

# Machine Learning in T- and B-Cell Epitope Prediction

Subjects: [Computer Science](#), [Interdisciplinary Applications](#) | [Immunology](#)

Contributor: Syed Nisar Hussain Bukhari

An antigenic determinant (AD) is a portion of an antigen molecule known as an epitope that is recognized by the human immune system, specifically by antibodies or T and B cells. Recognition of epitopes is considered important in EBPV design to contain pandemics, epidemics, and endemics due to the outbreak of infectious diseases. To design an effective and viable EBPV against different strains of a pathogen, it is important to identify the putative T- and B-cell epitopes. Using the wet-lab experimental approach to identify these epitopes is time-consuming and costly because the experimental screening of a vast number of potential epitope candidates is required. Fortunately, various available machine learning (ML)-based prediction methods have reduced the burden related to the epitope mapping process by decreasing the potential epitope candidate list for experimental trials. Moreover, these methods are also cost-effective, scalable, and fast.

machine learning   epitopes   B-Cell   T-Cell   antigenic determinant   antigen   antibody

immune-relevant determinants   epitope-based peptide vaccine   SARS-CoV-2   COVID-19

ensemble model

## 1. Introduction

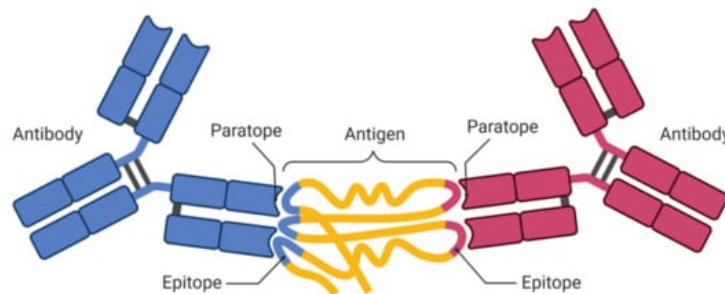
An antigenic determinant (AD) is a portion of an antigen molecule known as an epitope that is recognized by the human immune system, specifically by antibodies or T and B cells <sup>[1]</sup>. Recognition of epitopes is considered important in EBPV design to contain pandemics, epidemics, and endemics due to the outbreak of infectious diseases. The ongoing COVID-19 pandemic due to the SARS-CoV-2 outbreak is the latest among the major pandemics that have occurred in the last decade <sup>[1]</sup>. COVID-19 can be severe and has caused millions of deaths around the world. It is a respiratory illness and affects people according to the physiology and immune system of the human body. Affected people mostly develop mild to moderate illness and recover without hospitalization <sup>[1][2]</sup>. While the progress in COVID-19 vaccine design so far is remarkable, successfully vaccinating the worldwide population entails numerous hurdles, from manufacturing to distribution and deployment, and, most crucially, acceptability.

Due to the rate at which SARS-CoV-2 is circulating in the population, thereby causing unprecedented infections, its chances of mutating more and more have increased by now. The variant B.1.617.2, named Delta <sup>[3]</sup>, first identified during a serious wave of COVID-19 infections in India in April and May 2021 <sup>[4]</sup>, was declared a variant of concern (VOC) by the "US Centers for Disease Control and Prevention (CDC)" on 15 June 2021 <sup>[5]</sup>. Due to its partial resistance to existing vaccines, the infected cases per day increased to over 400,000 <sup>[6]</sup>. A study conducted by the Chinese Academy of Medical Sciences confirmed that viral loads in Delta infections are approximately 1000 times higher than those in previous SARS-CoV-2 variants <sup>[7]</sup>. The Mu variant, also known as B.1.621 <sup>[8]</sup>, first identified in January 2021 in Colombia, was declared a "variant of interest" (VOI) on 26 August 2021 by the European Centre for Disease Prevention and Control (ECDC) <sup>[9]</sup>. On August 30, "the Mu variant was added to the World Health Organization's (WHO's) watch list after being found to have a constellation of mutations that indicate potential properties of immune escape" <sup>[9]</sup>. The most recent variant, B.1.1.529, named Omicron, was first reported to WHO from South Africa on 24 November 2021 <sup>[8]</sup>. On 26 November 2021, WHO designated the variant B.1.1.529 a VOC on the advice of the Technical Advisory Group on Virus Evolution (TAG-VE) <sup>[8]</sup>. The hotspot of SARS-CoV-2 mutations is the spike S protein. The spike protein enables the pathogen to infect cells and is the basis for the majority of the vaccines. In <sup>[9]</sup>, it has been reported that "out of 10333 spike protein sequences analyzed, 8155 proteins comprised one or more mutations. A total of 9654 mutations were observed that correspond to 400 distinct mutation sites. The receptor binding domain (RBD) which is involved in the interactions with human angiotensin-converting enzyme-2 (ACE-2) receptor and causes infection leading to the COVID-19 comprised 44 mutations that included residues within 3.2 Å interacting distance from the ACE-2 receptor".

