OXFORD

## Systems biology

# DeepAntigen: a novel method for neoantigen prioritization via 3D genome and deep sparse learning

Yi Shi ᴵᴰ [1,2,3,*,†] Zehua Guo ᴵᴰ [2,4,†], Xianbin Su[1,†], Luming Meng ᴵᴰ [5,*], Mingxuan Zhang[6], Jing Sun[7], Chao Wu[7], Minhua Zheng[7], Xueyin Shang[1], Xin Zou ᴵᴰ [1], Wangqiu Cheng[2,3], Yaoliang Yu[8], Yujia Cai[1], Chaoyi Zhang[9], Weidong Cai[9], Lin-Tai Da[1,*], Guang He[2,3,*] and Ze-Guang Han[1,*]

[1]Key Laboratory of Systems Biomedicine (Ministry of Education), Shanghai Centre for Systems Biomedicine, Shanghai Jiao Tong University, Shanghai 200240, China, [2]Shanghai Jiao Tong University, Shanghai 200030, China, [3]Shanghai Key Laboratory of Psychotic Disorders, and Brain Science and Technology Research Center, Shanghai Jiao Tong University, Shanghai 200030, China, [4]Department of Instrument Science and Engineering, School of Electronic Information and Electrical Engineering, Shanghai Jiao Tong University, Shanghai 200240, China, [5]College of Biophotonics, South China Normal University, Guangzhou 510631, China, [6]Department of Mathematics, University of California San Diego, La Jolla, CA 92093-0112, USA, [7]Department of General Surgery & Shanghai Minimally Invasive Surgery Center, Ruijin Hospital, Shanghai Jiao Tong University, Shanghai 200025, China, [8]David R. Cheriton School of Computer Science, University of Waterloo, Waterloo, ON N2L3G1, Canada and [9]School of Computer Science, The University of Sydney, Darlington, NSW, 2008, Australia

*To whom correspondence should be addressed.

Associate Editor: Jonathan Wren

†The authors wish it to be known that, in their opinion, the first three authors should be regarded as Joint First Authors.

## Abstract

**Motivation:** The mutations of cancers can encode the seeds of their own destruction, in the form of T-cell recognizable immunogenic peptides, also known as neoantigens. It is computationally challenging, however, to accurately prioritize the potential neoantigen candidates according to their ability of activating the T-cell immunoresponse, especially when the somatic mutations are abundant. Although a few neoantigen prioritization methods have been proposed to address this issue, advanced machine learning model that is specifically designed to tackle this problem is still lacking. Moreover, none of the existing methods considers the original DNA loci of the neoantigens in the perspective of 3D genome which may provide key information for inferring neoantigens' immunogenicity.

**Results:** In this study, we discovered that DNA loci of the immunopositive and immunonegative MHC-I neoantigens have distinct spatial distribution patterns across the genome. We therefore used the 3D genome information along with an ensemble pMHC-I coding strategy, and developed a group feature selection-based deep sparse neural network model (DNN-GFS) that is optimized for neoantigen prioritization. DNN-GFS demonstrated increased neoantigen prioritization power comparing to existing sequence-based approaches. We also developed a webserver named deepAntigen (http://yishi.sjtu.edu.cn/deepAntigen) that implements the DNN-GFS as well as other machine learning methods. We believe that this work provides a new perspective toward more accurate neoantigen prediction which eventually contribute to personalized cancer immunotherapy.

**Availability and implementation:** Data and implementation are available on webserver: http://yishi.sjtu.edu.cn/deepAntigen.

**Contact:** yishi@sjtu.edu.cn or menglum@scnu.edu.cn or darlt@sjtu.edu.cn or heguang@sjtu.edu.cn or hanzg@sjtu.edu.cn

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

# 1. Introduction

The approval of several immunotherapies has led to dramatic changes in cancer therapy. In variety human malignancies, therapeutic efficacy was enhanced by immunotherapies via boosting the endogenous T cell's ability to destroy cancer cells (Schumacher and Schreiber, 2015). The 'checkpoint inhibitors' therapies work by blocking proteins that act as molecular breaks for T cells. With the breaks removed, T cells can better undertake their job to kill cancer cells. Despite the great success of checkpoint inhibitors, still many patients do not respond to the agents, and many that do temporarily respond, eventually relapse. Moreover, checkpoint inhibitors do not fully take advantage of the T cell's exquisite specificity, one of its most important characteristics (Sompayrac, 2019). This led many researchers pay more attention to the new immunotherapy strategies against tumor known as neoantigen therapies. T cells are potent at killing when they recognize 'foreign' antigens which could be some protein fragments from an invading virus or bacteria. The key ability of T cells in distinguishing foreign antigens from self and prevents autoimmunity, which on the contrast makes them less potent in recognizing tumor cells because they are our own but abnormal cells. The T cells overcome this dilemma in two ways. First, they tend to respond to tissue-specific antigens which are specific amino-acid fragments produced by cells of certain types. Second, T cells respond to neoantigens which are small peptides generated in tumor cells containing high level of DNA mutations. The non-synonymous mutations can be entirely absent from the human genome, leading the cancer cells vulnerable to T cells as they look 'foreign' (Sompayrac, 2019).

In several clinical practices, it has been demonstrated that endogenous T cells with mounted cancer-killing T-cell receptor (TCR) are able to recognize epitopes which are composed of the peptides displayed on major histocompatibility complexes (MHCs) on the surface of the cancer cells (Ott *et al.*, 2017; Schumacher and Schreiber, 2015). With the help of DNA- and RNA-sequencing technology, it has been revealed that tens to thousands of different somatic mutations can be generated during cancer initiation and progression, depending on different cancer types (Castro *et al.*, 2019; Prior *et al.*, 2019; Volkov *et al.*, 2020). Most of these mutations are often caused by genomic instability within the tumor cells and lead to no obvious cell growth advantage; they are also known as passenger mutations. On the contrast, a small percent of these mutations are known as driver mutations which interfere with normal cell regulation and help to drive cancer growth and resistance to targeted therapies (Yarchoan *et al.*, 2017). Both passenger and driver mutations can cause tumor to express abnormal proteins or polypeptides that cannot be found in normal cells as they can be non-synonymous mutations that alter protein-coding sequences. When cell metabolize, the proteins possessing abnormal sequences are cut into short peptides and are presented as epitopes on the cell surface by the MHC (also known as human leukocyte antigen, HLA, in human case) molecules, which have a chance to be recognized by T cells as foreign antigens (Yarchoan *et al.*, 2017). An effective neoantigen, which leads to the final immunological response, is determined by many factors. For instance, Dintzis *et al.* (1976) found that size-fractionated linear polymers of acrylamide substituted with hapten can affect the immunogenicity triggering. Other factors such as peptide degradation and transportation, peptide–MHC binding affinity and stability and pMHC–TCR interaction should also be considered (Blaha *et al.*, 2019).

