along the diagonal of the original matrix as submatrices with corner indices $[j, j], [j, j + w_{mat}], [j + w_{mat}, j + w_{mat}], [j + w_{mat}, j]$ in clockwise order, where $j = i + w_{mat}$. The idea here is that the first $0, 1, \dots, w_{mat}$ columns of each feature sample form the left flanking region of the training matrix, the next $w_{mat} + 1, w_{mat} + 2, \dots, 2w_{mat}$ correspond to the matrix' region and the last $2w_{mat} + 1, 2w_{mat} + 2, \dots, 3w_{mat}$ columns form the right flanking region. Since Hi-C matrices are symmetric by definition, only the upper triangular part of the submatrices was used, flattened into a $w \cdot (w + 1)/2$ vector.

Figure 6 exemplarily shows the sample generation process for a (16×16) -matrix with $w_{mat} = 4$ and $n = 3$ chromatin features. In this case, five training samples would be generated – the one encircled in green and four more to the right, as the window is sliding from left to right until the right flanking region hits the feature array boundary.

The sample generation process for predicting (unknown) matrices was the same as for training, except that no matrix window was generated. Due to the sliding window approach, the output of the network is a set of overlapping submatrices along the main diagonal of the actual target Hi-C matrix. To generate the final submatrix, all submatrices were added up in a position-aware fashion and finally all values were divided by the number of overlaps for their respective position. Figure 7 exemplarily shows the prediction process for $N = 5$ samples with window size $w = 4$ for a 16×16 matrix. Note that the first and last w bins in each row (matrix diagonal) always remain empty due to the flanking regions, as do all bins outside the main diagonal and the first $w - 1$ side diagonals.

The network and its learning algorithm was implemented in python using tensorflow deep learning framework with keras [XXX](#)citation, see [XXX](#)repository. The basic network was used with parameters $n = 14$ (cf. section 4.1.2), $w = 80$ and $b_{feat} = b_{mat} = 25$ kbp.

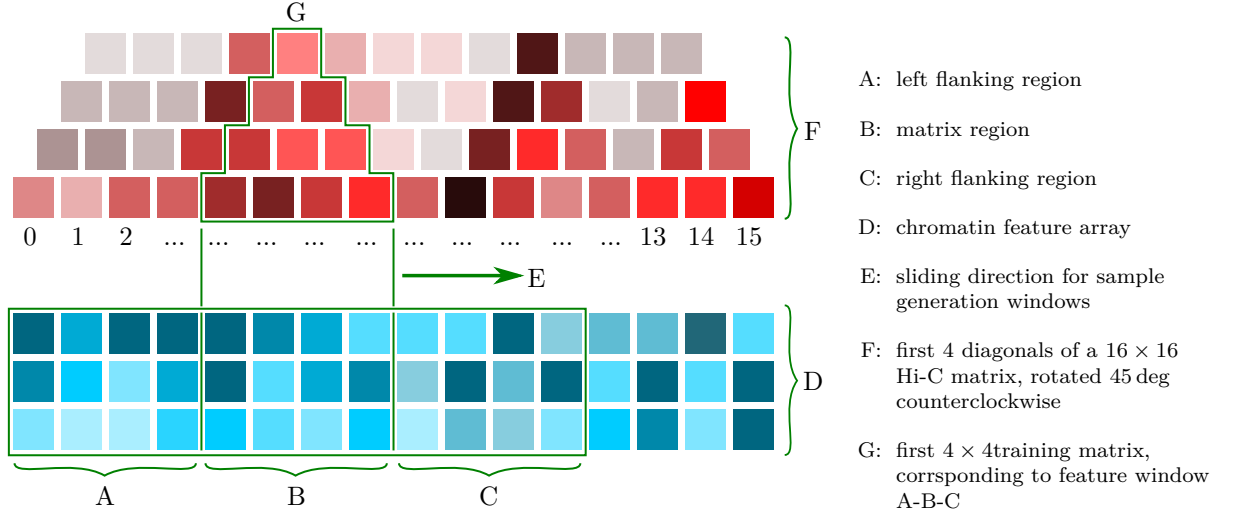


Figure 6: Sample generation process

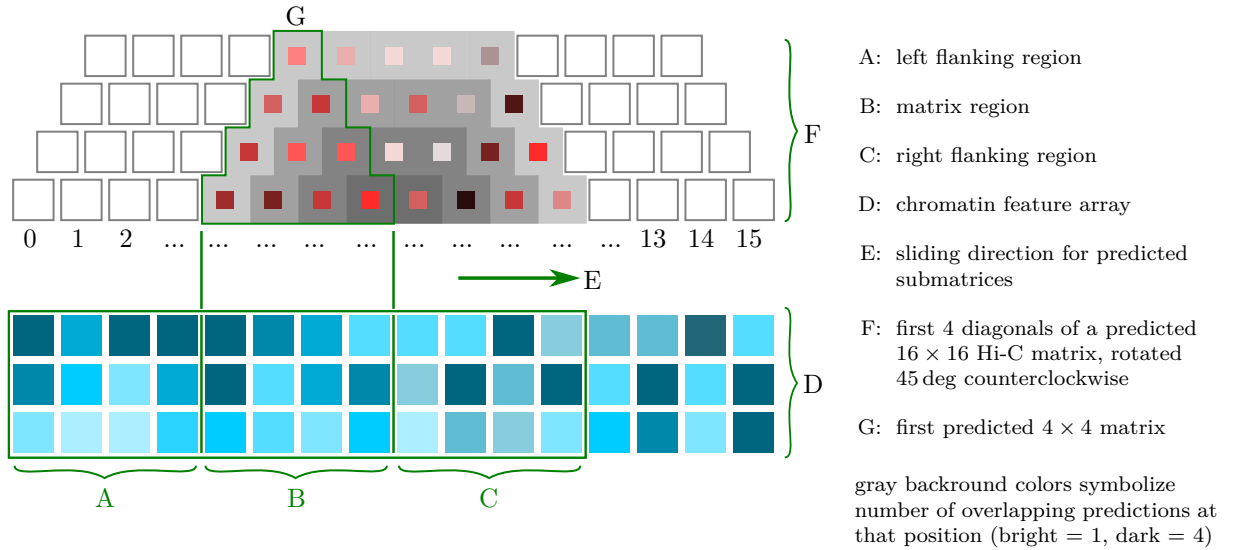


Figure 7: Prediction process

Training samples were drawn from GM12878, chromosomes 1, 2, 4, 7, 9, 11, 13, 14, 16, 17, 18, 20 and 22 of GM12878, validation samples from GM12878, chromosomes 6, 8, 12 and 15 and test samples from K562, chromosomes 3, 5, 10, 19 and 21, see table XXX for details. Test samples were processed through the network and the resulting submatrices – actually, vectors corresponding to their upper triangular parts – were subsequently re-assembled into full Hi-C matrices. Finally, the output matrices were scaled to value range [0...1000] for better visualisation and then stored in cooler format.

4.1.4 Sample generation for the cGAN

The samples for the cGAN were generated the same way as the samples for the dense neural network described in section 4.1.3, with two exceptions. First, Hi-C (sub-)matrices were not entered as vectors corresponding to their upper triangular part, but were instead taken as $w \times w \times 1$ grayscale images with value range [0...1] in 32-bit floating point format. Second, for the cGAN approach the input matrices need not only be square, but their sizes also needed to be powers of two. This ensured the required shapes for the connected up- and downsampling operations in the generator, which essentially are 2D-convolutions and transposed 2D-convolutions with strides two. Within the thesis, window sizes of $w = \{64, 128, 256\}$ were used, see XXX.

4.1.5 Generalization of feature binning

Most of the binning- and sample generation procedures described above also work for bin size relations $k = \frac{b_{mat}}{b_{feat}} \in \mathbb{N}^{>1}$.

The training matrices remain unchanged, i. e. $(l \times l)$ -arrays, from which training submatrices of size $w_{mat} \times w_{mat}$ can be extracted. With $k \in \mathbb{N}^{>1}$, one bin on the matrix diagonal corresponds to k bins of the feature array, so the feature window size needs to be k times the submatrix window size, $w_{feat} = k \cdot w_{mat}$. Since the first layer of all neural networks used in this thesis is a 1D convolution, this can be achieved by setting the filter width and strides parameters of the (first) convolutional layer to k , leaving the rest of the network unchanged. However, the number of bins along the matrix diagonal is generally not k times the number of bins in the feature array, see eq. (3).

$$l_{feat} = \left\lceil \frac{cs}{b_{feat}} \right\rceil = \left\lceil \frac{cs}{k \cdot b_{mat}} \right\rceil \neq k \cdot \left\lceil \frac{cs}{b_{mat}} \right\rceil = k \cdot l_{mat} \quad (3)$$

For the training process, this discrepancy was resolved by simply dropping the last training sample, if the feature window belonging to it had missing bins. For the prediction process, the feature array was padded with zeros on its end. This procedure ensures that no errors are introduced into the *training* process by imputing values, but keeps the size of the *predicted* matrix consistent with the training matrix binsizes.

Figure 8 shows the generalized training process with a (16×16) -training matrix and $k = 2$. If $15 \cdot b_{mat} + 1 \leq cs < 15 \cdot b_{mat} + b_{feat}$ held for the chromosome size cs in the example, then the number of bins on the matrix diagonal would be $l_{mat} = 16$, while the number of chromatin feature bins would be $l_{feat} = 31 \neq 2 \cdot l_{mat}$. In this case, the 5th sample would be dropped for

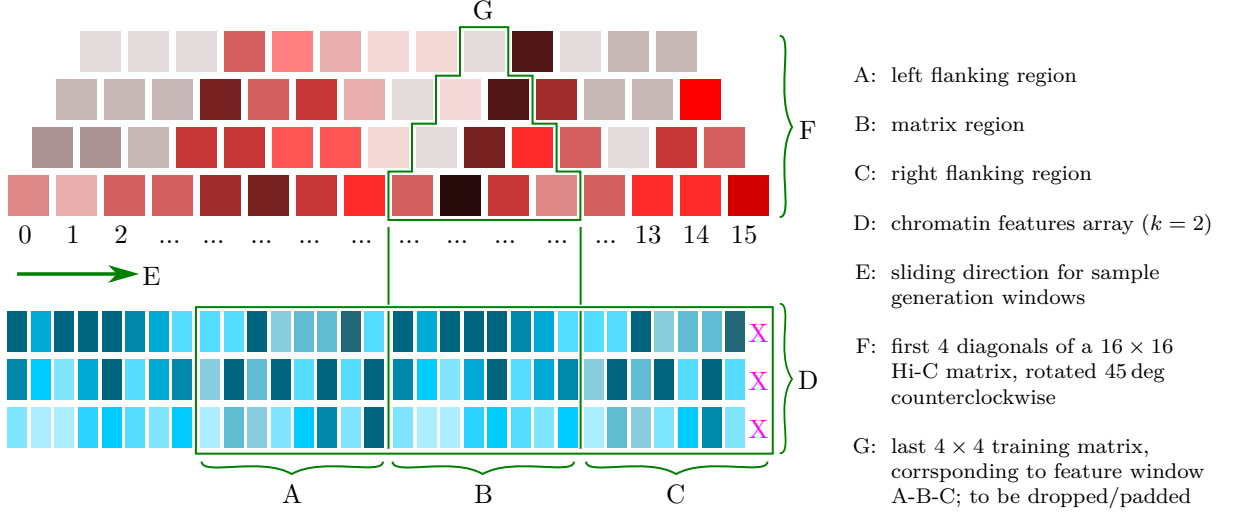


Figure 8: Generalized sample generation process f. $k = 2$

training, while a column with zero bins – symbolized by pink crosses in fig. 8 – would be added to the feature array for prediction so that the resulting matrix would still have the desired size of 16×16 .