### 1.1. Epitopes and Paratopes

An antigen is any substance that causes the immune system to produce antibodies against it. Its molecules are large biological polymers and introduce various molecular attributes that act as interaction sites between antibodies, T<sub>H</sub> cells and B cells, and antigen molecules. These interaction sites are called epitopes <sup>[10][11][12]</sup>. Epitopes are of two types: B-cell epitopes

(BCEs) and T-cell epitopes (TCEs). The fragment of an antigen that is attached to an antibody is called the B-cell epitope [13]. The BCEs are recognized by B cells and comprise a solvent region that is exposed to an antigen. On the other hand, T cells have a receptor on their surface, known as the T-cell receptor (TCR) [13]. When presented on the surfaces of APCs that are linked to MHC molecules, the TCR aids in antigen recognition. TCEs identified by CD8 and CD4 T cells are represented by MHC class I (MHC I) and class II (MHC II) molecules, respectively [13]. **Figure 1** shows an antibody containing two paratopes, indicating that these two paratopes can bind to two pathogens [14][15]. Chemical interactions between epitopes and paratopes that promote antigen–antibody binding are non-covalent [16][17][18].



**Figure 1.** Antigen recognition by antibodies.

## 1.2. Need for T- and B-Cell Epitope Prediction

The identification of epitopes is of great importance for many reasons, including EBPV design, antibody production, and immunodiagnostic tests. They also play a crucial role in activating the human immune system. Among the reasons listed, EBPV design is important for researchers, biologists, and scientists because there are numerous drawbacks to using whole-organism vaccines, particularly in immunocompromised patients [19][20]. EBPVs can be utilized to overcome the issues associated with heterogeneous and multicomponent vaccines and are seen as an alternative to traditional vaccines. They can act as powerful alternatives to conventional vaccines due to their low production cost, having less reactogenic and allergenic responses. A well-trained ML model of experimentally determined epitopes and non-epitopes can identify potential epitopes as vaccine candidates quickly and efficiently and can reduce the burden related to the epitope mapping process by decreasing the potential epitope candidate list for experimental trials. Using the wet-lab experimental approach to identify these epitopes is time-consuming and costly because the experimental screening of a vast number of potential epitope candidates is required. However, epitope prediction methods based on ML can prove to be cost-effective, scalable, and fast. The most recent vaccine technology is based on RNA vaccines, which have the distinct advantage of being simple to design and manufacture. Epitopes are critical, but often overlooked, for boosting the effectiveness of RNA vaccines. Although RNA vaccines can encode any gene of interest, even the most recent designs commonly encode sequences of original genes from the natural virus. Epitope prediction can be useful in assisting RNA vaccine design by guiding the sequence design and vaccine structure. RNA (mRNA) vaccines, on the other hand, can benefit from epitope-based design approaches, in which both B-cell and T-cell epitopes can be used for vaccine design. The epitope properties determine whether or not the RNA vaccine will elicit an immune response and which types of responses will be elicited.

## 2. ML-Based Studies for the Prediction of T- and B-Cell Epitopes

ML is concerned with the automated learning of machines that is not explicitly programmed. It focuses on making data-driven predictions and has several applications in bioinformatics [21]. Bioinformatics deals with applying computational techniques to derive knowledge from biological data. It covers the collection, retrieval, storage, manipulation, and data modeling for analysis or prediction using various algorithms and software [21]. Earlier, one had to explicitly program bioinformatics algorithms, which was an extremely laborious task for predicting protein structures [21]. However, with the advent of ML algorithms, such problems have become much easier to solve. In recent years, the exponential growth of T- and B-cell epitope data has become the primary motivation for researchers to develop ML-based methods for the prediction of ADs or IRDs, i.e., B- and T-cell epitopes. ML applied to experimentally determined peptide sequence data of pathogens (virus, bacteria, etc.) opens up new frontiers for areas such as EBPV design, antibody production, and immunodiagnostic tests. The ML-based in silico approach has emerged as a promising field for epitope prediction [22]. Accordingly, various ML-based studies and methods exist that utilize the physicochemical properties of amino acids as features or descriptors for the prediction of epitopes (Table 1).