Based on the earlier knowledge, in ideal situation, after the DNA-sequencing procedure, potential neoantigens can be synthesized *in vitro* and their efficacy can be validated *in vivo* via either cancer cell-line or animal model, before conducting in clinical practice (Schumacher and Schreiber, 2015; Yarchoan *et al.*, 2017). Indeed, the cancers with a single dominant mutation can often be effectively treated by focusing on the driver mutation (O'Brien *et al.*, 2003; Yarchoan *et al.*, 2017). Nevertheless, in many other cancer situations, the somatic mutations are usually abundant, which lead to a computationally challenging task to efficiently prioritize the potential neoantigen candidates according to their ability to activate the T cell's immunoresponse (Hackl *et al.*, 2016). In the past decade, many prediction methods have been proposed to address the neoantigen prioritization problem (Jurtz *et al.*, 2017; Lundegaard *et al.*, 2008; Nielsen and Andreatta, 2016). These methods can be categorized into two major classes: the protein spatial conformation-based approaches which consider the pMHC and TCR 3D structures, and the protein sequence-based approaches which consider the amino-acid combinatorial characters. For the protein spatial conformation-based approaches, when high-quality pMHC 3D structures are available, methods such as molecular dynamic (MD) can be adopted to explore the complex interaction between TCR and pMHC (Blevins *et al.*, 2016; Riley *et al.*, 2018; Wang *et al.*, 2017). If high-quality pMHC spatial information is lacking, by sacrificing computational complexity and spatial model accuracy, computational pMHC modeling can be adopted, followed by 3D to 1D feature transformation and machine learning approaches (Riley *et al.*, 2019). Most neoantigen prediction methods belong to the sequence-based class because they can usually be set up efficiently (Gupta *et al.*, 2016; Hackl, et al., 2016), and there are much larger datasets available for training and validation (Vita *et al.*, 2019; Zhang *et al.*, 2011).

Early sequence-based methods such as BIMAS (Parker *et al.*, 1994) and SYFPEITHI(Schuler *et al.*, 2007) utilized the position-specific scoring matrices (PSSMs), which are defined from experimentally confirmed peptide binders of a particular MHC allele (Hackl *et al.*, 2016). More sophisticated approaches based on machine learning techniques were later developed which were demonstrated to perform better than the PSSM-based methods; these approaches capture and utilize the non-linear nature of the pMHC–TCR interaction. In recent years, consensus approaches such as CONSENSUS (Moutaftsi *et al.*, 2006) and NetMHCcons (Karosiene *et al.*, 2012) were exploited which combine the results of multiple neoantigen prediction tools, aiming to obtain more robust and accurate outcomes, and their efficacies were supported by experimental results. Nonetheless, the performance gain of these methods is determined by the weighting scheme among different prediction components, which lead to increased computational complexity (hyperparameter tuning). Because the peptide MHC binding can be affected by HLA allele variety, most recently, the pan-specific methods, such as NetMHCpan (Jurtz *et al.*, 2017; Nielsen and Andreatta, 2016), were developed which allow the HLA-type independent prioritization. In NetMHCpan, a neural network is first trained based on multiple public datasets, then the binding affinity for a given peptide–MHC complex is predicted according to the trained neural network, with the polymorphic HLA types, e.g. HLA-A, HLA-B or HLA-C being considered. Even compared to HLA allele-specific approaches (Hackl *et al.*, 2016; Trolle *et al.*, 2015), both NetMHCpan (Jurtz *et al.*, 2017) and NetMHCIIpan (Karosiene *et al.*, 2013) could perform remarkably better. Although methods such as NetMHC or NetMHCpan were designed to predict peptide–MHC binding affinity, they were either considered as strong indicators for neoantigens' effectiveness (Harndahl *et al.*, 2012; Lundegaard *et al.*, 2011; Rasmussen *et al.*, 2016), or were adopted as important features in the state-of-the-art neoantigen-predicting methods such as Neopepsee and pTuneos (Kim *et al.*, 2018; Zhou *et al.*, 2019). More recently, Wu *et al.* (2019) proposed a recurrent-neural-network-based approach DeepHLApan which considered both pMHC binding and potential immunogenicity, yet sequence information of both peptide and HLA were still adopted as training features.

For all the existing neoantigen prediction methods, although several evaluation criteria were proposed for a more fair and robust comparison (Peters *et al.*, 2006; Trolle *et al.*, 2015; Wang *et al.*, 2008), independent benchmark studies that can be used to recommend specific tools are still lacking. More importantly, although there are abundant previous researches indicating that somatic mutations, including point mutations, gene fusions and copy number abnormalities do not occur at random in the perspective of genome 3D conformation (Berger *et al.*, 2011; Branco and Pombo, 2006; Engreitz *et al.*, 2012; Mani *et al.*, 2009; Mathas *et al.*, 2009; Meaburn *et al.*, 2007; Nikiforova, 2000; Roix *et al.*, 2003; Wijchers and de Laat, 2011), for which we also did a thorough study and discovered the somatic comutation hotspot (SCH) in 3D genome (Shi *et al.*, 2016), none of the existing neoantigen prediction methods considers this spatial genomic information of somatic mutations, i.e.

the DNA loci of these mutations in the perspective of high-order genome 3D conformation. We believe that the 3D genome information could contain much richer information compared to the existing amino-acid sequence-based neoantigen prediction methods. Therefore, in this work, we retrospect the DNA origin of the neoantigens, both immunopositive and immunonegative, in the context of the genome 3D conformation, and demonstrate some discoveries that worth paying attention to. We adopted the 3D genome information into an ensemble peptide feature coding scheme, and developed a group feature selection-based deep sparse neural network (DNN-GFS) model that is customized and optimized for the neoantigen prediction task. We also developed an off-the-shelf webserver that implements the DNN-GFS method along with other machine learning methods; the webserver takes sequencing result (vcf file) and produces prioritized neoantigens as well as some useful intermediate functions such as vcf annotation and candidate neoantigen enumeration, etc. The whole workflow is illustrated in Figure 1, where the adoption of 3D genome information, ensemble feature coding and the DNN-GFS algorithm are keys for distinguishing our neoantigen prediction from all the existing methods.

# 2 Materials and methods

## 2.1 Immunogenicity data curation and reference genome mapping

The neoantigen peptide sequences and the immune responses were collected from the IEDB database under the T-Cell Assay category (Vita *et al.*, 2019). For the cross-validation experiments, we collected training data before 2018 in IEDB; After collecting 337 248 peptide records in the primary dataset, we performed filtering under *Homo sapiens* and MHC-I subtypes and restrained the peptide length nine, as well as merging identical records and mapping to human reference genome hg19. When mapping the peptides to the reference genome, we first applied the PANDAS library to create a data frame object for subsequent processing. Then we assigned the column name by importing a name dictionary and filtered the dataset so that the only entries left have *H.sapiens* as their hostname. The dataset was further cleaned up by applying two functions we developed, Letter_check and Drop_legal, which checks for amino-acid alphabet legitimacy. We developed a pipeline to query the BLAST (Boratyn *et al.*, 2013) web server and map the gene names to chromosomes and starting positions. The dataset was divided into 711 partitions where each partition contains 100 sequences. To set up BLAST queries, we restricted the search to *H.sapiens* using the entrez ID keywords and used the PAM30 matrix to find matches; the gap costs were adjusted to regulate gap penalty. We then queried BLAST iteratively. For each match, we adopted the accession and raw bit score for the first hit. After obtaining the accessions, we used
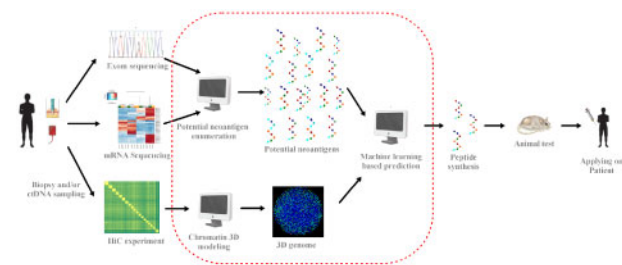