4.2 Dense neural network approach

4.2.1 Basic setup

The basic setup for the dense neural network approach follows the proposal by Farré et al. [29]. It consists of five trainable layers; one convolutional layer and four densely connected layers with 7,486,436 trainable weights in total, fig. 9. For brevity, the sigmoid activation of the 1D convolution and the ReLU activations of the Dense layers are not shown. The dense layers all use a “L2” kernel regularizer and the dropout layers following all but the last dense layer have a dropout probability of 10%.

The training goal for the neural network G is to find weights ω^* for its five trainable layers such that the mean squared error L_2 between the predicted submatrices $M_s = G_s(\omega)$ and the training submatrices T_s becomes minimal for all N training samples $s \in (1, 2, \dots, N)$. Here, M_s is given by the activations of the last dense layer, which are to be interpreted as the upper triangular part of a symmetric $(w \times w)$ Hi-C matrix. Formally, one needs to optimize

$$\omega^* = \arg \min_{\omega} L_2(\omega) = \arg \min_{\omega} \frac{1}{N} \sum_{s=0}^N (M_s - T_s)^2 \quad (4)$$

For the thesis at hand, stochastic gradient with learning rate $1e^{-5}$ was used to find ω^* . Initial values ω_{init} were drawn from a Xavier initializer, a uniform distribution with limits depending on the in- and output shapes. Following [29], optimization was performed on minibatches of 30 samples, assembled randomly from the N training samples to prevent location-dependent effects and improve generalization. For training, the last batch was dropped, if $N/30 \notin \mathbb{N}$.

4.2.2 Modifying kernel size, number of filter layers and filters

For the “wider” variant, the kernel size of the of the 1D convolutional layer was increased to $4k$ with strides k , where $k = b_{mat}/b_{feat}$ is the relation between matrix- and feature binsize. Mirror-type padding was used to maintain the output dimensions of the basic network, which now had 7 486 478 trainable weights for $k = 1$.

For the “longer” variant three 1D-convolutional layers with 4, 8 and 16 filters were used in place of the basic network’s single convolutional layer. This replacement was also made for the “wider-longer”-variant, additionally increasing the respective kernel sizes to $4k$ (with strides k), 4 and 4. In both cases, the dropout rate was increased to 20%. The “longer” variant had 9 142 665 trainable weights and the “wider-longer” had 9 143 313 for $k = 1$.

For the variant with generalized feature binning, features were binned at $b_{feat} = 5$ kbp while keeping the matrix binsize $b_{mat} = 25$ kbp, so the factor for the relation between the two binsizes was $k = 25/5 = 5$. This yields input feature arrays of size $3w \cdot k \times n = 3 \cdot 80 \cdot 5 \times 14 = 1200 \times 14$. Replacing the first convolutional layer of the basic network by a 1D convolutional filter with kernel size $k = 5$ and strides $k = 5$, this input was again compressed to a $3w \times 1 = 240 \times 1$ tensor and fed into the flatten layer, cf fig. 9. The rest of the network remained the same so that the number of trainable parameters increased only slightly to 7 486 492

4.2.3 Custom loss function based on TAD insulation score

- computation details, custom layer
- full insulation score takes too long, simplified

4.2.4 Combination of mean squared error, perception loss and TV loss

For computing the combined MSE, perception- and TV-loss, input features were first passed through the network as normal, and mean squared error was computed on the predicted “upper triangular vectors” vs. the real vectors. Next, both the output- and target vectors were converted to symmetric grayscale images by embedding them into $w \times w \times 1$ tensors, where w is again the window size. TV loss was then computed by simply using the appropriate tensorflow function on the predicted images. For perception loss, the predicted images and the true images were first fed through a pre-trained VGG16 network truncated at layer **XXX**. The loss was then computed as mean squared error between the “predicted” and “true” activations at the output of that network.

Let $M_s = G_s(\omega)$ again be the output of the neural network G described above, and T_s the true matrices for training samples s in vector form, and let M'_s and T'_s be their grayscale image counterparts as described above. Furthermore, let $tv(x)$ be the total variation of image x and $vgg(x)$ the output of the perception loss network for image x . The goal of the modified network was now to find weights w^* such that

$$\omega^* = \arg \min_{\omega} (\lambda_{MSE} \frac{1}{N} \sum_{s=0}^N (M_s - T_s)^2 + \lambda_{TV} \sum_{s=0}^N tv(M'_s) + \lambda_{VGG} \frac{1}{N} \sum_{s=0}^N (vgg(M'_s) - vgg(T'_s))^2) \quad (5)$$

Weight initialization for network G and minibatching was done as described in section 4.2.1. The weights for the VGG16 network were taken from the pre-trained keras implementation and can here be considered as constants. In addition, the usage of VGG16 imposes the restriction $w \geq 32$ on the window size w , which is not problematic, since again $w = 80$ was chosen for all experiments.

The network $G(\omega)$ could in principle be any of the variants described above in section 4.2.2, but for the thesis at hand, only the initial network from section 4.2.1 was used.

4.2.5 Combination of mean squared error and TAD-score-based loss

Formally, the combination between mean squared error and TAD-score-based loss can be defined as the following optimization task:

$$\omega^* = \arg \min_{\omega} (\lambda_{MSE} \frac{1}{N} \sum_{s=0}^N (M_s - T_s)^2 + \lambda_{score} \frac{1}{N} \sum_{s=0}^N (score(M'_s, ds) - score(T'_s, ds))^2) \quad (6)$$

where M_s is again the Hi-C submatrix predicted by the network $G(\omega)$ for sample s and T_s is the corresponding true Hi-C submatrix.

To compute the TAD-score-based loss $score(\cdot, \cdot)$, the predictions and true Hi-C matrices M and T in vector form (upper triangular part) were first converted back to complete, symmetric Hi-C matrices M' , T' . Next, in a custom network layer, all diamonds with size ds inside the submatrices of size w were cut out using tensor slicing and the values inside the diamonds were reduced to their respective mean. This yields score vectors – more exactly, tensors with shape $(w - 2 * ds, 1)$. After computing the latter for both predicted- and real Hi-C submatrices, the mean squared error between them was computed as usual and weighted with a user-selected loss weight λ_{score} , see eq. (6).

While it would also have been possible, and probably faster, to slice the outputs of the networks directly in vector form, this is rather unintuitive and was therefore not done for the thesis.

4.3 HiC-GAN approach

4.3.1 Pix2Pix network layout

4.3.2 Using a DNN for 1D-2D conversion

4.3.3 Using a pre-trained DNN for 1D-2D conversion

4.3.4 Using a CNN for 1D-2D conversion

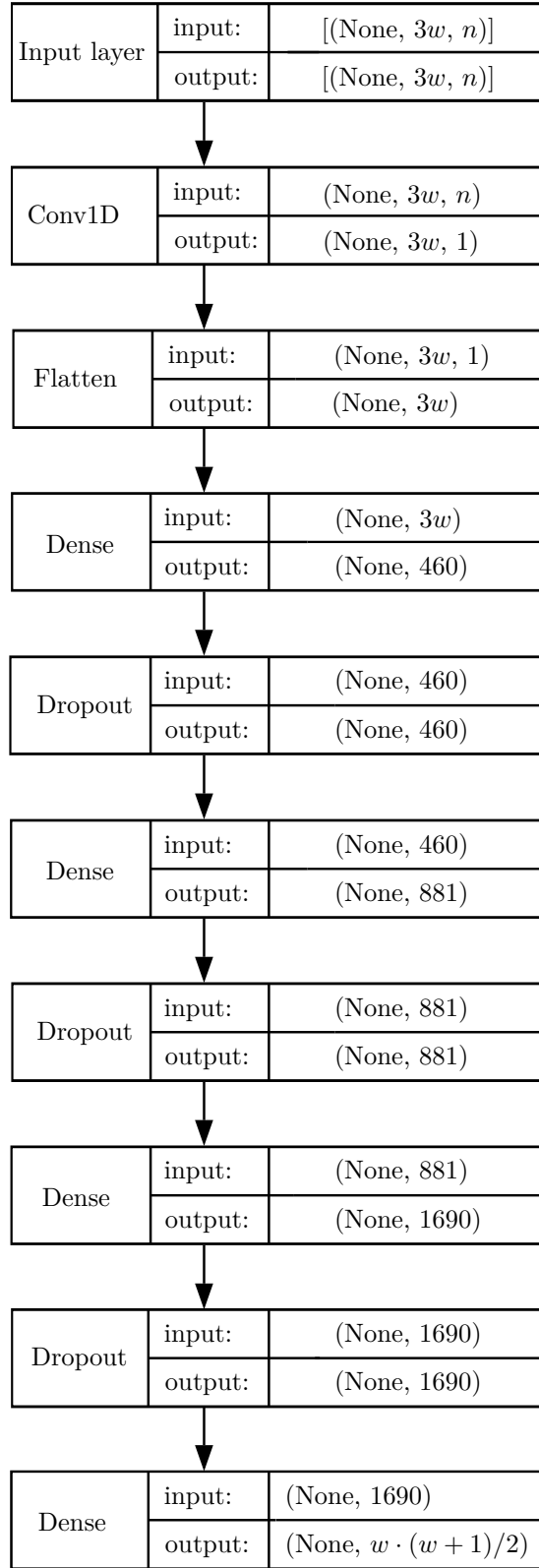


Figure 9: Basic dense neural network

5 Results

5.1 Dense Neural Network approaches

5.1.1 Initial results for comparison

The basic network was trained for 3000 epochs as explained in XXX. The validation error (MSE) for the basic neural network reached a minimum of about 150 000 after 500 epochs, fig. 10. Beyond that, the learning curve indicates overfitting, but differences between predictions at 500 and 1000 epochs are hard to spot in the matrix plots from the test set, fig. 13 and 14. To asses