**Table 1.** Existing studies for T- and B-cell epitope prediction.

Study Conducted	Methodology Adopted	Strengths/Limitations
T. Liu et al. <a href="#">[23]</a>	A feedforward deep neural network-based ensemble of 11 classifiers was created to predict BCEs. IEDB was used to obtain the BCE peptide dataset. On the test set, the model was evaluated using the AUROC metric.	Model reports peptide as an epitope if classified by all 11 classifiers. It would provide the best results if simple majority voting was used for classification.
Fatoba, A. J. et al. <a href="#">[24]</a>	In <a href="#">[24]</a> , potential epitope-based vaccine candidates were explored. After retrieving 600 genome sequences of SARS-CoV-2 from the ViPR repository, CD8+ and CD4+ epitopes and B-cell (linear) epitopes were generated and screened for immunogenicity, antigenicity, and non-allergenicity.	The results of <a href="#">[25]</a> reported 19 candidate T-cell epitopes (CD8+), which were found to overlap strongly with 8 B-cell epitopes. The results provide the basis for an experimental design for a suitable peptide vaccine against SARS-CoV-2.
R. Moody et al. <a href="#">[26]</a>	Authors used IEDB prediction tools for predicting B-cell epitopes and those with high scores in terms of prediction were selected as candidate epitopes. The epitopes were then matched to human proteins using NCBI Blast technology.	The findings showed eleven (11) novel B-cell epitopes in the host that were capable of explaining key elements of COVID-19 extrapulmonary disease that previous research had not been able to explain.
Jespersen MC et al. <a href="#">[27]</a>	The authors employed feedforward neural networks (FFNN) with two hidden layers, each with 25 neurons, an activation function (sigmoid) at all neurons, and an ADAM as an optimizing function to predict antibody-specific epitopes (B cell) or epitope targets of provided cognate antibodies. The dataset was obtained from the IEDB database. PCA was used for dimensionality reduction before the model was trained.	It was shown that a simple set of attributes retrieved from the cognate antibody boosted the rate of accuracy in predicting individual epitopes. Furthermore, sophisticated features such as Zernike Moments can improve the model's predictive potential. When compared to DiscoTope 2.0, this model performs better in finding patches overlapping with an actual patch of an epitope in cross-validation and on an independent dataset.
Ling-yun Liu et al. <a href="#">[28]</a>	The authors used PCA and RNN networks. They converted the physicochemical properties into digital vectors, intending to have high-dimensional feature space, and later PCA was applied to process them. The output from PCA was used as an input to the RNN for predicting epitopes.	Prediction results obtained by this process demonstrated that PCA reduced dimensions, but at the same time, original features of the main component were retained, and the rate of prediction was also improved.
Bin Cheng et al. <a href="#">[29]</a>	Authors introduced a novel scale to measure feature importance, called the relevance of amino acid pair (RAAP). RAAP was calculated by decomposing the sequences of amino acids based on their physicochemical properties.	The successful prediction rate was drastically improved here by using LSTM. It does not suffer from gradient instability and is good enough for textual classification sequences. Fivefold cross-validation was used to test and validate the models.
Balachandran Manavalan et al. <a href="#">[30]</a>	Here, a non-redundant dataset was constructed containing 5500 BCEs experimentally validated, and 6893 non-B-cell epitopes were retrieved from IEDB. Then, an ensemble model to predict B-cell epitopes based on ERT (extremely randomized tree) and a classifier called GB (gradient boosting) was developed. The model works based on the physicochemical properties, AA composition, and combination of dipeptides and PCP as the input features.	After performing cross-validation on a benchmark dataset, it was shown that this model performed far better than the individual classifiers such as ERT and GB, with an MCC (Matthews correlation coefficient) of 0.454.
Yuh-Jyh Hu et al. <a href="#">[31]</a>	A cost-sensitive strategy based on bagging MDT was suggested, which integrates two ensemble-based learning algorithms. Without employing the prediction of a pre-trained single predictor, it makes it independent of multiple prediction tools. It can also learn a meta-classification architecture with varied data, without being constrained by a particular hierarchy.	It was demonstrated that the performance of prediction is superior as compared to a single epitope predictor. However, epitope prediction based on meta-learning is purely dependent upon the predictive strength of various other pre-trained linear and conformational epitope prediction tools, which cannot be retained directly by users. Hence, this limits the flexibility and applicability of these meta-classifiers.
Jing Ren et al. <a href="#">[32]</a>	The authors proposed a novel staged heterogeneity-based learning model. The model learns both heterogeneity and characteristics of data in a phased manner to identify residue of antigens of conformational B-cell type epitopes that are heterogeneous, purely based on sequences of antigens. In the first stage, the model is made to learn the generic epitope pattern with propensities, and in the second stage, the same model is made to	It was demonstrated that if heterogeneity was learned well, the transferability of the model improved remarkably in handling new data. It was tested and validated on two different datasets: one with epitopes determined experimentally and another with computationally defined. It showed outstanding performance that was around twice that of existing predictors, including CBTOPE.