**Fig. 1.** Workflow of neoantigen therapy supported by 3D genome information. Left to right: tumor sample collection from patient; Whole-exome sequencing and mRNA sequencing for somatic mutations calling and gene expression estimation (whether the mutated DNA is expressed into mRNA and could potentially be translated into protein/peptide), respectively; Hi-C data curation to obtain 3D genome information; candidate peptides determined by NGS are generated and by combining 3D genome information immunopositive peptides are predicted machine learning methods; the top ranked peptides are screened by conducting animal experiments; the final peptide penal can be applied back to the target patient. This work aims to solve the tasks within the dashed red frame

the DAVID tool (Huang *et al.*, 2009) to obtain the gene names composed with gene symbols and the chromosome positions are also obtained. The final results contain a tuple of peptides, HLA subtype, chromosome number and chromosome position. For identical peptides with multiple immune experiments, we define peptides with positive rate >80% as immunopositive samples and with positive rate <20% as immunonegative peptides. Finally, we obtained 3909 peptides, with 809 immunopositive peptides and 3100 immunonegative peptides. We also collected a standalone validation dataset from IEDB dated after 2018 and performed the same operation mentioned herein. In the end, 430 validation peptides were obtained with 125 positive samples and 305 negative samples.

## 2.2 Hi-C data curation and A/B compartment determination

For the chromatin 3D conformation data, we used two well-known Hi-C data resources (Dixon *et al.*, 2012; Rao *et al.*, 2015), and obtained eight Hi-C datasets, i.e. hESC, IMR90, GM12878, HUVEC, IMR90-Rao, NHEK, K562 and KBM7. The Knight–Ruiz normalization (KR-norm) was applied on both intrachromosomal and the interchromosomal (genomewise) Hi-C contact maps. Bin sizes of 40, 100 and 500 kb were adopted for intrachromosomal contact frequency analyses, A/B compartment analyses and interchromosomal contact frequency analyses and chromatin 3D modeling. To determine the compartment activeness (compartment A: active, compartment B: inactive) of each chromosome bin, we used individual chromosome Hi-C contact maps. We first diagonal normalized each contact map by dividing the contact frequencies by their corresponding off-diagonal mean. Then, we computed the Pearson correlation coefficient (PCC) matrices for each chromosome, and the compartment type was jointly determined by the sign of the eigenvector corresponding to the first eigenvalue of the PCC matrices and the signal of the epigenetic marker H3k4me1.

## 2.3 Chromatin 3D modeling

We used MD and developed a human genome 3D conformation modeling approach with resolution 500 kb (bin size) for all eight Hi-C datasets. The bins were coarse-grained as beads and intact genome was represented by bead-on-the-string structures consisting of 23 polymer chains. The beads' spatial positioning is affected by both chromatin connectivity that constrains linearly neighboring beads in close 3D proximity and chromatin activity that ensures active regions tend to be located closer to the nucleus center. The chromatin activity was determined according to compartment degree that can be directly calculated from Hi-C matrix as described earlier and also in previous work (Xie *et al.*, 2017). Based on compartment degree index, beads were assigned distance values with respect to the nuclear center; the conformation of chromatin was then optimized from random structures with MDs approach by applying bias potential to satisfy these distance constraints. For each cell linage, 300 feasible conformation structures were optimized from random ones to reduce possible variation for further analysis.

## 2.4 Deep sparse neural network methods

The deep feedforward networks, also known as multilayer perceptrons (MLPs) were used in this work as the basic neural network architecture (Goodfellow *et al.*, 2016). For a single unit, its basic form is $y = f(x; \theta)$, where $x$ is the input, $y$ is the output and $\theta$ represents the parameters of the network that need to be optimized by adaptable methods. For a single middle-layer neural network, a generic form can be given as:

$$\mathbf{y}_k = g_k(\mathbf{W}_k\mathbf{x}_k + \mathbf{b}_k) \tag{1}$$

where $\{\mathbf{W}_k, \mathbf{b}_k\}$ are the optimized parameters of the layer, corresponding to $\theta$ in basic form, and $g_k(.)$ is the activation function of the layer for which we chose the widely adopted linear unit (ReLU) and the sigmoid unit in our model. Their function forms are $g(z) = \max\{0, z\}$ and $g(z) = \sigma(z)$; $\mathbf{x}_k$ is the input and $\mathbf{y}_k$ is the output. Note that, an important prerequisite in our model is $\mathbf{x}_{k+1} = \mathbf{y}_k$,

which makes all layers form the whole network, and specially, there is no input $\mathbf{x}_k$ for input layer. To obtain a set of adaptable $\{\mathbf{W}_k, \mathbf{b}_k\}$, the network should be trained multiple times by minimizing the regularized objective function $\tilde{J}$ (Goodfellow *et al.*, 2016):

$$\tilde{J}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = J(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \lambda R(\boldsymbol{\theta}) \tag{2}$$

In practice, only the weights ($\mathbf{W}$) of $\theta$ at each layer are penalized, and to simplify the equation, $\theta$ can be replaced by $\mathbf{w}$ (Goodfellow *et al.*, 2016):

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = J(\mathbf{w}; \mathbf{X}, \mathbf{y}) + \lambda R(\mathbf{w}) \tag{3}$$

where $J(\mathbf{w}; \mathbf{X}, \mathbf{y})$ is the standard objective function, $R(\mathbf{w})$ is the parameter norm penalty, and $\lambda \in [0, \infty]$ is a hyperparameter that weights the two terms. Larger values of $\lambda$ correspond to more regularization and setting $\lambda =$ results in no regularization. In this work, we set $J(.)$ as the cross-entropy loss. Many effective regularization strategies have been previously studied. The most common regularization strategy is the $L_2$ norm penalization, which is usually adopted to avoid overfitting. Its general form is:

$$R(\mathbf{w}) = \|\mathbf{w}\|_2^2 \tag{4}$$

also known as Tikhonov regularization or ridge regression. Another common practice is the $L_1$ regularization, which has a similar presentation:

$$R(\mathbf{w}) = \|\mathbf{w}\|_1 = \sum_i |\mathbf{w}_i| \tag{5}$$

that is the sum of absolute values of all weights. Particularly, the least absolute shrinkage and selection operator (LASSO) is a typical model that uses a $L_1$ penalization. The $L_1$ regularization can not only avoid overfitting, but also obtain a sparser solution than $L_2$, by making a subset of the weights to become zero (or very close to zero), suggesting that the corresponding features may safely be discarded. Due to this important property and the ability of preventing overfitting, $L_1$ regularization is used in feature selection scenario extensively (Goodfellow *et al.*, 2016). Note that, recent study has revealed that sparsity is the key to imitate human brain for the neural network (Dettmers and Zettlemoyer, 2019).