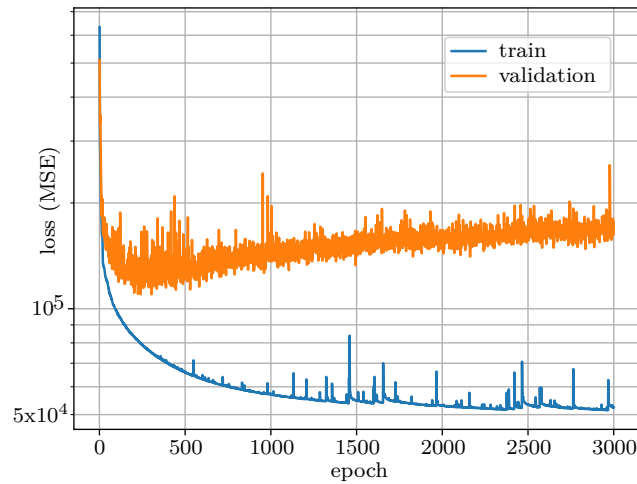


Figure 10: learning progress for basic DNN

the quality of the predictions, distance-stratified Pearson correlations were computed between predicted and real matrices, up to the maximum distance from the diagonal given by window size and bin size, i.e. for distances up to $80 \cdot 25 \text{ kbp} = 2 \text{ Mbp}$ in the basic setting. While the suitability of Pearson correlation as metric for seems debatable for Hi-C matrices XXX, it is quite common in the field and thus useful at least for comparisons. Figure 11 shows the Pearson correlations alongside area under the correlation curve (AUC) for the five test chromosomes at hand. The red lines in each correlation plot show the correlation between the corresponding training matrix from GM12878 and the target matrix from K562. It is obvious that all predicted test matrices have a strictly positive Pearson correlation, but are not better than transferring data from the training cell line.

The predicted matrices themselves looked modest when plotted with pygenometricks XXX. While some of the highly interacting regions, for example between XXX and XXX of chromosome 19 were well predicted, other structures, especially larger ones like the one between XXX and XXX in chromosome 19 or between XXX and XXX in chromosome 21 are not predicted at all.

5.1.2 Results for variations of the convolutional part

The predictions from the “wider” variant were generally similar to the initial results, both in terms of Pearson correlations and in terms of matrix plots, fig. 15 and XXX. Given the small

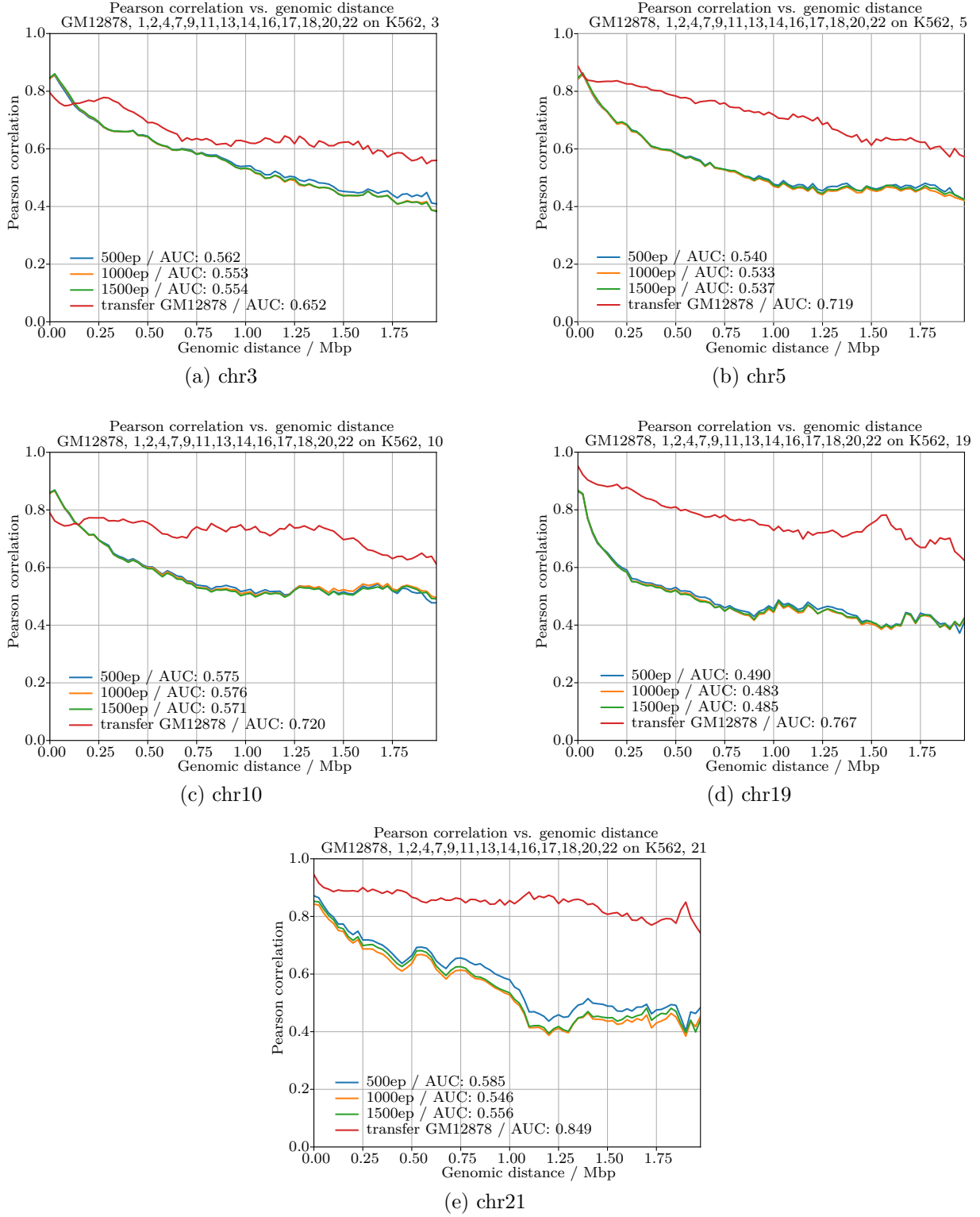


Figure 11: Pearson correlations, basic DNN, test chromosomes

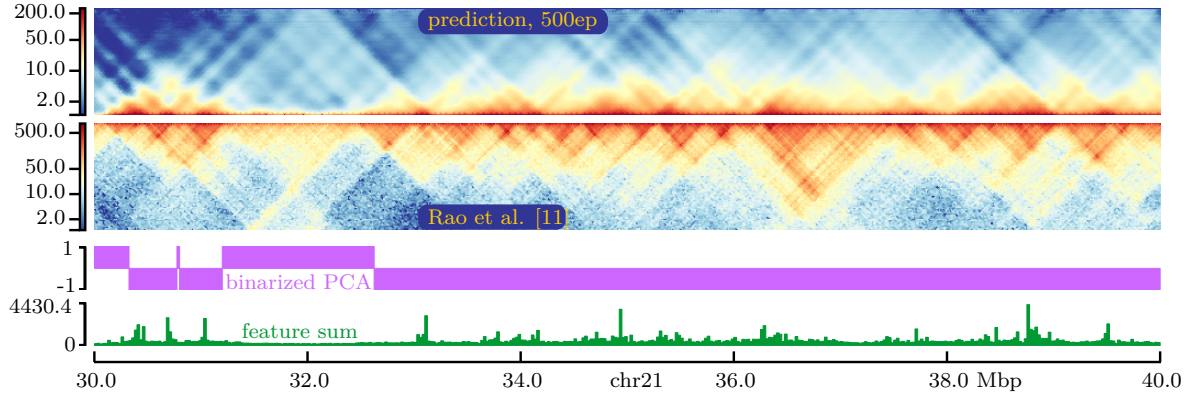


Figure 12: example prediction, 500 epochs

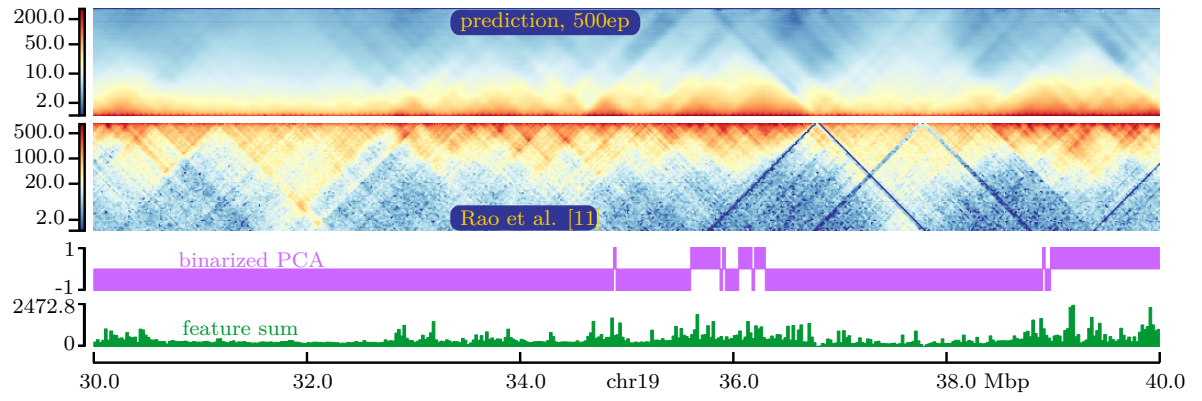


Figure 13: example prediction, 500 epochs

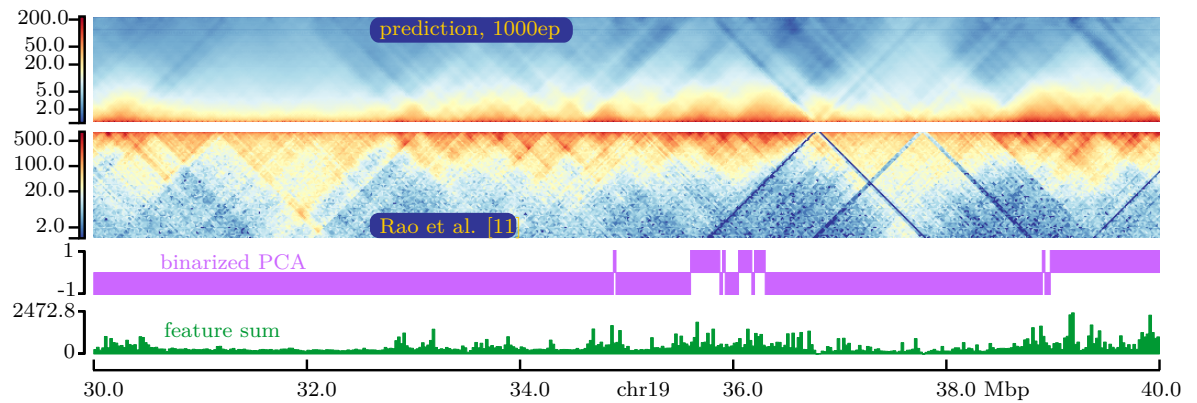


Figure 14: example prediction, 1000 epochs

increase in the number of trainable parameters at overall similar network topology, this is not surprising. However, it is interesting that hardly any improvement was found. Overfitting was less obvious than with the initial setup and the training process looked more smooth overall, but the remaining validation error was slightly higher than for the initial approach, fig. 16.