Study Conducted	Methodology Adopted	Strengths/Limitations
	learn the complementarity of the propensities used in the first stage, which is heterogeneous but this time on a small dataset of experimentally verified epitopes.	
Georgios A. et al. [33]	A novel method, "SEPIa", has been proposed here to predict B-cell epitopes from protein sequences and is sufficiently faster, and it can also be applied to large-scale datasets. The model is the combination of two classifiers, random forest and naïve Bayes algorithm.	The average prediction accuracy of SEPIa is limited. The AUC score is 0.65 in both 10-fold cross-validation and on the independent test dataset, which is higher than other approaches tested on the same test dataset.
Gene Sher et al. [25]	Authors proposed a novel, analytically trained DREEP (Deep Ridge Regressed Epitope Predictor) based on string kernels using a deep neural network tailored to predict continuous epitopes.	The model was tested with input as long sequences of proteins from datasets such as AntiJen, Pellequer, and HIV. The results were compared with epitope predictors such as DMNLBE, LBtope, etc. Using the area under the curve (AUC) metric, the model achieved performance improvements over SARS by 13.7%, HIV by 8.9%, and Pellequer by 1.5%.
Wen Zhang et al. [34]	Authors attempted to differentiate immunogenic epitopes from non-immunogenic epitopes based purely on their primary structure. To effectively utilize various features, an ensemble method based on a genetic algorithm was proposed.	The model was tested on two benchmark datasets: IMMA2, PAAQD. The model was compared with methods such as POPI, PAAQD, and POPISK, which are considered state-of-the-art in nature. The model performed better, with an AUC score on IMMA2 of 0.846 and 0.829 on PAAQD.
Wei Zheng et al. [35]	The authors used ensemble learning to improve the prediction of BCEs. Their ensemble method combined twelve SVMs. To handle imbalanced datasets, resampling and AdaBoost methods were used.	The proposed ensemble model achieved an AUC score of 0.642–0.672 on the training dataset with five-fold cross-validation and an AUC score of 0.579–0.604 on the test dataset.
Jian Zhang et al. [36]	To predict antigenic determinants, the authors devised a cost-sensitive ensemble approach, and a spatial clustering-based algorithm was used to identify probable epitopes.	The model performed admirably in terms of prediction. AUC scores of 0.721 and 0.703 were obtained using leave-one-out cross-validation (LOOCV) on two benchmark datasets: bound and unbound.
Kavitha K V et al. [37]	PCA was used to reduce dimensions and to filter out the essential features; for prediction purposes, a random forest algorithm was used.	Experimental results showed that the random forest-based classifier had an improved prediction accuracy rate as compared to BCPred, AAP, etc.
Wen Zhang et al. [38]	The authors used sequence-derived features and developed an ensemble model based on random forest to predict epitopes accurately.	The model was evaluated using the leave-one-out cross-validation procedure, and an AUC score of 0.687 and 0.651 on bound and unbound datasets was obtained.
Ping Chen et al. [39]	Authors reviewed various prediction models for epitopes, such as models based on SVM, neural network, random forest, etc., to defend computational approaches in the prediction of epitopes as in silico methods require a lot of effort and time.	Apart from defending the computational approaches, it was also concluded that there is a limitation to current models as it is impossible to devise an exact model without having complete knowledge of the immune system, and current models are simply best at approximation.
Claus Lundegaard et al. [40]	Here, an artificial neural network was used. The standard feedforward neural network with backpropagation was employed to predict epitopes. The dataset was retrieved from the SYFPEITHI database.	The model efficiently and accurately predicts MHC class I type peptides and outperforms the existing methods.

## References

- Immunology Guidebook|ScienceDirect. Available online: <https://www.sciencedirect.com/book/9780121983826/immunology-guidebook> (accessed on 25 September 2021).
- COVID Live Update: 270,426,226 Cases and 5,321,864 Deaths from the Coronavirus—Worldometer. Available online: <https://www.worldometers.info/coronavirus/> (accessed on 10 December 2021).
- Centers for Disease Control and Prevention (CDC). SARS-CoV-2 Variant Classifications and Definitions. The primary basis for T-cell epitope prediction is peptide–MHC binding prediction. A number of tools and methodologies for predicting T-cell epitopes have been developed and are freely available online. We hereby provide a categorized review of