Apparently, the regularization can prevent overfitting, but its contribution is not limited to that. Scardapane *et al.* (2017) considered group-level sparsity, a weight grouping strategy was achieved by grouping all outgoing connections from a single neuron, which may induce the property of pruning the corresponding neuron from the network. As introduced in group LASSO (Simon and Tibshirani, 2012), group sparse regularization, e.g. $L_{2,1}$ norm, can be written as:

$$R_{\ell_{21}}(\mathbf{w}) \triangleq \sum_{\mathbf{g} \in G} \sqrt{|\mathbf{g}|} \|\mathbf{g}\|_2 \tag{6}$$

where $|\mathbf{g}|$ is the dimensionality of the vector $\mathbf{g}$, vector $\mathbf{g}$ corresponds to weight matrix $\mathbf{W}$, every $\mathbf{g}$ is one row of a matrix $\mathbf{W}$, denoting all outgoing connections from an input neuron. $G$ is the set of $\mathbf{g}$, $\mathbf{g} \in G$, which is the result of grouping $\mathbf{W}$ by row. Furthermore, sparse group LASSO penalization was proposed by combining LASSO and group LASSO (Scardapane *et al.*, 2017; Simon and Tibshirani, 2012; Simon *et al.*, 2013)

$$R_{\text{SGL}}(\mathbf{w}) \triangleq R_{\ell_{21}}(\mathbf{w}) + R_{\ell_1}(\mathbf{w}) \tag{7}$$

which can increase the sparsity above group sparse regularization. In addition, the hyperparameter can be used to weight the two terms, that is (Friedman *et al.*, 2010):

$$R_{\text{SGL}}(\mathbf{w}) \triangleq (1 - \alpha) R_{\ell_{21}}(\mathbf{w}) + \alpha R_{\ell_1}(\mathbf{w}) \tag{8}$$

where $\alpha = 1$ corresponds to the $L_1$ term and $\alpha = 0$ corresponds to the $L_{2,1}$ term. This form gives users more choice for their problem.

## 2.5 Group feature selection-based DNN (DNN-GFS)

Traditional DNN and some relevant sparse DNNs have a good performance but remain to be improved in many research fields (Goodfellow *et al.*, 2016). When real problems are handled by deep learning, there are usually some prior knowledge neglected, leading to an unideal performance. If we only consider the datasets, it is difficult to obtain the optimal model and the corresponding parameters we expect. Moreover, the situation will get worse with decreasing sample size, especially in biology problems with more features than samples. But when the prior information is imposed on models, the model will be closer to our expectation and generalization may be improved.

For our neoantigen prioritization problem, based on the existing sparse DNN models (Friedman *et al.*, 2010; Simon and Tibshirani, 2012; Simon *et al.*, 2013), we develop a new regularization strategy that aims to tackle both feature selection and the group sparse regularization challenge, which is an extension of the $L_2$ and $L_1$ penalization. Specifically, the feature grouping nature is considered in group sparse regularization, forming a new regularization strategy. We term it group feature selection (GFS) regularization. In the feature vector of our neoantigen prediction problem, some groups contain multiple features and some groups contain a single feature. In the former cases, features of the same group need to be either all selected or all rejected, simultaneously. This means that all outgoing connections from all neurons in one group should be either simultaneously all zeros, or all non-zeros (Scardapane *et al.*, 2017). The GFS regularization can be written as follows:

$$R_{\text{GFS}}(\mathbf{w}) \triangleq \sum_{\bar{\mathbf{g}} \in G_f} |\mathbf{F}_s| \sqrt{|\bar{\mathbf{g}}|} \|\bar{\mathbf{g}}\|_2 \tag{9}$$

where vector $\bar{\mathbf{g}}$ is the average of the squares of $\mathbf{g}$ vectors of a feature group, which can efficiently reduce computational complexity. $|\bar{\mathbf{g}}|$ is the dimensionality of the vector $\bar{\mathbf{g}}$, and $G_f$ is the result of grouping again by feature group information based on $G$ (groups of group LASSO). As Figure 2a illustrates, some features form new groups. $|\mathbf{F}_s|$ is the corresponding feature number matrix of $G_f$. Note that, when a group contains a single feature, the expression can be simplified as $L_{2,1}$. Moreover, for one-dimensional groups, it can also be reduced to the standard LASSO, while for all features in a new group, it is closer to $L_2$ regularization. These regularization terms other than $L_2$ are convex but non-smooth, since their gradient is not defined when $\|\bar{\mathbf{g}}\|_2 = 0$, which is illustrated in Figure 2c.

The GFS devised here is a flexible regularization strategy as $G_f$ can be customized according to different preferences to adapt various requirements. Furthermore, $|\mathbf{F}_s|$ is also chosen skillfully in this work, i.e. $|\mathbf{F}_s|$ can be replaced or rectified by other coefficients, which is able to enlarge or narrow the differences among groups. When imposing $R_{\text{GFS}}(\mathbf{W})$ on the $\mathbf{W}$, we achieve the feature selection effect, as illustrated in Figure 2b. Results that are based on other regularization strategies are shown in Supplementary Figures S5 and S6, and it is
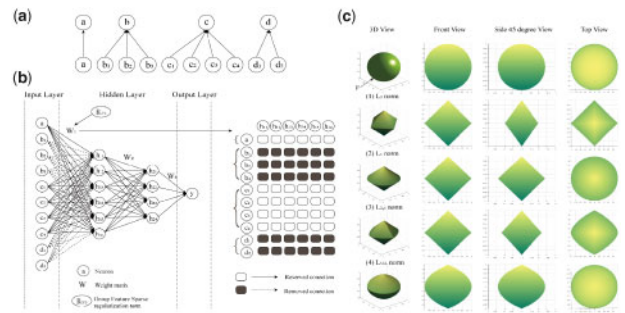


**Fig. 2.** The DNN-GFS method. (**a**) Illustration of features belonging to groups of different sizes. All features belong to at most one group. A group can contain a single feature or multiple features. (**b**) Illustration of the DNN-GFS architecture and the GFS effect. (**c**) Illustration of the geometric principles of different regularization terms applied on the weighted neural network wiring and 2D projection from three representative views. F denotes Front view in c (1)

demonstrated that only GFS can achieve the GFS effect. The detailed comparisons including network structures (Supplementary Fig. S7), sparse effects of different strategies (Supplementary Fig. S8) and tuning processes (Supplementary Figs S9–S14) are also given in Supplementary Materials. The geometric interpretations of different approaches in 3D space and 2D projection from three representative views is shown in Figure 2c, and more details can be found in Supplementary Materials. Similar to $L_1$, GFS achieves sparsity and avoids overfitting, and moreover, the performance is improved by exploiting group information.

## 3. Results

### 3.1 The distribution of neoantigens' DNA loci in 3D genome

For all the peptides (both immunopositive and immunonegative) included in this study, we first generated a pool that contains all the peptide pairs. Then we classified all the peptide pairs in this pool into three categories: positive–positive pairs (Pos–Pos), negative–negative pairs (Neg–Neg) and positive–negative pairs (Pos–Neg). For each peptide pair, we computed contact frequencies for each Hi-C datasets, i.e. hESC, IMR90, GM12878, HUVEC, IMR90-Rao, NHEK, K562 and KBM7, respectively (Dixon *et al.*, 2012; Rao *et al.*, 2015). The contact frequency distribution of the three categories are shown in Figure 3a. It is demonstrated that on all the Hi-C datasets, immunopositive peptide pairs are more proximate to each other comparing to immunonegative peptide pairs; the corresponding *T*-test and Wilcoxon rank sum test *P*-values, i.e. Pos–Pos versus Neg–Neg, are all smaller than $10^{-99}$ and $10^{-18}$, respectively. This indicates that the immunopositive peptide's DNA loci tend to be more proximate in genome spatial space. We then computed the A/B compartment type (A: active; B: inactive) for each chromosomal region (bin), based on both Hi-C dataset and epigenetic markers, shown in Figure 3b and c. The whole genome contact maps of the eight Hi-C datasets are shown in Supplementary Figure S2 and the A/B compartment results of each chromosome are shown in Supplementary Figure S3. Then, we assigned the corresponding DNA loci of the positive and negative peptides with their A/B compartment type. We found that in certain chromosomes, immunopositive neoantigens tend to be located on compartment A, comparing to immunonegative neoantigens, as shown in Figure 3d and Supplementary Figure S4. This indicates that the DNA loci of the immunopositive or immunonegative peptides are positively correlated to chromosome compartment type, either A or B, depending on which chromosome.