The predictions from the “longer” variant were partially better for the test set than the initial ones in terms of Pearson correlations, fig. 17. Interestingly, no predictions were available for certain distances after 250 and 500 epochs, while predictions for all distances were available after 1000 epochs. The reason for this behavior is unknown, but due to the network setup, comparatively few neurons are responsible for certain distances. Since the longer network setup has considerably more trainable parameters, 500 epochs might not be enough to fully adjust the weights of these (outer) neurons. The learning process in itself looked more smooth and reached a lower validation error than before, fig. 18.

The Pearson correlations for predictions from the “wider-longer” variant are shown in fig. 19. While improvements can be seen for 3 of 5 test chromosomes compared to the initial network, the results were worse than the ones from the highly similar “longer”-variant alone, and the remaining validation error was also higher.

The Pearson correlations for the variant with feature binsize 5 kbp and matrix binsize 25 kbp are shown in fig. 21. Much like the “wider” variant, the results did not improve compared to the initial predictions. The learning curve was smooth and showed signs of slight overfitting beyond 300 epochs, fig. 22.

5.1.3 Results for variations of the loss function

The Pearson correlations for a combined loss function according to eq. (5) with weighting parameters $\lambda_{MSE} = 0.8999$, $\lambda_{VGG} = 0.1$, $\lambda_{TV} = 0.0001$ are shown in fig. 23. For all test chromosomes, the results were highly similar to the initial network’s results.

While manually finding parameters λ for MSE- and VGG-loss that made the results better was not successful, it was found that the TV loss weight needed to be much smaller than the two other weights. Otherwise, many true interactions outside the first few matrix diagonals were considered as noise and optimized away early in the training process. XXXmaybe put one figure here from 2020-12-05_tvLoss

XXXRuntime score based loss about 7:00 per epoch (excl. validation)

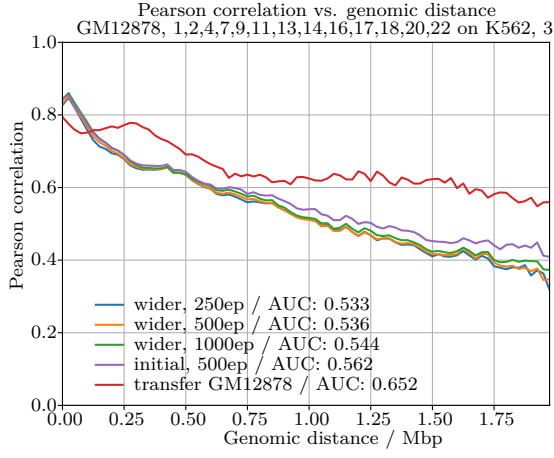
5.1.4 Results for different binsizes and windowsizes

5.2 Hi-cGAN approaches

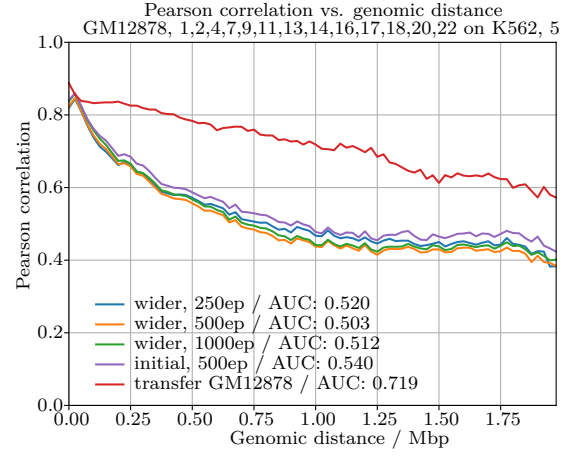
5.2.1 Initial results

The results at window size 64 bins were encouraging XXX

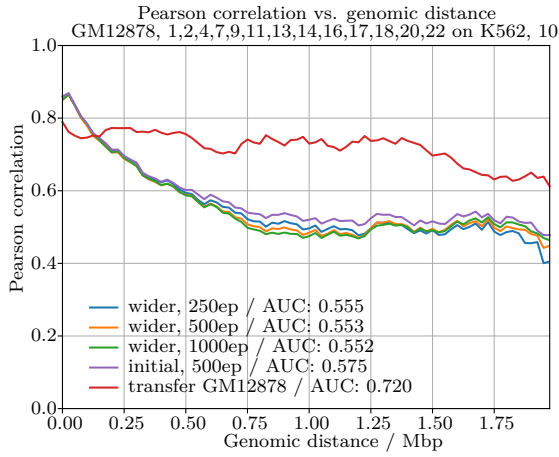
The results at 128k were even better XXX