the August 2021). on the methods they use for prediction. The methods used are structure-based (SB), motif matrix (MM), sequence motif (SM), quantitative affinity matrix (QAM), artificial neural network (ANN), support vector machine (SVM), the 4. WHO Director-General's opening remarks at the 8th meeting of the IHR Emergency Committee on COVID-19—14 July 2021. Available online: <https://www.who.int/director-general/speeches/detail/who-director-general-s-opening-remarks-at-the-8th-meeting-of-the-ihr-emergency-committee-on-covid-19-14-july-2021> (accessed on 10 December 2021). mentioned the URL and which class of MHC binding prediction is supported (class I or II or both). These tools only assess a peptide's binding capability. It is still difficult for these methods to estimate deterministically whether a given peptide is an epitope or not. CTLpred [41], one of the servers, works in this category; however, it is limited to peptides with a length of up to 9mers only. However, the benefit of using ML algorithms for epitope prediction for the methods illustrated is that they address two distinct problems: the differentiation of MHC binders from non-binders and the prediction of the binding affinity of a peptide to MHC molecules. The first issue has been addressed by using classifiers such as ANNs, SVMs, decision trees (DT), 6. Canaway, E. Delta Coronavirus Variant: Scientists Brace for Impact. *Nature* 2021, 595, 175–18. and Hidden Markov models (HMMs). All of these classifiers have been trained on data containing peptides that have or do not have binding affinity to the MHC molecule. ML classifiers were developed on a dataset of peptides with an affinity to the MHC 7. Li, B.; Deng, A.; Li, K.; Hu, Y.; Li, Z.; Xiong, Q.; Liu, Z.; Guo, Q.; Zou, L.; Zhang, H.; et al. Viral infection and transmission in a large, well-traced outbreak caused by the SARS-CoV-2 Delta variant. *MedRxiv* 2021. ngov/index.html (accessed on 7 August 2021). to predict T-cell epitopes, 8. Canaway, E. Delta Coronavirus Variant: Scientists Brace for Impact. *Nature* 2021, 595, 175–18. by training ANNs on data containing MHC residues [43]. Furthermore, it has been established that combining different approaches and providing a 9. Guruprasad, L. Human SARS CoV-2 spike protein mutations. *Proteins Struct. Funct. Bioinform.* 2021, 89, 569–576. consensus prediction improves peptide–MHC prediction [44].

### 3.2. Tools for B-Cell Epitope Prediction

12. Abbas, A.K.; Lichtman, A.H.; Pillai, S. *Basic Immunology: Functions and Disorders of the Immune System*; Elsevier Slanders Publishing: Amsterdam, The Netherlands, 2015; ISBN 9780323401192. Allergy Asthma Clin. Immunol. 2018, 14, 49. The goal of predicting BCEs is to make it easier to identify a BCE for antigen replacement in an antibody production process. 14. Abbas, A.K.; Lichtman, A.H.; Pillai, S. *Basic Immunology: Functions and Disorders of the Immune System*; Elsevier Slanders Publishing: Philadelphia, PA, USA, 2007; p. 566. BCEs are classified into two types: conformational and linear. As shown in Figure 2, linear BCEs are composed of consecutive peptides and residues. Conformational ones, on the other hand, are formed of patches of solvent-exposed atoms 12. Abbas, A.K.; Lichtman, A.H.; Pillai, S. *Basic Immunology: Functions and Disorders of the Immune System*; Elsevier Slanders Publishing: Philadelphia, PA, USA, 2007; p. 566. from non-sequential residues. In addition, discontinuous and linear BCEs are also known as discontinuous and continuous BCEs. 12. Abbas, A.K.; Lichtman, A.H.; Pillai, S. *Basic Immunology: Functions and Disorders of the Immune System*; Elsevier Slanders Publishing: Philadelphia, PA, USA, 2007; p. 566. Apphen aan den Rijn, The Netherlands, 2012; ISBN 9781451109375.

13. Abbas, A.K.; Lichtman, A.H.; Pillai, S. *Basic Immunology: Functions and Disorders of the Immune System*; Elsevier Slanders Publishing: Amsterdam, The Netherlands, 2015; ISBN 9780323401192.

14. Barlow, D.J.; Edwards, M.S.; Thornton, J. Continuous and discontinuous protein antigenic determinants. *Nature* 1986, 322, 747–748.

15. BioRender Template. Available online: <https://app.biorender.com/biorender-templates> (accessed on 26 September 2021).

16. Mix, E.; Goertsches, R.; Zettl, U.K. Immunoglobulins—Basic considerations. *J. Neurol.* 2006, 253 (Suppl.

5), V9–V17, Erratum in *J. Neurol.* 2008, 255, 308. Regarding Linear BCEs, although being in the minority, their prediction has received more attention. A few existing 17. A Compact Vocabulary of Antigen–Antibody Interactions Enables Predictability of Antibody Antigen Binding. *Elife* 2018, 7, e34608. Available online: <https://doi.org/10.7554/eLife.34608> (accessed on 4 September 2021). The tool is based on a machine learning model that uses the accessibility, hydrophobicity, and flexibility properties of amino acids to predict the binding of an antibody to an antigen. It also takes into account the assessment of  $\beta$ -turns. However, it is not clear how the model was trained and whether it has been shown that the amino acid propensity scale is suitable for predicting epitopes. <https://doi.org/10.7554/eLife.34608> (accessed on 4 September 2021).