We then developed a novel MD-based chromatin 3D modeling method and mapped the immunopositive and immunonegative peptides' corresponding chromosomal loci to the constructed 3D genome structure and calculated their radius distance to the nucleus center, as shown in Figure 3e. We found that the immunopositive peptide's corresponding loci tend to locate closer to the nuclear periphery (more far away from the nucleus center), compared to the immunonegative ones, as Figure 3f demonstrates. We found that by adopting the radius position information, the prediction power of the existing methods such as NetMHCPan and NetMHC can be elevated. In detail, prediction scores defined as $Y_{pred} = S_{NetMHCPan} \times r^2$ or $Y_{pred} = S_{NetMHC} \times r^2$ can significantly better discriminate the immunopositive peptides from the immunonegative peptides, comparing to using NetMHCPan or NetMHC alone. We thus believe that the DNA loci's radius positions of the immunopositive and immunonegative peptides are significantly differently distributed and can play an important role in predicting pMHC-I immunogenicity.

### 3.2 Peptide encoding and predictions

A reasonable and proper peptide-encoding strategy is a key to the downstream predictions as it can include and quantify more features that are plausibly related to the outcome. But by including more features into the prediction model, we also increase the risk of adding noisy (irrelevant) features into the feature pool and making the prediction prone to overfitting. To overcome this dilemma, we propose
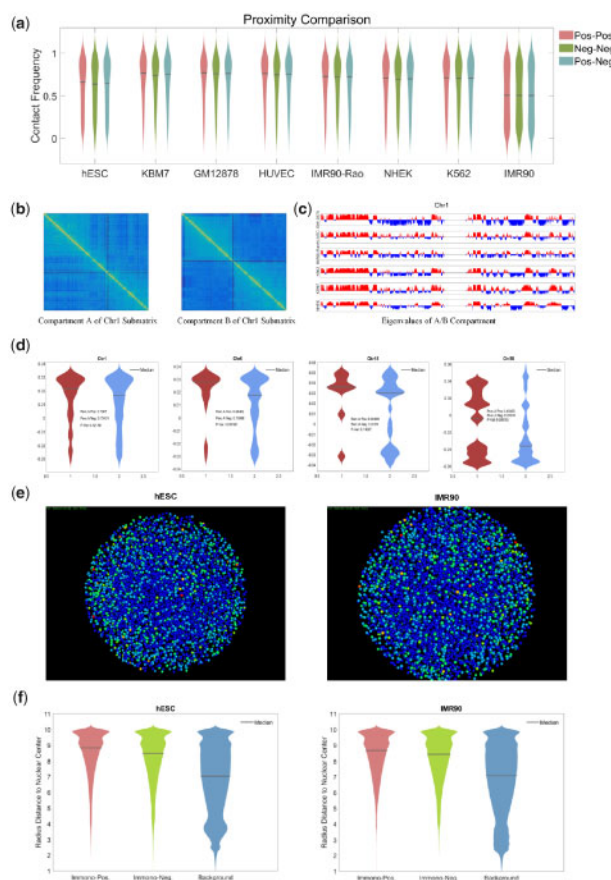


**Fig. 3.** The DNA loci of neoantigens in 3D genome. (**a**) Distribution of proximities between peptide pairs of different types. Immunopositive peptide pairs tend to be more proximate to each other comparing to immunonegative ones, while immunopositive–negative pairs lie in between (all the *P*-values of the *T*-test comparison are smaller than $10^{-99}$). (**b**) Illustration of Hi-C submatrices of compartment A and B on chromosome 1. (**c**) Illustration of eigenvalues of compartment A (red) and B (blue) on chromosome 1. (**d**) Comparison of percentages of immunopositive peptide belonging to compartment A (red) and immunonegative peptide belonging to compartment A (blue). (**e**) The 3D genome molding results based on hESC and IMR90 Hi-C datasets and the distribution of the DNA loci of immunopositive (yellow to red color spectrum, depending on positive occurrence on the same 500k bin) and immunonegative peptides (green). (**f**) Radius position comparison of the immunopositive and immunonegative peptides' DNA loci 3D genome. The positive loci (red) are significantly closer to the nuclear periphery (more far away from the nucleus center), compared to the immunonegative ones (green); they are all closer to the nuclear periphery comparing to the background distribution (blue). All *T*-test *P*-values are smaller than $10^{-99}$

to first enumerate as many features as possible and then perform feature selection within the training process of the prediction modeling. Previous neoantigen prediction methods adopted one or more coding schemes such as amino-acid (AA) composition, AA sparse coding, BLOSM, BLOMAP, and so on. In this work, based on the earlier observation that chromatin 3D information may significantly contribute to discriminating immunopositive peptides from immunonegative ones, we adopted this piece of information in the peptide-encoding strategy. In detail, the 3D coordinates and the radius positions of the Hi-C data based 3D modeling results, the HLA subtype encoding, the amino-acid compositions, the sparse coding, BLOMAP coding and BLOSUM coding of the peptides, the AAindex2 coding of the peptides are adopted and collected as features. At the end, we obtained a training matrix with 3909 peptides and 5459 features, shown in Figure 4a. Note that as Figure 4a demonstrates, there is no obvious pattern that a single feature or a group of features are correlated to the true label vector.
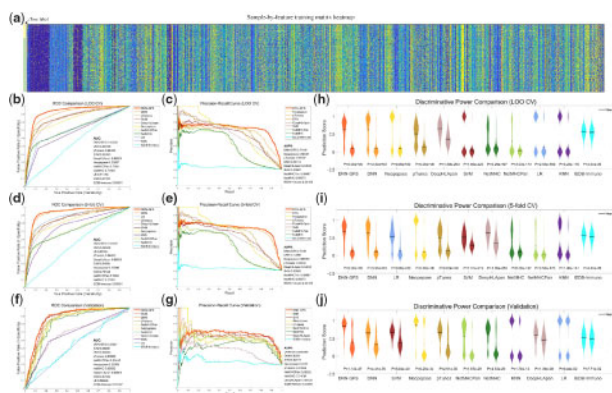
**Fig. 4.** Prediction results comparison. (**a**) The leftmost column vector indicates the true labels of the immunopositive (yellow) and immunonegative (green) for each of the 3909 peptides. The matrix heatmap indicates the columnwise normalized feature values of the 3909 peptides by 5459 features. (**b**) and (**c**) are the ROC plot comparison for DNN-GFS, DNN, SVM, LR, KNN, Neopesee, pTuneos, DeepHLApan, NetMHCpan, NetMHC and IEDB-immuno, under 5-fold and leave-one-out (LOO) cross-validation, respectively. (**d**) and (**e**) are the precision–recall plot comparison for different prediction methods under 5-fold and LOO cross-validation, respectively. (**f**) and (**g**) are the prediction score (normalized) distribution comparison for immunopositive (left violins) and immunonegative peptides (right violins); all the *P*-values of the *T*-tests are equal to or very close to zero

In theoretical deep neural network (DNN) studies, there have been plenty of evidences pointing to the fact that the majority of weights in most deep networks are redundant and may jeopardize the prediction accuracy (Han *et al.*, 2015; Sainath *et al.*, 2013; Scardapane *et al.*, 2017). It is possible to learn only a small percentage of the weights, while still preserving the prediction accuracy (Han *et al.*, 2015). Nevertheless, studies focusing on the input feature selection-based neural network are limited. Moreover, in the neoantigen prediction problem, the features that encode the peptides come in groups, e.g. the 3D coordinates $<x, y, z>$ of a peptide's DNA loci are in one group, or the sparse coding for an amino acid is a group of 20 binary features, etc. Therefore, when imposing feature selection on the DNN, it should be in a group fashion, i.e. features belonging to the same group should be either all selected or all rejected. The DNN-GFS is introduced in detail in the Section 2.