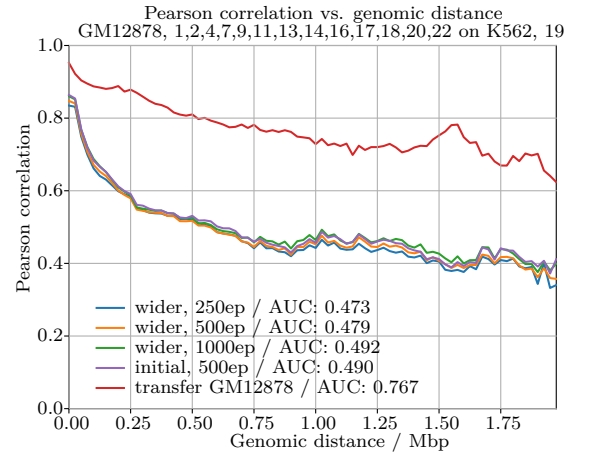
(a) chr3



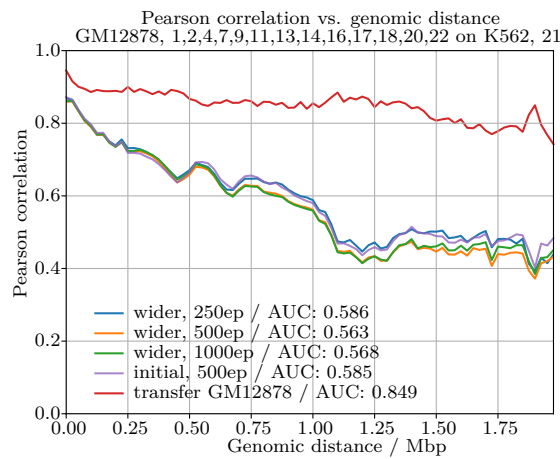
(b) chr5



(c) chr10



(d) chr19



(e) chr21

Figure 15: Pearson correlations, “wider” variant of DNN, test chromosomes

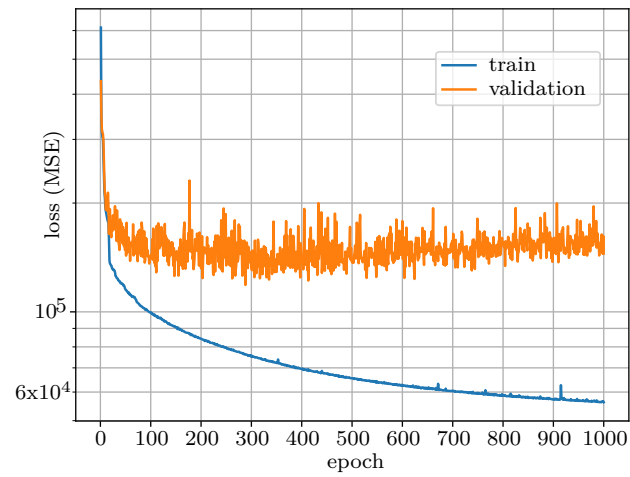


Figure 16: learning progress for wider DNN

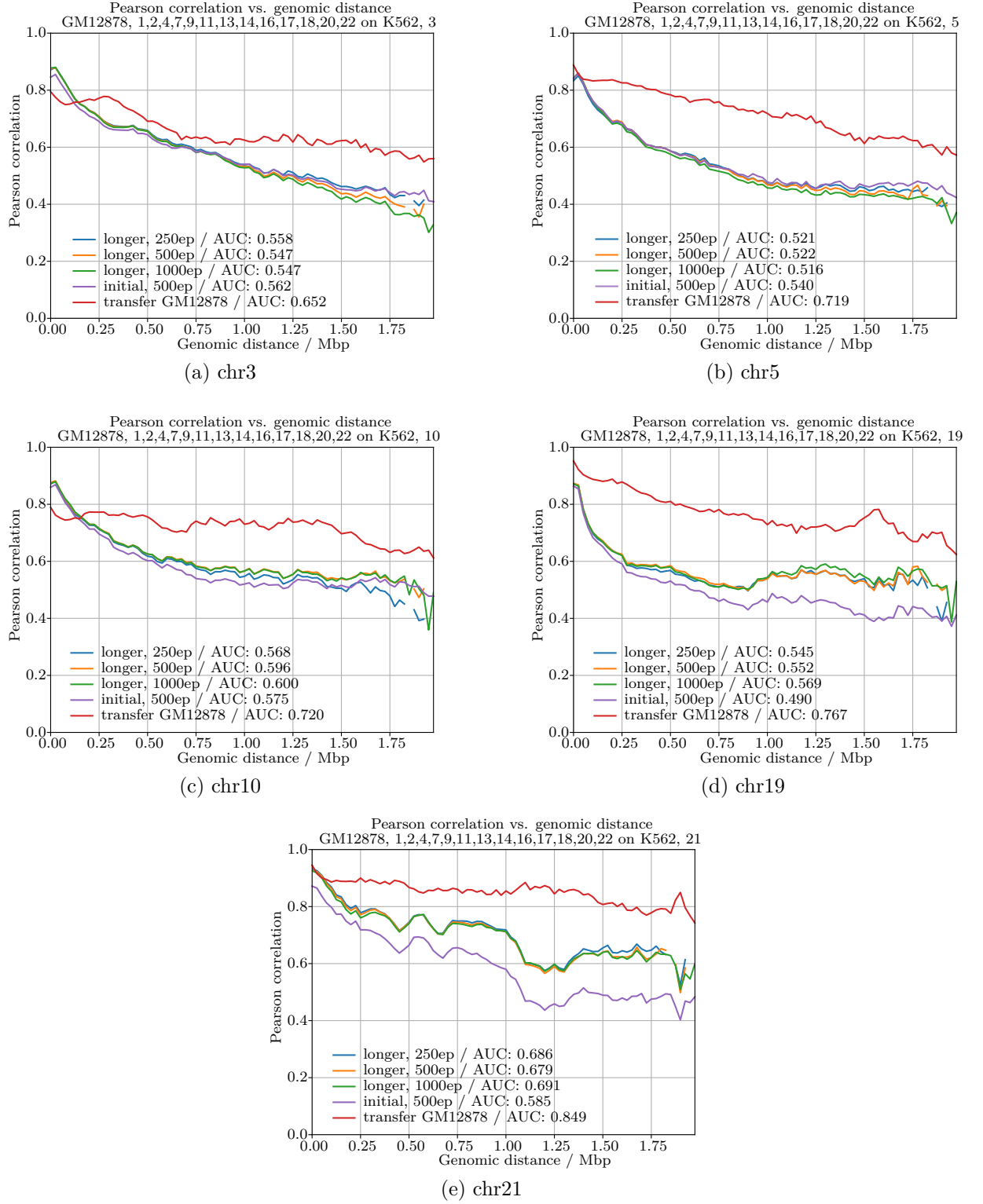


Figure 17: Pearson correlations, “longer” variant of DNN, test chromosomes

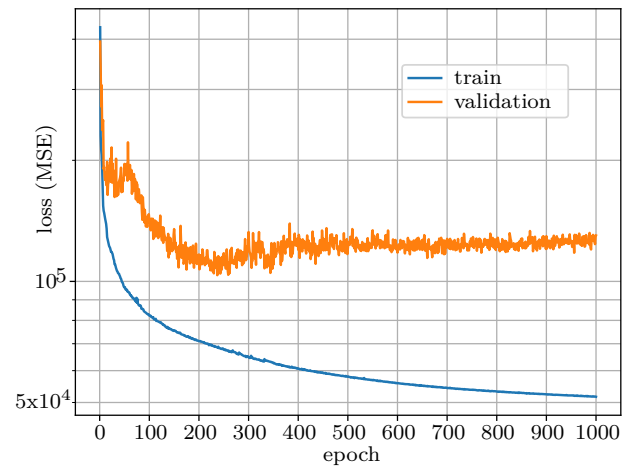


Figure 18: learning progress for longer DNN

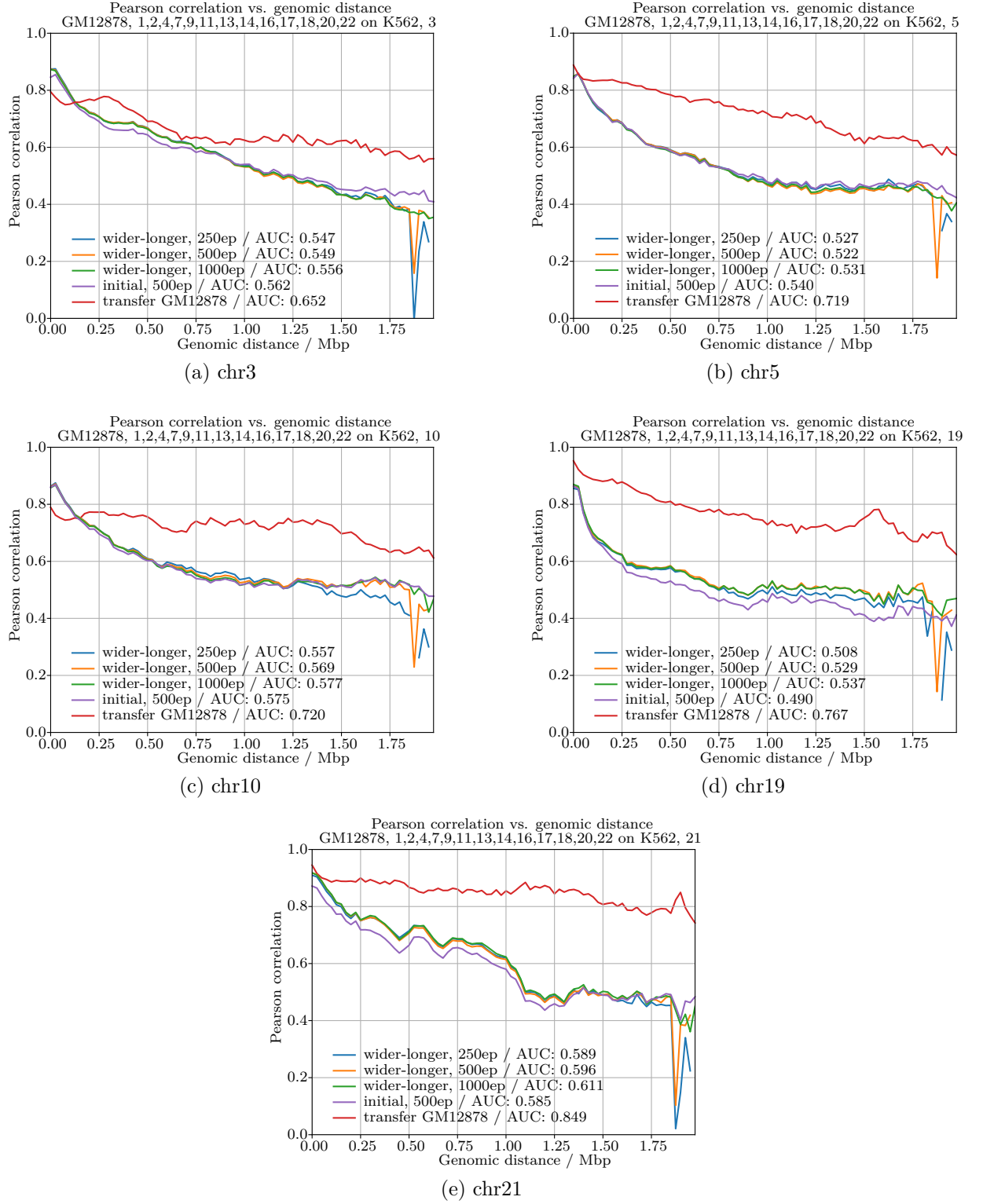


Figure 19: Pearson correlations, “wider-longer” variant of DNN, test chromosomes