18. Ravetch, J.V.; Bolland, S. IgG Fc Receptors. *Annu. Rev. Immunol.* 2001, 19, 275–290.

The unreliability issue in predicting BCEs due to amino acid scales has been mitigated using ML algorithms. To differentiate 19. Canaway, E. Delta Coronavirus Variant: Scientists Brace for Impact. *Nature* 2021, 595, 175–18. ML methods based on heads and tails of the protein sequence (Canaway, E. Delta Coronavirus Variant: Scientists Brace for Impact. *Nature* 2021, 595, 175–18. and BepiPred 2.0 [45]) have been shown to be more reliable than methods based on the whole protein sequence (Canaway, E. Delta Coronavirus Variant: Scientists Brace for Impact. *Nature* 2021, 595, 175–18. and BepiPred 2.0 [45]).

20. Al-Qaraghul, M.M.; Kubiak-Ossowska, K.; Ferro, V.A.; Mulheran, P.A. Antibody–protein binding and conformational changes: Identifying allosteric signalling pathways to engineer a better effector response. *Sci. Rep.* 2020, 10, 13696.

21. Introduction to Antigen–Antibody Reactions. Available online: <https://microbenotes.com/introduction-to-antigen-antibody-reactions/> (accessed on 4 September 2021).

## 4. Predicting SARS-CoV-2 Epitopes

22. An Introduction to Antibodies: Antibody–Antigen Interaction. Available online:

<https://www.biorxiv.org/content/10.1101/2021.09.14.458114v1.full.pdf> (accessed on 4 September 2021). and colleagues, it has been found that, among 26 viral proteins

of SARS-CoV-2, a few proteins on its surface, such as the spike protein (S), are more variable, while others are more 23. Roper, R.L.; Rehm, K.E. SARS vaccines: Where are we? *Expert Rev. Vaccines* 2009, 8, 887–898.

conserved and internal, such as the nucleocapsid protein (N). It has been found that the spike protein (S) is responsible for 24. Shang, W.; Yao, Y.; Rao, Y.; Rao, X. The outbreak of SARS-CoV-2 pneumonia calls for viral vaccines. *NPJ Vaccines* 2020, 5, 18.



28. Liu, T.; Shi, K.; Li, W. Deep learning methods improve linear B-cell epitope prediction. *BioData Min.* 2020, 13, 10. [CrossRef]

Sr. No.	Method Name	Usage
01	NetMHC <sup>[63]</sup>	To predict HLA I class or CD8+ SARS-CoV-2 T-cell epitopes
02	NetMHCpan <sup>[64]</sup>	
03	NetCTLpan_1.1 <sup>[65]</sup>	
04	NetMHC_4.0 <sup>[66]</sup>	
05	HLAthena <sup>[67]</sup>	
06	MHCflurry <sup>[68]</sup>	To predict HLA II class or CD4+ SARS-CoV-2 T-cell epitopes
07	NetMHCII_2.3 <sup>[69]</sup>	
08	NetMHCIIpan_3.0 <sup>[70]</sup>	
09	NetMHCIIpan_4.0 <sup>[71]</sup>	
10	NeonMHC2 <sup>[72]</sup>	
11	MARIA <sup>[73]</sup>	

A few techniques listed in Table 2 have "pan" as a suffix, which indicates an ability to "predict the binding of HLA peptides for a wide range of HLA alleles".

34. A few studies have collected various HLA alleles, and a compilation of an ensemble approach for predicting the binding of peptide antigens to HLA class I and class II molecules. [\[72\]](#) A few studies have

35. **Conformational B-cell Epitope Prediction: ENET-CBioMed** [\[74\]](#) and **NetMHC2.1** [\[75\]](#), which are extra- and intracellular variables responsible for the presentation of HLA antigens were integrated to improve the prediction accuracy of the binding of peptide

36. Ren, J.; Song, J.; Ellis, J.; Li, J. Staged heterogeneity learning to identify conformational B-cell epitopes from HLA. The methods NetC1.2 [\[76\]](#) and NetChop [\[74\]](#) have also been utilized in a few studies, where extra- and intracellular variables have been integrated, which are responsible for presenting HLA antigens. It is essential to mention here that almost

37. all known B-cell epitope prediction systems use a known sequence-driven algorithm for predicting conformational B-cell epitopes from the amino acid sequences (BMC Bioinform. 2013, 18, 18).

38. The spike protein of the original virus bind to the ACE2 receptor on human cells. It has been reported in [\[77\]](#) that the D614G mutation alters the genetic code of the spike protein of SARS-CoV-2, where a change in a single amino acid takes place, and most of the COVID-19 vaccines are based on this spike protein. Due

39. to this mutation, the virus spreads faster and the spikes become more stable than those in the original virus. As a result, more functional spikes are available to bind to the ACE2 receptors, making the virus more infectious. Crooke et al. [\[2015\]](#) developed a computational model using various open-source algorithms and web-based tools to analyze the SARS-CoV-2 proteome so as

40. to identify potential B-cell epitopes. They used a peptide library as a potential vaccine target for predicting B-cell epitopes. Selection criteria to filter out the B-cell epitopes, the study described B1, T-cell epitopes (5 HLA class I, 36 HLA class II) and six B-cell

41. epitopes that have the potential to serve as primary targets for epitope-based peptide vaccine development against SARS-CoV-2.