To compare the prediction efficacy, in addition to DNN-GFS, we also applied traditional L2 norm DNN, support vector machine (SVM), logistic regression (LR) and k-nearest neighbor (KNN) classifiers on the 5459 encoded feature matrix. Moreover, we included the widely adopted methods IEDB-immunogenicity, NetMHCpan and NetMHC into the comparison, as well as the most recent popular methods Neopesee, pTuneos and DeepHLApan. The comparison was conducted in the framework of both cross-validation (5-fold or leave-one-out) and validation alone. The ROC curves are shown in Figure 4b, d and f; the precision–recall curves are shown in Figure 4c, e and g; the prediction score distributions for the immunopositive and immunonegative samples are shown in Figure 4h, i and j. Note that, in the ten prediction methods, the KNN and LR output binary values, so for precision–recall curve comparison, we excluded them. Detailed prediction statistics are shown in Supplementary Tables S1-1, -2 and S2. As the comparison results demonstrate, the deep learning-based approaches DNN-GFS and DNN outperform the rest of the methods and DNN-GFS, due to its feature selection potency, is better than traditional DNN. The SVM, Neopesee, pTuneos and DeepHLApan are also effective methods and ranked second tier among the ten methods. Although NetMHC and NetMHCpan were initially designed to predict peptide–MHC binding affinity, their capability in predicting neoantigen cannot be neglected and they are ranked third tier. The logistic regression and KNN classifiers, although performs reasonably well in cross-validation experiments, are not very stable when applied on the standalone validation set. The IEDB-immunogenicity prediction method, does not catch up with other prediction methods, possibly

due to the fact that the immunogenicity scoring function is too simple to capture subtle sequence features that only advanced non-linear machine learning methods can. We also implemented other well-known sparse learning neural network models and compared their efficacy with DNN-GFS, as introduced in Supplementary Tables S1 and S2, and the results indicate that DNN-GFS outperforms existing sparse neural network methods in terms of prediction statistics.

### 3.3 Features selected by DNN-GFS
Based on the whole training dataset, the DNN-GFS model selected 2693 features out of the 5459 features, achieving a feature sparsity ratio 49.33%. Features belonging to the same group are either all selected or all excluded. Among the selected features, all the 3D genome-related features are selected, including radius position, HLA subtype, 3D coordinates of peptides' DNA loci, and so on. For HLA subtype-encoding, all features are selected and cross-validation performance is improved about 3–4% compared to dataset of not containing HLA subtype information, which illustrates their importance. For nine AA peptides, the sparse coding of the peptide's position one to five and seven to eight are all selected but not position six and nine. BLOSUM coding features are all excluded while BLOMAP coding features for AA position one to four are selected. Except AA position five, other side chain polarity features are all selected, and side-chain charge features for position one to three are selected. For the hydropathy features, AA position five and nine are selected, and for molecular weight, feature of AA position two, six and nine are selected. Other selected features are mostly AAindex2-related features. The DNN-GFS model thus suggests that the combination of these grouped features play an important role in building the prediction model and we believe that the importance of these features in neoantigen prediction is worth further investigating. Detailed feature selection and model sparsity analyses can be found in Supplementary Materials.

## 4. Discussion
From the association study of peptides' immunogenicity and their 3D genome information, we found that immunopositive peptides' DNA loci tend to be more proximate to each other and locate closer to the nuclear periphery, i.e. greater radius value to the nuclear center, comparing to immunonegative ones. This implies that if a non-synonymous mutation happens closer to some non-synonymous mutation that were already proven to produce immunopositive peptides, or if it is located closer to the nuclear periphery, the mutation is more likely to generate immunopositive neoantigens. This association can be further enhanced if the A/B compartment information of the mutation is provided. In practice, the whole genome spatial organization is more conserved across different cell types and even in mutated cancer cells. While A/B compartment characters of certain chromatin regions may flip across cell lines or in cancer cells, i.e. more transient, we only adopted the 3D coordinates and radius position of peptides' DNA loci in the prediction model, but if the A/B compartment information can also be included if the real-time cancer cell's chromatin 3D experiment can be performed in the future.

To explain such intriguing relationship between 3D genome and the neoantigens' immunogenicity, factors of at least three aspects should be considered: First, the non-random nature of coding sequence distribution in 3D genome: during evolution, wild-type coding sequences where neoantigens of different immunogenicity characters originate are located in different regions of the nucleus (Gorkin *et al.*, 2019; Svozil *et al.*, 2008). Second, the gene expressions affected by 3D genome: the missense mutations need to be transcribed to generate potential neoantigens and the gene expressions are known to be affected by high-order genome organization (Gorkin *et al.*, 2014). Third, the non-random occurrences of somatic mutations in 3D genome: previous discoveries indicated that somatic mutations may not occur at random, and we systematically studied and discovered in our prior work that comutations may

occur in a spatial clustering fashion in genome 3D space (spatial comutation hotspot, SCH), possibly due to abnormal chemical concentration or a systematic DNA repair protein failure at certain chromatin 3D loci (Shi *et al.*, 2016). This leads to a straightforward hypothesis that mutations in different chromosomal regions may carry different immunogenicity character, affected by wild-type coding sequences, somatic mutation patterns and gene transcriptions. We thus believe that it is worth considering these aspects when studying the underlying mechanism of how high-order genome organization affects neoantigens' immunogenicity. For example, to explain our discovery that immunopositive neoantigens' corresponding DNA sequences tend to locate closer to nucleus periphery (greater radius to the nuclear center), one may consider the fact that their transcribed missense mRNAs enter cytoplasm more easily (a shorter path from transcription loci to nuclear envelope). Core genes during evolution, if mutated, require stronger TCR responses, because otherwise it causes greater cancerous impact, and these genes are usually expressed across cell types and their distribution in 3D genome also worth further study. Therefore, it is also worth to investigate the relationship between neoantigen immunogenicity and gene evolutionary essentiality in the perspective of high-order genome conformation.

When building the prediction models, due to the fact that most MHC-I presented peptides are of nine amino-acid long, the features we used to encode the peptides are all based on 9mer peptides, and the predictions are targeted on the 9mers as well. Nevertheless, our approach is not restricted to 9mers and can be easily extended to peptides of other length. For example, if a target peptide is longer than nine amino acids, a sliding window of length nine can be used to enumerate all possible 9mers, and the prediction score can be estimated by taking the maximum or average of each individual scores. In the cases where a target peptide is shorter than length nine, we only need to consider length eight as peptides presented by MHC-I shorter than or equal to length seven is very rare. So, for the 8mer cases, we can compensate an extra amino acid to the beginning or to the end of the sequence and enumerate all possible peptides and again take the maximum or average of each individual one's prediction score.