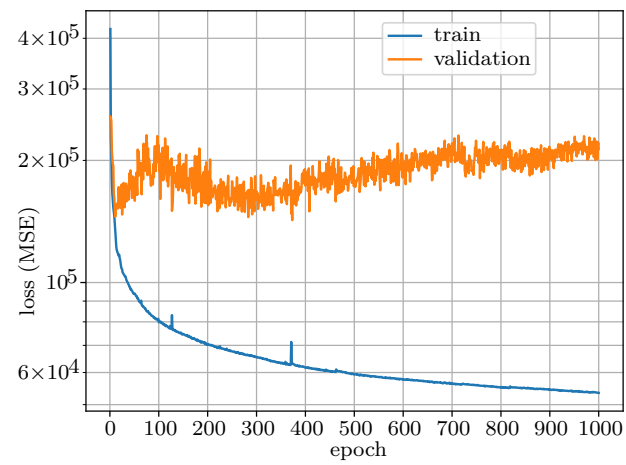


Figure 20: learning progress for “wider-longer” DNN

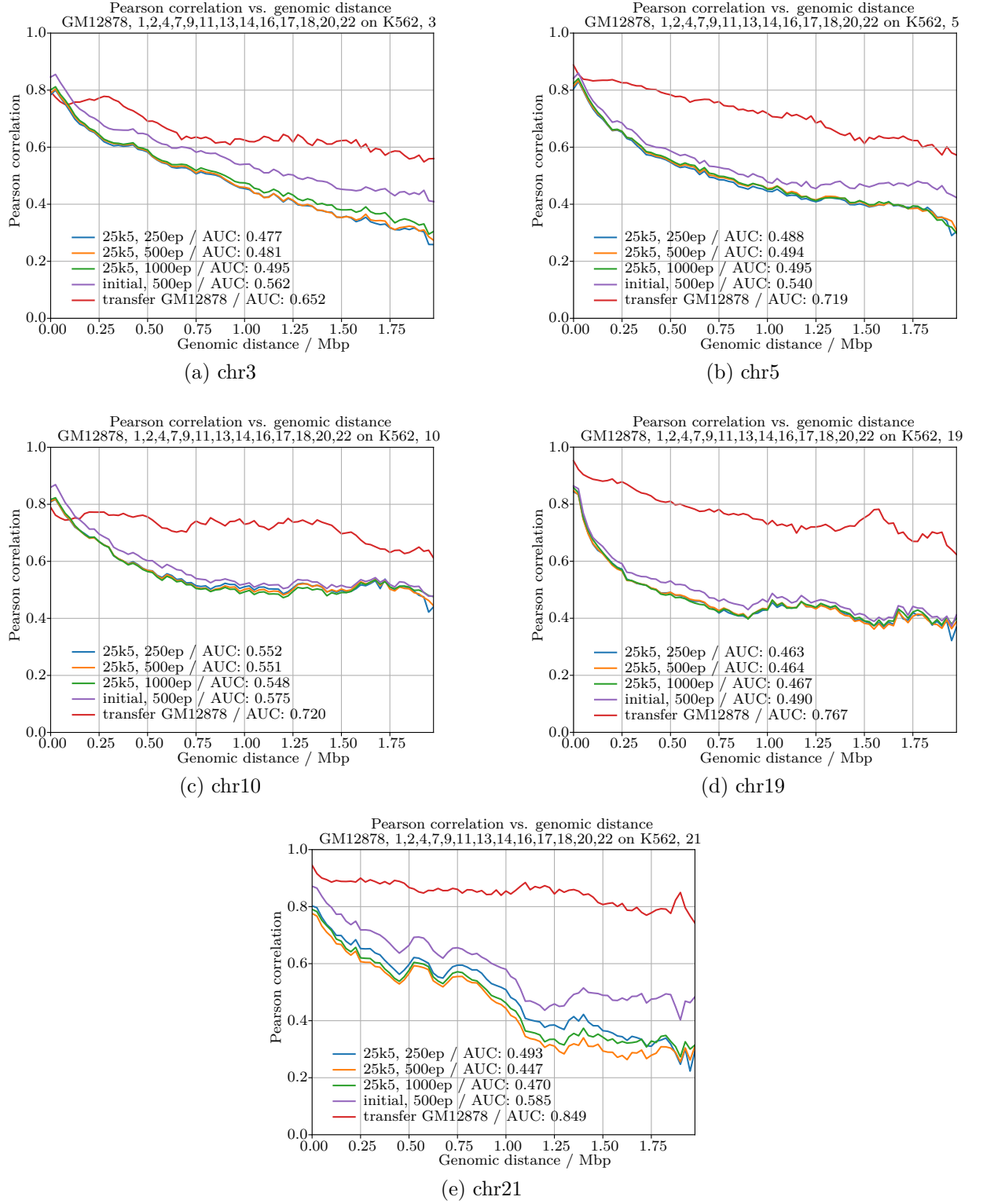


Figure 21: Pearson correlations, “5k – 25k” variant of DNN, test chromosomes

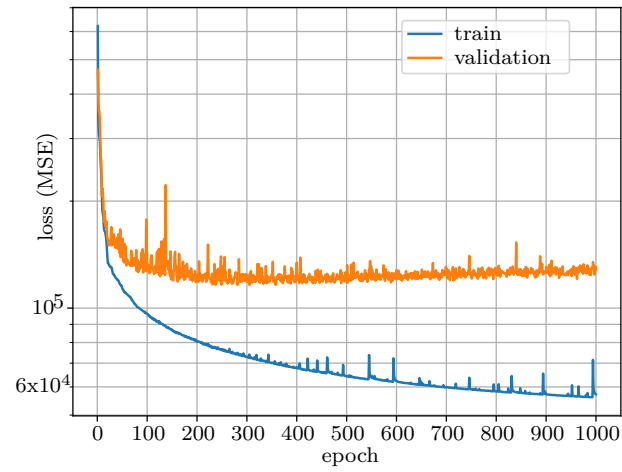


Figure 22: learning progress for DNN with $b_{feat} = 5$ kbp and $b_{mat} = 25$ kbp

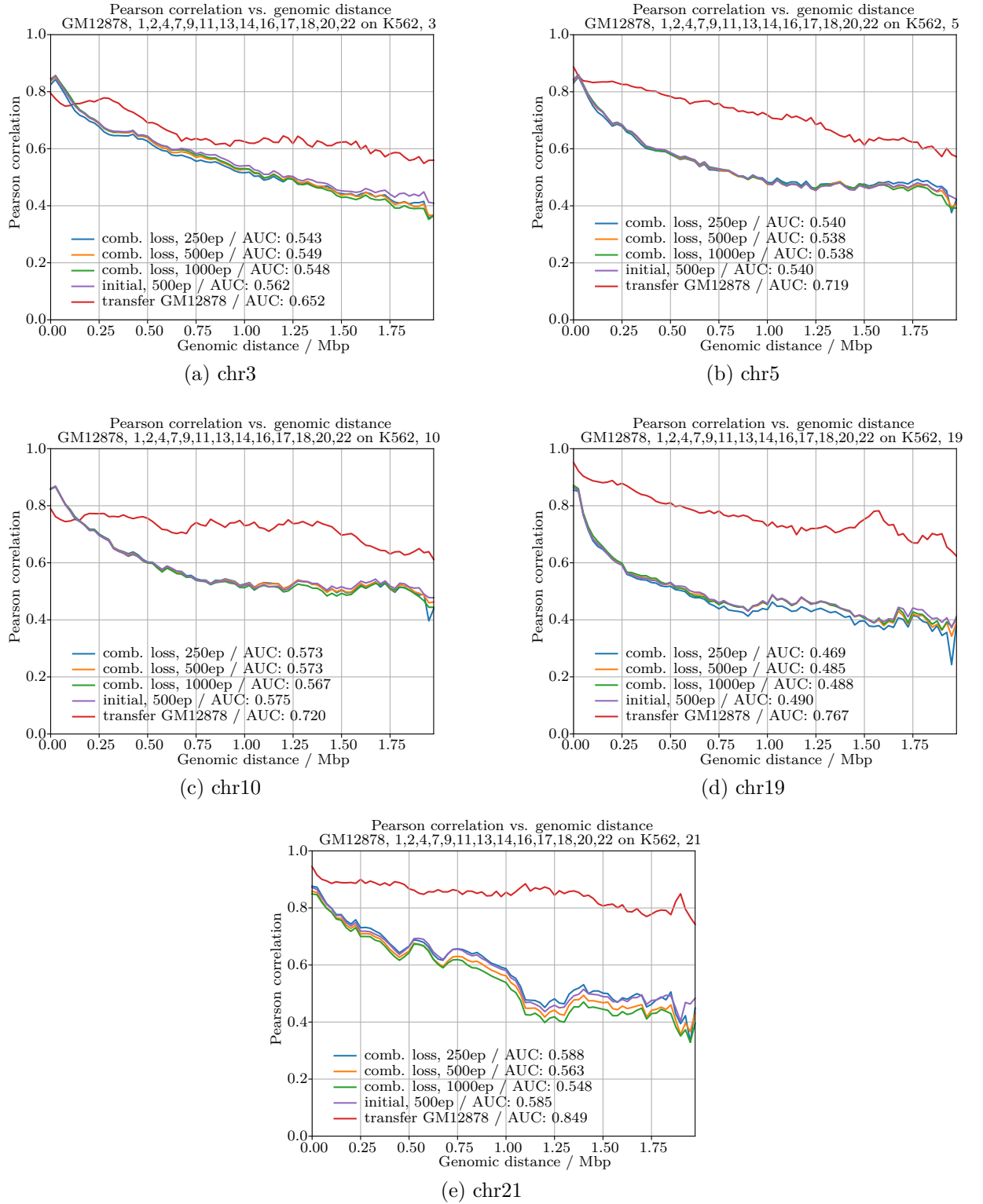


Figure 23: Pearson correlations, combined loss function (MSE, TV, VGG-16), test chromosomes

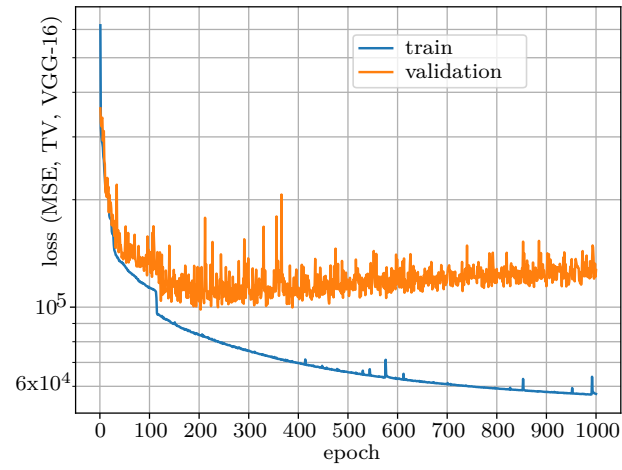


Figure 24: learning progress for DNN with combined loss function (MSE, TV, VGG-16)