42. Zhang, J.; Zhao, X.; Sun, P.; Gao, B.; Ma, Z. Conformational B-Cell Epitopes Prediction from Sequences Using Cost-Sensitive Ensemble Classifiers and Spatial Clustering. BioMed Res. Int. 2014, 24, 689219.

By efficiency and proteasol cleavage predictions. *Immunoform* 2009, **35**, 279-301. Several studies have been performed to predict BCEs, and TCEs. Sequence variability **Table 1**. The methods used to predict SARS-CoV-2 epitopes are listed in **Table 2**, again, these predict only the peptide-binding capacity. This is a limitation with these methods, instead of predicting the binding capability of a peptide, predicting epitopes deterministically is desired. Because viruses continue to mutate, as with SARS-CoV-2, existing algorithms may prove to be somewhat less effective against new variants. Either the vaccine's composition has to be changed or a new vaccine needs to be developed to protect against these variants [81]. Time being the critical factor, BCGs can be a great option. BCGs are, however, not a silver bullet. BCGs are a systemic TLR agonist and can be used as adjuvants in combination with peptide immunoprecipitation and even a universal adjuvant approach. Pavesio et al. [82] proposed the future research conditions for epitope prediction as predicting epitopes is a sensitive task and needs due attention in order to improve it.

45. Alex, A.J.R. Predictive estimation of protein linear epitopes by using the program PEP-2. *Appl. Bioinform.* 1999; 10: 311–314.
46. Pellequer, J.-L.; Westhof, E. PREDITOP: A program for antigenicity prediction. *J. Mol. Graph.* 1993; 11, 204–210.
47. Blythe, M.J.; Flower, D.R. Benchmarking B cell epitope prediction: Underperformance of existing methods. *Proteins Sci.* 2005; 14, 1246–1248.
48. Saha, S.; Raghava, G.P.S. Prediction of continuous B-cell epitopes in an antigen using recurrent neural network. *Proteins Struct. Funct. Bioinform.* 2006; 65, 40–48.
49. El-Manzalawy, Y.; Dobbs, D.; Honavar, V. Predicting linear B-cell epitopes using string kernels. *J. Mol. Recognit.* 2008; 21, 243–255.
50. Singh, Y.; Ansan, H.R.; Raghava, G.P.S. Improved Method for Linear B-Cell Epitope Prediction Using Antigen's Primary Sequence. *PLOS ONE* 2013; 8, e62266.
51. Yao, B.; Zhang, L.; Liang, S.; Zhang, C. SVMTrip: A Method to Predict Antigenic Epitopes Using Support Vector Machine to Integrate Tri-Peptide Similarity and Propensity. *PLoS ONE* 2012; 7, e45152.
52. Jager, J.; van der, P.; Buisson, B.; Marcatelli, P. Biophysical 2D modelling of sequence-based B-cell epitopes prediction using combinatorial epitopes models. *ACS Nano* 2017; 11, 2434–2442.
53. Greenbaum, J.A.; Andersen, P.H.; Blythe, M.; Bui, H.-H.; Cachau, R.E.; Crowe, J.; Davies, M.; Kolaskar, A.S.; Lund, O.; Morrison, S.; et al. Towards a consensus on datasets and evaluation metrics for developing B-cell epitope prediction tools. *J. Mol. Recognit.* 2007; 20, 75–82.
54. Leveson, M. Nature of the protein universe. *Proc Natl Acad Sci USA* 2009; 106, 11079–11084.
55. Su, S.; Wong, G.; Shi, W.; Liu, J.; Lai, A.C.K.; Zhou, J.; Liu, W.; Bi, Y.; Gao, G.F. Epidemiology, Genetic Recombination, and Pathogenesis of Coronaviruses. *Trends Microbiol.* 2016; 24, 490–502.
56. Hughes, K.E.; Goals, E.J.; Swaminathan, S.; Chaitz, L.; de la, D.S.; Sze, C.; Skuse, H.; Benker, R.; Bui, T.; Marks, R.; et al. Sensing and Targeting of SARS-CoV-2 by Multiple Base Learners (hologram) and one-pass learning with self-diverse seasonal prior virus similarity. *Immunity* 2021; 54, 1055–1065.
57. Zhang, X.; Fan, Y.; Ling, Y.; Lu, G.; Liu, F.; Yi, Z.; Jia, X.; Wu, M.; Shi, B.; Xu, S.; et al. Viral and host factors related to the clinical outcome of COVID-19. *Nature* 2020; 583, 437–440.
58. Schmidt, M.F.; Varga, S.M. The CD8 T Cell Response to Respiratory Virus Infections. *Front. Immunol.* 2018; 9, 678.
59. Yang, C.; Liu, C.; Fan, A.; Fan, A.; Fan, A.; Fan, A.; Fan, A.; Fan, A.; Fan, A.; Fan, A.; et al. A comprehensive framework for predicting T-cell epitopes using ensemble learning. *PLoS ONE* 2016; 11, e0158111.
60. Channanavar, P.; Berlan, S. Pathogenic human coronavirus infections: Causes and consequences of cytokine storm and immunopathology. *Semin. Immunopathol.* 2017; 39, 529–539.
61. Huber, S.E.; Beek, J.E.; de Jonge, J.; Eijssjes, W.; Baane, D.E. T Cell Responses to Viral Infections—Opportunities for Peptide Vaccination. *Front. Immunol.* 2014; 5, 171.
62. Seder, R.A.; Darrah, P.A.; Finkelstein, M. T-cell quality in memory and protection: Implications for vaccine design. *Nat. Rev. Immunol.* 2008; 8, 247–258.
63. Zhang, G.L.; DeLuca, D.S.; Keskin, D.B.; Chitkushev, L.; Zlateva, T.; Lund, O.; Reinherz, E.L.; Brusic, V. MULTIPRED2: A computational system for large-scale identification of peptides predicted to bind to HLA supertypes and alleles. *J. Immunol. Methods* 2011; 374, 53–61.
64. Nielsen, M.; Lundegaard, C.; Worning, B.; Lamberth, K.; Buus, S.; Brunak, S.; Lund, O. Reliable prediction of T-cell epitopes using neural networks with novel sequence representations. *Protein Sci.* 2003; 12, 1007–1017.
65. Ström, M.; Lundegaard, C.; Nielsen, M. NetCTLpan: Pan-specific MHC class I pathway epitope predictions. *Immunogenetics* 2010; 62, 357–368.
66. Paul, S.; Croft, N.P.; Purcell, A.W.; Tschärke, D.C.; Sette, A.; Nielsen, M.; Peters, B. Benchmarking predictions of MHC class I restricted T cell epitopes in a comprehensively studied model system. *PLoS Comput. Biol.* 2020; 16, e1007757.