Most existing machine learning algorithms for the classification problem usually assume that the feature across different training examples is independent and obey the same distribution, and the links among them are usually neglected which is not reasonable for an unbalanced problem. In many real-world applications, however, the small sample issue is ubiquitous and the features are usually correlated. The DNN-GFS developed here provides a new way of exploiting these links for feature selection in addition to traditional neural networks. In the machine learning area, quite a few studies have exploited introducing sparse regularization into DNN framework, but most of these models only focus on reducing complexity of the network as a whole, resulting in pruning edges and nodes of the network, but not specifically targeting on the input layer, i.e. the input feature vector. In this work therefore, due to the scenario that peptides are represented in an ensemble encoding which may introduce noise or redundant features into the learning process, the proposed DNN-GFS model focus on reducing the features of the input layer. Moreover, due to the nature that certain features are grouped and should be either all selected or all rejected, we considered selecting features in a grouped fashion in the model, by imposing group-specific regularization. As shown in Figure 4, the DNN-GFS model not only exceeds the widely adopted methods NetMHCpan and NetMHC, but also exceeds other existing machine learning methods such as DNN, LR, SVM and KNN that are performed based on the same 5459 feature encoding strategy. Moreover, DNN-GFS outperforms other sparse learning DNN models, as shown in Supplementary Tables S3–S10. This agrees with our conjecture that DNN-GFS is a better DNN heuristic designed specifically for the neoantigen prediction in the specific 5459 encoding scenario. Although DNN-GFS outperforms the widely adopted NetMHCpan and NetMHC methods to a large extend, due to its ability of capturing subtle non-linear relationships of features in a grouped fashion, the prediction power can be further improved once more immunogenicity training data are provided, especially for each HLA subtypes. We also believe that DNN-GFS can also be applied in other problems where GFS is demanded.

To facilitate practical usage, we developed a webserver deepAntigen (Supplementary Fig. S1). In the current version, if the end user only provides sequencing result vcf file, the candidate peptides will be generated by only considering non-synonymous point mutations, i.e. 9mer peptides surrounding the mutated amino acid, while small insertions or deletions (INDEL) can also be considered as *rankPep* function is independent and user can provide their own plausible peptides for prediction. For the prediction method, we suggest to use DNN-GFS as its power of discriminating immunopositive peptides from immunonegative ones are most potent, but other machine learning approaches can also be considered and the consensus result maybe of more interest to an end user.

Although the mechanism of under what conditions certain specific neoantigens activate T-cell immunogenicity is still under studying, this work focuses on the machine learning challenge of effectively and efficiently predict/prioritize immunopositive neoantigens. We found that the spatial distributions of the immunopositive and immunonegative peptides' corresponding DNA loci follow different patterns, i.e. immunopositive peptides' DNA loci tend to be located more proximate to the nuclear periphery and tend to be more clustered in 3D genome space, compared with immunonegative peptides' DNA loci; the peptides' DNA loci distribution is also related to the A/B compartment of the chromatin. It is therefore salient that utilizing the 3D genome information of the peptides' corresponding DNA loci can significantly contribute to the prediction of immunopositive neoantigens. To utilize the most of 3D genome information, we customized a DNN-GFS model, which takes not only the 3D genome information, but also a combinatorial peptide sequence features represented by an ensemble peptide-encoding strategy. The DNN-GFS selected 3D genome related features as well as some other important peptide sequence features and position specific amino-acid features; the comparison studies demonstrated that DNN-GFS outperforms the widely adopted methods NetMHCpan and NetMHC, and other machine learning prediction models including DNN, SVM, LR and KNN. DNN-GFS is implemented in the webserver deepAntigen along with other machine learning methods. To the best of our knowledge, this is the first time that the DNA origins' 3D genome perspective is considered in the neoantigen study and we hope that our work contributes novel insights to neoantigen study and eventually benefits personalized cancer immunotherapy. Although close-up studies are needed to uncover the relationship between 3D genome and neoantigen immunogenicity, in this work, we only demonstrate the contributes of 3D genome information in more accurate neoantigen prediction, as well as providing plausible explanation that it is evolution that places sequences of different immunogenicity characters in different locations in the 3D genome while different locations are prone to mutations of different causes.

## References

Berger,M.F. *et al.* (2011) The genomic complexity of primary human prostate cancer. *Nature*, **470**, 214–220.

Blaha,D.T. *et al.* (2019) High-throughput stability screening of neoantigen/HLA complexes improves immunogenicity predictions. *Cancer Immunol. Res.*, **7**, 50–61.

Blevins,S.J. *et al.* (2016) How structural adaptability exists alongside HLA-A2 bias in the human alpha beta TCR repertoire. *Proc. Natl. Acad. Sci. USA*, **113**, E1276–E1285.

Boratyn,G.M. *et al.* (2013) BLAST: a more efficient report with usability improvements. *Nucleic Acids Res.*, **41**, W29–W33.

Branco,M.R. and Pombo,A. (2006) Intermingling of chromosome territories in interphase suggests role in translocations and transcription-dependent associations. *PLoS Biol.*, **4**, e138.

Castro,A. *et al.* (2019) Elevated neoantigen levels in tumors with somatic mutations in the HLA-A, HLA-B, HLA-C and B2M genes. *BMC Med. Genomics*, **12**, 107.

Dettmers,T. and Zettlemoyer,L. (2019) Sparse networks from scratch: faster training without losing performance. *arXiv preprint arXiv : 1907.04840*.

Dintzis,H.M. *et al.* (1976) Molecular determinants of immunogenicity: the immunon model of immune response. *Proc. Natl. Acad. Sci. USA*, **73**, 3671–3675.

Dixon,J.R. *et al.* (2012) Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature*, **485**, 376–380.

Engreitz,J.M. *et al.* (2012) Three-dimensional genome architecture influences partner selection for chromosomal translocations in human disease. *PLoS One*, **7**, e44196.

Friedman,J. *et al.* (2010) A note on the group lasso and a sparse group lasso. *arXiv preprint arXiv : 1001.0736*.

Goodfellow,I. *et al.* (2016) *Deep Learning*. MIT Press. https: //www.deeplearningbook.org/.

Gorkin,D.U. *et al.* (2014) The 3D genome in transcriptional regulation and pluripotency. *Cell Stem Cell*, **14**, 762–775.

Gorkin,D.U. *et al.* (2019) Common DNA sequence variation influences 3-dimensional conformation of the human genome. *Genome Biol.*, **20**, 255.

Gupta,S.K. *et al.* (2016) Personalized cancer immunotherapy using systems medicine approaches. *Brief. Bioinf.*, **17**, 453–467.

Hackl,H. *et al.* (2016) Computational genomics tools for dissecting tumour-immune cell interactions. *Nat. Rev. Genet.*, **17**, 441–458.

Han,S. *et al.* (2015) Learning both weights and connections for efficient neural networks. In *Neural Information Processing Systems (NIPS)*, pp. 1135–1143.