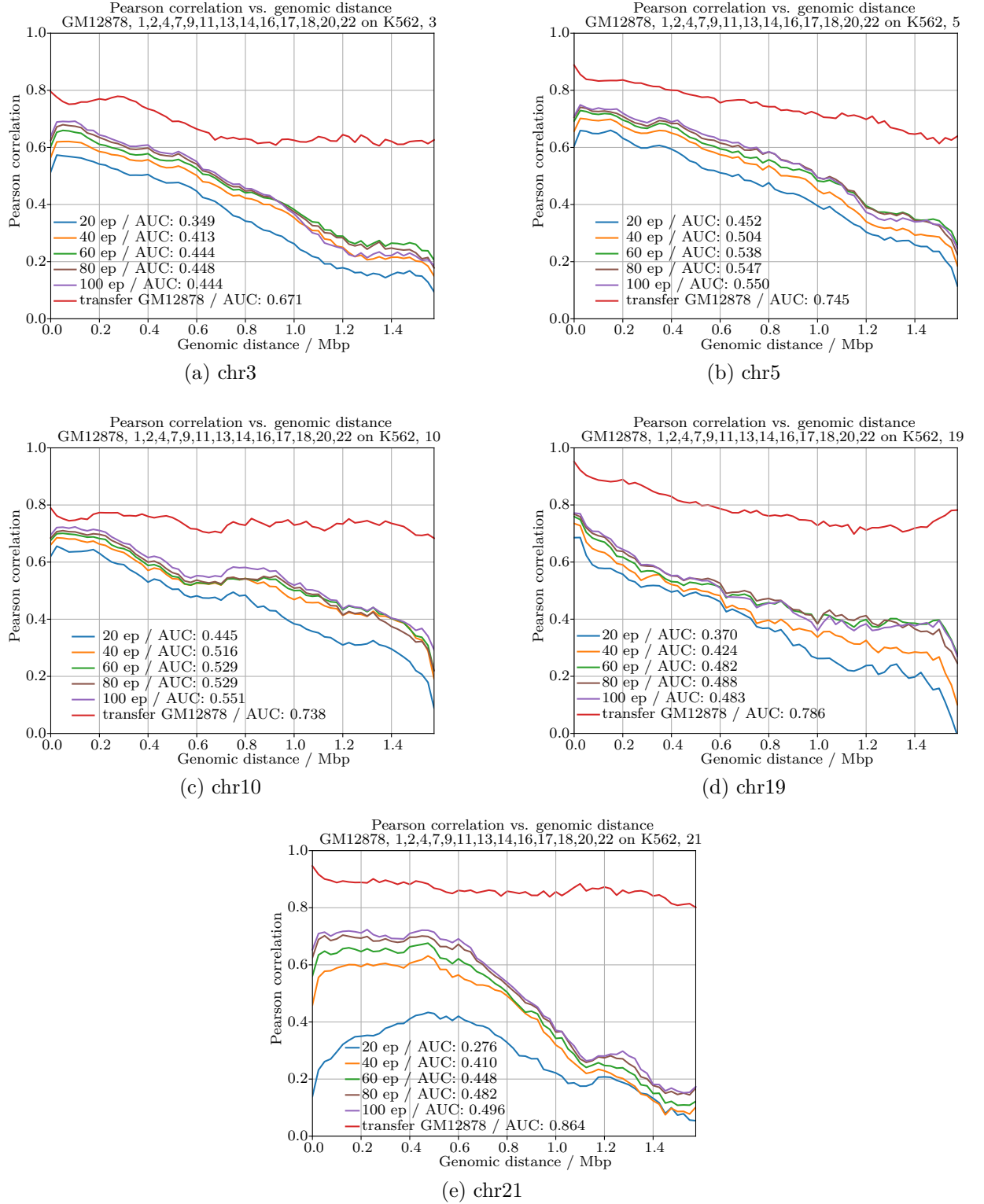


Figure 25: Pearson correlations, GAN, windowsize 64, test chromosomes

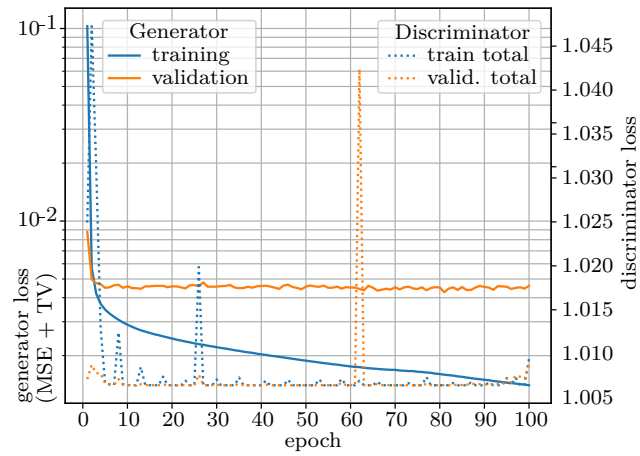


Figure 26: learning progress cGAN, window size 64

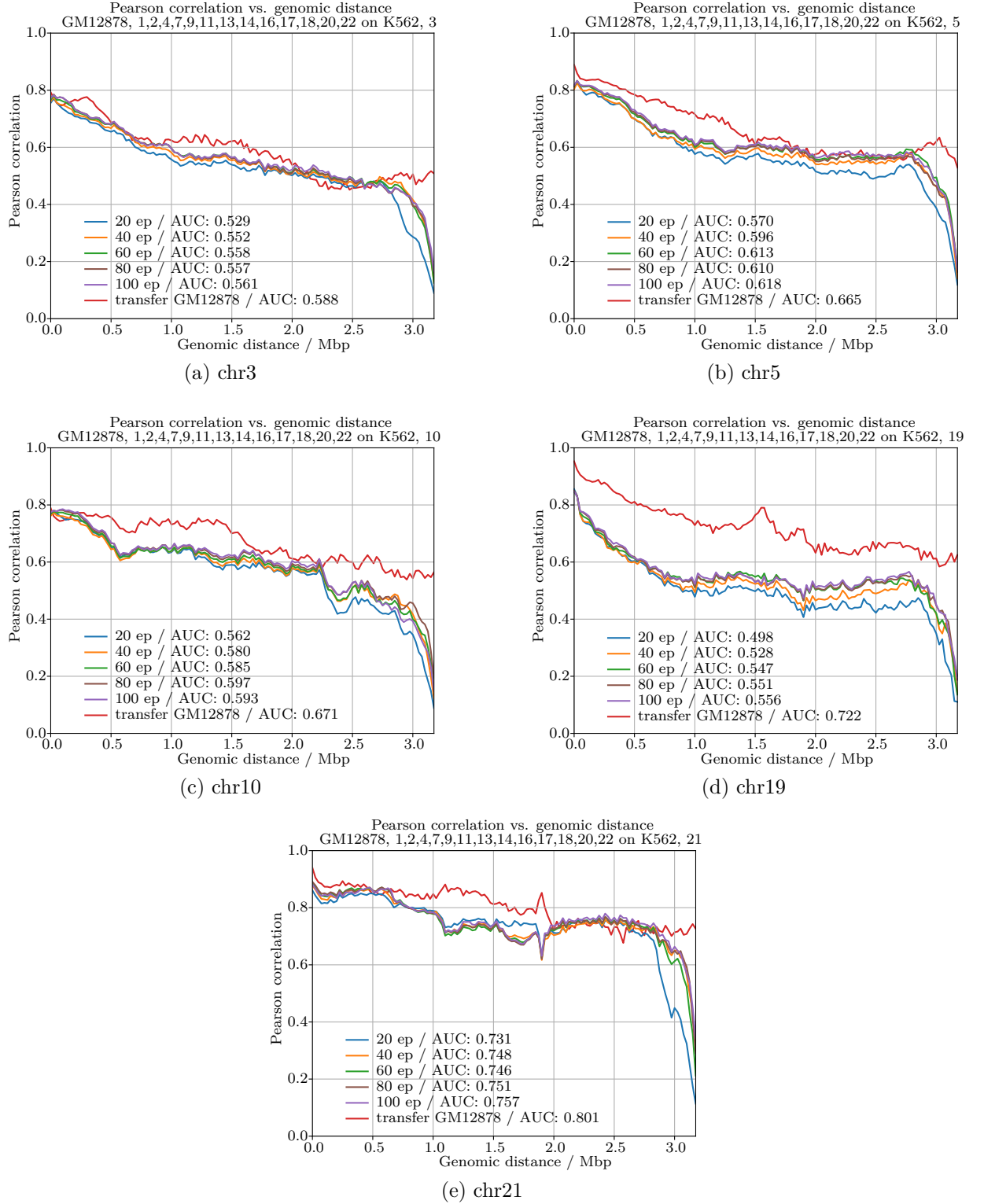


Figure 27: Pearson correlations, GAN, window size 128, test chromosomes

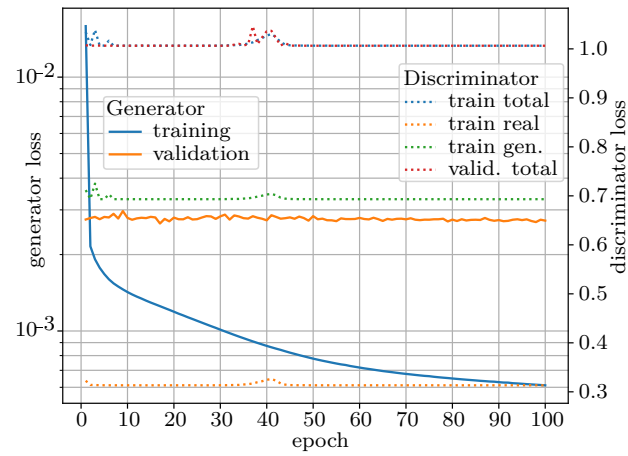


Figure 28: learning progress cGAN, window size 128

6 Discussion and Outlook

Approach by Farre et al also working for human data. Changes might bring benefits in certain situations, but are not a big improvement.

cGAN seems promising, but the discriminator remains to improve. Improvements still not good enough to beat transfer GM12878-K562 in many cases.

If possible, compare to other approaches. Maybe we can use the one by Zhang et al. for which data from study project exists, but it might take too long.

References

- [1] J. D. Watson and F. H. C. Crick. “Molecular Structure of Nucleic Acids: A Structure for Deoxyribose Nucleic Acid”. In: *Nature* 171.4356 (Apr. 1953), pp. 737–738. DOI: 10.1038/171737a0.
- [2] J.Y. Lee and T.L. Orr-Weaver. “Chromatin”. In: *Encyclopedia of Genetics*. Elsevier, 2001, pp. 340–343. DOI: 10.1006/rwgn.2001.0199.
- [3] Andrea Smallwood and Bing Ren. “Genome organization and long-range regulation of gene expression by enhancers”. In: *Current Opinion in Cell Biology* 25.3 (June 2013), pp. 387–394. DOI: 10.1016/j.ceb.2013.02.005.
- [4] David U. Gorkin, Danny Leung and Bing Ren. “The 3D Genome in Transcriptional Regulation and Pluripotency”. In: *Cell Stem Cell* 14.6 (June 2014), pp. 762–775. DOI: 10.1016/j.stem.2014.05.017.
- [5] Keerthi T. Chathoth and Nicolae Radu Zabet. “Chromatin architecture reorganization during neuronal cell differentiation in Drosophila genome”. In: *Genome Research* 29.4 (Feb. 2019), pp. 613–625. DOI: 10.1101/gr.246710.118.
- [6] Haoyue Zhang et al. “Chromatin structure dynamics during the mitosis-to-G1 phase transition”. In: *Nature* 576.7785 (Nov. 2019), pp. 158–162. DOI: 10.1038/s41586-019-1778-y.
- [7] Jennifer E. Phillips and Victor G. Corces. “CTCF: Master Weaver of the Genome”. In: *Cell* 137.7 (June 2009), pp. 1194–1211. DOI: 10.1016/j.cell.2009.06.001.
- [8] Jennifer E. Phillips-Cremins et al. “Architectural Protein Subclasses Shape 3D Organization of Genomes during Lineage Commitment”. In: *Cell* 153.6 (June 2013), pp. 1281–1295. DOI: 10.1016/j.cell.2013.04.053.
- [9] Jesse R. Dixon et al. “Chromatin architecture reorganization during stem cell differentiation”. In: *Nature* 518.7539 (Feb. 2015), pp. 331–336. DOI: 10.1038/nature14222.
- [10] E. Lieberman-Aiden et al. “Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome”. In: *Science* 326.5950 (Oct. 2009), pp. 289–293. DOI: 10.1126/science.1181369.
- [11] Suhas S.P. Rao et al. “A 3D Map of the Human Genome at Kilobase Resolution Reveals Principles of Chromatin Looping”. In: *Cell* 159.7 (Dec. 2014), pp. 1665–1680. DOI: 10.1016/j.cell.2014.11.021.
- [12] Houda Belaghzal, Job Dekker and Johan H. Gibcus. “Hi-C 2.0: An optimized Hi-C procedure for high-resolution genome-wide mapping of chromosome conformation”. In: *Methods* 123 (July 2017), pp. 56–65. DOI: 10.1016/j.ymeth.2017.04.004.
- [13] Ralf Krauth. *Improving predictions of Hi-C matrices from ChIP-seq data*. Tech. rep. Albert-Ludwigs Universität Freiburg, 2020. URL: <https://github.com/MasterprojectRK/HiCPrediction>.
- [14] D. S. Johnson et al. “Genome-Wide Mapping of in Vivo Protein-DNA Interactions”. In: *Science* 316.5830 (June 2007), pp. 1497–1502. DOI: 10.1126/science.1141319.
- [15] Gordon Robertson et al. “Genome-wide profiles of STAT1 DNA association using chromatin immunoprecipitation and massively parallel sequencing”. In: *Nature Methods* 4.8 (June 2007), pp. 651–657. DOI: 10.1038/nmeth1068.