86. Gagniuc, P.A.; Ionescu-Tirgoviste, C.; Gagniuc, E.; Militaru, M.; Nwabudike, L.C.; PavaloIU, B.I.; VasilăJeanu, A.; Goga, N.; Drăgoi, G.; Popescu, I.; et al. Spectral forecast: A general purpose prediction model as an alternative to classical neural networks. *Chaos Interdiscip. J. Nonlinear Sci.* 2020, 30, 033119.
  87. Bukhari, S.N.H.; Jain, A.; Haq, E.; Khder, M.A.; Neware, R.; Bhola, J.; Najafi, M.L. Machine Learning-Based Ensemble Model for Zika Virus T-Cell Epitope Prediction. *J. Health Eng.* 2021, 2021, 9591670.
  88. Huang, F.; Xie, G.; Xiao, R. Research on Ensemble Learning. In *Proceedings of the 2009 International Conference on Artificial Intelligence and Computational Intelligence*, Shanghai, China, 7–8 November 2009; Volume 3, pp. 249–252.
  89. A Gentle Introduction to Ensemble Learning Algorithms. Available online: <https://machinelearningmastery.com/tour-of-ensemble-learning-algorithms> (accessed on 8 September 2021).
  90. Dong, X.; Yu, Z.; Cao, W.; Shi, Y.; Ma, Q. A survey on ensemble learning. *Front. Comput. Sci.* 2020, 14, 241–258.
  91. Why Use Ensemble Learning? Available online: <https://machinelearningmastery.com/why-use-ensemble-learning/> (accessed on 10 July 2021).
  92. Osorio, D.; Rondón-Villarreal, P.; Torres, R.T.R. Peptides: A Package for Data Mining of Antimicrobial Peptides. *Small* 2015, 12, 44–444.
  93. Hofmann, H.; Hare, E.; GGobi Foundation. Peptider: Evaluation of Diversity in Nucleotide Libraries. R Package Version 0.2.2. 2015. Available online: <https://CRAN.R-project.org/package=peptider> (accessed on 10 September 2021).
  94. Jain, P.; Chawla, P. A Novel Smart Healthcare System Design for Internet of Health Things. In *Proceedings of the 2021 International Conference on Innovative Computing, Intelligent Communication and Smart Electrical Systems (ICSES)*, Chennai, India, 24–25 September 2021; pp. 1–8.
  95. Bukhari, S.N.H.; Jain, A.; Haq, E.; Mehbodniya, A.; Webber, J. Ensemble Machine Learning Model to Predict SARS-CoV-2 T-Cell Epitopes as Potential Vaccine Targets. *Diagnostics* 2021, 11, 1990.
- 

Retrieved from <https://encyclopedia.pub/entry/history/show/46425>