Harndahl,M. *et al.* (2012) Peptide-MHC class I stability is a better predictor than peptide affinity of CTL immunogenicity. *Eur. J. Immunol.*, **42**, 1405–1416.

Huang,D.W. *et al.* (2009) Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat. Protoc.*, **4**, 44–57.

Jurtz,V. *et al.* (2017) NetMHCpan-4.0: improved peptide-MHC class I interaction predictions integrating eluted ligand and peptide binding affinity data. *J. Immunol.*, **199**, 3360–3368.

Karosiene,E. *et al.* (2012) NetMHCcons: a consensus method for the major histocompatibility complex class I predictions. *Immunogenetics*, **64**, 177–186.

Karosiene,E. *et al.* (2013) NetMHCIIpan-3.0, a common pan-specific MHC class II prediction method including all three human MHC class II isotypes, HLA-DR, HLA-DP and HLA-DQ. *Immunogenetics*, **65**, 711–724.

Kim,S. *et al.* (2018) Neopepsee: accurate genome-level prediction of neoantigens by harnessing sequence and amino acid immunogenicity information. *Ann. Oncol.*, **29**, 1030–1036.

Lundegaard,C. *et al.* (2008) NetMHC-3.0: accurate web accessible predictions of human, mouse and monkey MHC class I affinities for peptides of length 8–11. *Nucleic Acids Res.*, **36**, W509–W512.

Lundegaard,C. *et al.* (2011) Prediction of epitopes using neural network based methods. *J. Immunol. Methods*, **374**, 26–34.

Mani,R.S. *et al.* (2009) Induced chromosomal proximity and gene fusions in prostate cancer. *Science*, **326**, 1230–1230.

Mathas,S. *et al.* (2009) Gene deregulation and spatial genome reorganization near breakpoints prior to formation of translocations in anaplastic large cell lymphoma. *Proc. Natl. Acad. Sci. USA*, **106**, 5831–5836.

Meaburn,K.J. *et al.* (2007) Spatial genome organization in the formation of chromosomal translocations. *Semin. Cancer Biol.*, **17**, 80–90.

Moutaftsi,M. *et al.* (2006) A consensus epitope prediction approach identifies the breadth of murine TCD8+-cell responses to *Vaccinia virus*. *Nat. Biotechnol.*, **24**, 817–819.

Nielsen,M. and Andreatta,M. (2016) NetMHCpan-3.0; improved prediction of binding to MHC class I molecules integrating information from multiple receptor and peptide length datasets. *Genome Med.*, **8**, 33.

Nikiforova,M.N. (2000) Proximity of chromosomal loci that participate in radiation-induced rearrangements in human cells. *Science*, **290**, 138–141.

O'Brien,S.G. *et al.* (2003) Imatinib compared with interferon and low-dose cytarabine for newly diagnosed chronic-phase chronic myeloid leukemia. *N. Engl. J. Med.*, **348**, 994–1004.

Ott,P.A. *et al.* (2017) An immunogenic personal neoantigen vaccine for patients with melanoma. *Nature*, **547**, 217–221.

Parker,K.C. *et al.* (1994) Scheme for ranking potential HLA-A2 binding peptides based on independent binding of individual peptide side-chains. *J. Immunol.*, **152**, 163–175.

Peters,B. *et al.* (2006) A community resource benchmarking predictions of peptide binding to MHC-I molecules. *PLoS Comput. Biol.*, **2**, e65.

Prior,L. *et al.* (2019) Genomic profiling of a dedifferentiated mucosal melanoma following exposure to immunotherapy. *Melanoma Res.*, **30**, 213–218.

Rao,S.S.P. *et al.* (2015) A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell*, **162**, 687–688.

Rasmussen,M. *et al.* (2016) Pan-specific prediction of peptide-MHC class I complex stability, a correlate of T cell immunogenicity. *J. Immunol.*, **197**, 1517–1524.

Riley,T.P. *et al.* (2018) T cell receptor cross-reactivity expanded by dramatic peptide-MHC adaptability. *Nat. Chem. Biol.*, **14**, 934–942.

Riley,T.P. *et al.* (2019) Structure based prediction of neoantigen immunogenicity. *Front. Immunol.*, **10**, 1–14.

Roix,J.J. *et al.* (2003) Spatial proximity of translocation-prone gene loci in human lymphomas. *Nat. Genet.*, **34**, 287–291.

Sainath,T.N. *et al.* (2013) Low-rank matrix factorization for deep neural network training with high-dimensional output targets. 2013 *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 6655–6659.

Scardapane,S. *et al.* (2017) Group sparse regularization for deep neural networks. *Neurocomputing*, **241**, 81–89.

Schuler,M.M. *et al.* (2007) SYFPEITHI: database for searching and T-cell epitope prediction. *Methods Mol. Biol.*, **409**, 75–93.

Schumacher,T.N. and Schreiber,R.D. (2015) Neoantigens in cancer immunotherapy. *Science*, **348**, 69–74.

Shi,Y. *et al.* (2016) Chromatin accessibility contributes to simultaneous mutations of cancer genes. *Sci. Rep.*, **6**, 35270.

Simon,N. *et al.* (2013) A sparse-group lasso. *J. Comput. Graph. Stat.*, **22**, 231–245.

Simon,N. and Tibshirani,R. (2012) Standardization and the group lasso penalty. *Stat. Sin.*, **22**, 983.

Sompayrac,L. (2019) *How the Immune System Works*. Wiley-Blackwell, Hoboken, NJ.

Svozil,D. *et al.* (2008) DNA conformations and their sequence preferences. *Nucleic Acids Res.*, **36**, 3690–3706.

Trolle,T. *et al.* (2015) Automated benchmarking of peptide-MHC class I binding predictions. *Bioinformatics*, **31**, 2174–2181.

Vita,R. *et al.* (2019) The Immune Epitope Database (IEDB): 2018 update. *Nucleic Acids Res.*, **47**, D339–D343.

Volkov,N.M. *et al.* (2020) Efficacy of immune checkpoint blockade in MUTYH-associated hereditary colorectal cancer. *Invest. New Drugs*, **38**, 894–898.

Wang,P. *et al.* (2008) A systematic assessment of MHC class II peptide binding predictions and evaluation of a consensus approach. *PLoS Comput. Biol.*, **4**, e1000048.

Wang,Y. *et al.* (2017) How an alloreactive T-cell receptor achieves peptide and MHC specificity. *Proc. Natl. Acad. Sci. USA*, **114**, E4792–E4801.

Wijchers,P.J. and de Laat,W. (2011) Genome organization influences partner selection for chromosomal rearrangements. *Trends Genet.*, **27**, 63–71.

Wu,J. *et al.* (2019) DeepHLApan: a deep learning approach for neoantigen prediction considering both HLA-peptide binding and immunogenicity. *Front. Immunol.*, **10**, 2559.

Xie,W.J. *et al.* (2017) Structural modeling of chromatin integrates genome features and reveals chromosome folding principle. *Sci. Rep.*, **7**, 2818.

Yarchoan,M. *et al.* (2017) Targeting neoantigens to augment antitumour immunity. *Nat. Rev. Cancer*, **17**, 209–222.

Zhang,G.L. *et al.* (2011) Dana-Farber repository for machine learning in immunology. *J. Immunol. Methods*, **374**, 18–25.

Zhou,C. *et al.* (2019) pTuneos: prioritizing tumor neoantigens from next-generation sequencing data. *Genome Med.*, **11**, 67.