- [16] ENCODE Project Consortium. “An integrated encyclopedia of DNA elements in the human genome”. In: *Nature* 489.7414 (Sept. 2012), pp. 57–74. DOI: 10.1038/nature11247. URL: <https://www.encodeproject.org>.
- [17] Carrie A. Davis et al. “The Encyclopedia of DNA elements (ENCODE): data portal update”. In: *Nucleic Acids Research* 46.D1 (Nov. 2017), pp. D794–D801. DOI: 10.1093/nar/gkx1081.
- [18] Boyan Bonev and Giacomo Cavalli. “Organization and function of the 3D genome”. In: *Nature Reviews Genetics* 17.11 (Oct. 2016), pp. 661–678. DOI: 10.1038/nrg.2016.112.
- [19] Chris A. Brackley et al. “Predicting the three-dimensional folding of cis-regulatory regions in mammalian genomes using bioinformatic data and polymer models”. In: *Genome Biology* 17.1 (Mar. 2016). DOI: 10.1186/s13059-016-0909-0.
- [20] Quinn MacPherson, Bruno Beltran and Andrew J. Spakowitz. “Bottom-up modeling of chromatin segregation due to epigenetic modifications”. In: *Proceedings of the National Academy of Sciences* 115.50 (Nov. 2018), pp. 12739–12744. DOI: 10.1073/pnas.1812268115.
- [21] Michele Di Pierro et al. “De novo prediction of human chromosome structures: Epigenetic marking patterns encode genome architecture”. In: *Proceedings of the National Academy of Sciences* 114.46 (Oct. 2017), pp. 12126–12131. DOI: 10.1073/pnas.1714980114.
- [22] Yifeng Qi and Bin Zhang. “Predicting three-dimensional genome organization with chromatin states”. In: *PLOS Computational Biology* 15.6 (June 2019). Ed. by Jian Ma, e1007024. DOI: 10.1371/journal.pcbi.1007024.
- [23] Pau Farré and Eldon Emberly. “A maximum-entropy model for predicting chromatin contacts”. In: *PLOS Computational Biology* 14.2 (Feb. 2018). Ed. by Alexandre V. Morozov, e1005956. DOI: 10.1371/journal.pcbi.1005956.
- [24] Jian Zhou and Olga G. Troyanskaya. “Probabilistic modelling of chromatin code landscape reveals functional diversity of enhancer-like chromatin states”. In: *Nature Communications* 7.1 (Feb. 2016). DOI: 10.1038/ncomms10528.
- [25] Ziad Al Bkhetan and Dariusz Plewczynski. “Three-dimensional Epigenome Statistical Model: Genome-wide Chromatin Looping Prediction”. In: *Scientific Reports* 8.1 (Mar. 2018). DOI: 10.1038/s41598-018-23276-8.
- [26] Yan Kai et al. “Predicting CTCF-mediated chromatin interactions by integrating genomic and epigenomic features”. In: *Nature Communications* 9.1 (Oct. 2018). DOI: 10.1038/s41467-018-06664-6.
- [27] Shilu Zhang et al. “In silico prediction of high-resolution Hi-C interaction matrices”. In: *Nature Communications* 10.1 (Dec. 2019). DOI: 10.1038/s41467-019-13423-8.
- [28] Laura D. Martens et al. “Identifying regulatory and spatial genomic architectural elements using cell type independent machine and deep learning models”. In: (Apr. 2020). DOI: 10.1101/2020.04.19.049585.
- [29] Pau Farré et al. “Dense neural networks for predicting chromatin conformation”. In: *BMC Bioinformatics* 19.1 (Oct. 2018). DOI: 10.1186/s12859-018-2286-z.
- [30] Shashank Singh et al. “Predicting enhancer-promoter interaction from genomic sequence with deep neural networks”. In: *Quantitative Biology* 7.2 (June 2019), pp. 122–137. DOI: 10.1007/s40484-019-0154-0.

-
- [31] Rui Peng. *Predicting High-order Chromatin Interactions from Human Genomic Sequence using Deep Neural Networks*. Tech. rep. Carnegie Mellon University, 2017. URL: <https://www.ml.cmu.edu/research/dap-papers/F17/dap-peng-rui.pdf>.
 - [32] Shashank Singh et al. “Predicting Enhancer-Promoter Interaction from Genomic Sequence with Deep Neural Networks”. In: *bioRxiv* (Nov. 2016). DOI: 10.1101/085241.
 - [33] Jacob Schreiber et al. “Nucleotide sequence and DNaseI sensitivity are predictive of 3D chromatin architecture”. In: *bioRxiv* (Jan. 2017). DOI: 10.1101/103614.
 - [34] Geoff Fudenberg, David R. Kelley and Katherine S. Pollard. “Predicting 3D genome folding from DNA sequence with Akita”. In: *Nature Methods* 17.11 (Oct. 2020), pp. 1111–1117. DOI: 10.1038/s41592-020-0958-x.
 - [35] Geoff Fudenberg, David R. Kelley and Katherine S. Pollard. “Predicting 3D genome folding from DNA sequence”. In: *bioRxiv* (Oct. 2019). DOI: 10.1101/800060.
 - [36] Ron Schwessinger et al. “DeepC: Predicting chromatin interactions using megabase scaled deep neural networks and transfer learning”. In: *bioRxiv* (Aug. 2019). DOI: 10.1101/724005.
 - [37] Sarvesh Nikumbh and Nico Pfeifer. “Genetic sequence-based prediction of long-range chromatin interactions suggests a potential role of short tandem repeat sequences in genome organization”. In: *BMC Bioinformatics* 18.1 (Apr. 2017). DOI: 10.1186/s12859-017-1624-x.
 - [38] A. Tsirikoglou, G. Eilertsen and J. Unger. “A Survey of Image Synthesis Methods for Visual Machine Learning”. In: *Computer Graphics Forum* 39.6 (Sept. 2020), pp. 426–451. DOI: 10.1111/cgf.14047.
 - [39] Tong Liu and Zheng Wang. “HiCNN2: Enhancing the Resolution of Hi-C Data Using an Ensemble of Convolutional Neural Networks”. In: *Genes* 10.11 (Oct. 2019), p. 862. DOI: 10.3390/genes10110862.
 - [40] Ian Goodfellow et al. “Generative Adversarial Nets”. In: *Advances in Neural Information Processing Systems*. Ed. by Z. Ghahramani et al. Vol. 27. Curran Associates, Inc., 2014, pp. 2672–2680. URL: <https://proceedings.neurips.cc/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf>.
 - [41] Mehdi Mirza and Simon Osindero. *Conditional Generative Adversarial Nets*. 2014. arXiv: 1411.1784 [cs.LG].
 - [42] Qiao Liu, Hairong Lv and Rui Jiang. “hicGAN infers super resolution Hi-C data with generative adversarial networks”. In: *Bioinformatics* 35.14 (July 2019), pp. i99–i107. DOI: 10.1093/bioinformatics/btz317.
 - [43] Hao Hong et al. “DeepHiC: A generative adversarial network for enhancing Hi-C data resolution”. In: *PLOS Computational Biology* 16.2 (Feb. 2020). Ed. by Ferhat Ay, e1007287. DOI: 10.1371/journal.pcbi.1007287.
 - [44] Michael C. Dimmick, Leo J. Lee and Brendan J. Frey. “HiCSR: a Hi-C super-resolution framework for producing highly realistic contact maps”. In: *bioRxiv* (Feb. 2020). DOI: 10.1101/2020.02.24.961714.
 - [45] Kai Huang, Vadim Backman and Igal Szleifer. “Interphase chromatin as a self-returning random walk: Can DNA fold into liquid trees?” In: (Sept. 2018). DOI: 10.1101/413872.

- [46] Artemi Bendandi et al. “Chromatin Compaction Multiscale Modeling: A Complex Synergy Between Theory, Simulation, and Experiment”. In: *Frontiers in Molecular Biosciences* 7 (Feb. 2020). DOI: 10.3389/fmolb.2020.00015.
- [47] Andre Bajorat. *Hi-C Predictions based on protein levels*. Tech. rep. Albert-Ludwigs Universität Freiburg, 2019. URL: https://www.bioinf.uni-freiburg.de/Lehre/Theses/TP_Andre_Bajorat.pdf.
- [48] Phillip Isola et al. “Image-to-Image Translation with Conditional Adversarial Networks”. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, July 2017, pp. 5967–5976. DOI: 10.1109/CVPR.2017.632.
- [49] Yingjing Lu. “The Level Weighted Structural Similarity Loss: A Step Away from the MSE”. In: (Apr. 30, 2019). arXiv: 1904.13362 [cs.CV].
- [50] Z. Wang, E. P. Simoncelli and A. C. Bovik. “Multiscale structural similarity for image quality assessment”. In: *The Thrity-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*. IEEE. DOI: 10.1109/acssc.2003.1292216.
- [51] Hang Zhao et al. “Loss Functions for Image Restoration With Neural Networks”. In: *IEEE Transactions on Computational Imaging* 3.1 (Mar. 2017), pp. 47–57. DOI: 10.1109/tci.2016.2644865.
- [52] Karen Simonyan and Andrew Zisserman. “Very Deep Convolutional Networks for Large-Scale Image Recognition”. In: *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*. Ed. by Yoshua Bengio and Yann LeCun. 2015. URL: <http://arxiv.org/abs/1409.1556>.
- [53] Justin Johnson, Alexandre Alahi and Li Fei-Fei. “Perceptual Losses for Real-Time Style Transfer and Super-Resolution”. In: (Mar. 27, 2016). arXiv: 1603.08155 [cs.CV].
- [54] Leonid I. Rudin, Stanley Osher and Emad Fatemi. “Nonlinear total variation based noise removal algorithms”. In: *Physica D: Nonlinear Phenomena* 60.1-4 (Nov. 1992), pp. 259–268. DOI: 10.1016/0167-2789(92)90242-f.
- [55] Rola Dali and Mathieu Blanchette. “A critical assessment of topologically associating domain prediction tools”. In: *Nucleic Acids Research* 45.6 (Mar. 2017), pp. 2994–3005. DOI: 10.1093/nar/gkx145.
- [56] Marie Zufferey et al. “Comparison of computational methods for the identification of topologically associating domains”. In: *Genome Biology* 19.1 (Dec. 2018). DOI: 10.1186/s13059-018-1596-9.
- [57] Emily Crane et al. “Condensin-driven remodelling of X chromosome topology during dosage compensation”. In: *Nature* 523.7559 (June 2015), pp. 240–244. DOI: 10.1038/nature14450.
- [58] Hanjun Shin et al. “TopDom: an efficient and deterministic method for identifying topological domains in genomes”. In: *Nucleic Acids Research* 44.7 (Dec. 2015), e70–e70. DOI: 10.1093/nar/gkv1505.
- [59] Joachim Wolff et al. “Galaxy HiCExplorer: a web server for reproducible Hi-C data analysis, quality control and visualization”. In: *Nucleic Acids Research* 46.W1 (June 2018), W11–W16. DOI: 10.1093/nar/gky504.

Acronyms

ChIA-PET chromatin interaction analysis by paired-end tag sequencing.

ChIP-seq chromatin immunoprecipitation followed by sequencing.

CNN convolutional neural network.

DamID DNA adenine methyltransferase identification.

DNN dense neural network.

GAN generative adversarial network.

LSTM long-short-term memory.

TAD topologically associating domain.