

SOFTWARE

Open Access



# USMPep: universal sequence models for major histocompatibility complex binding affinity prediction

Johanna Vielhaben, Markus Wenzel, Wojciech Samek and Nils Strodthoff\* 

\*Correspondence:  
nils.strodthoff@hi.fraunhofer.de  
Fraunhofer Heinrich Hertz Institute,  
Einsteinufer 37, 10587 Berlin,  
Germany

## Abstract

**Background:** Immunotherapy is a promising route towards personalized cancer treatment. A key algorithmic challenge in this process is to decide if a given peptide (neoepitope) binds with the major histocompatibility complex (MHC). This is an active area of research and there are many MHC binding prediction algorithms that can predict the MHC binding affinity for a given peptide to a high degree of accuracy. However, most of the state-of-the-art approaches make use of complicated training and model selection procedures, are restricted to peptides of a certain length and/or rely on heuristics.

**Results:** We put forward *USMPep*, a simple recurrent neural network that reaches state-of-the-art approaches on MHC class I binding prediction with a single, generic architecture and even a single set of hyperparameters both on IEDB benchmark datasets and on the very recent HPV dataset. Moreover, the algorithm is competitive for a single model trained from scratch, while ensembling multiple regressors and language model pretraining can still slightly improve the performance. The direct application of the approach to MHC class II binding prediction shows a solid performance despite of limited training data.

**Conclusions:** We demonstrate that competitive performance in MHC binding affinity prediction can be reached with a standard architecture and training procedure without relying on any heuristics.

**Keywords:** Major histocompatibility complex, Binding affinity prediction, Peptide data, Recurrent neural networks, Language modeling

## Background

Immunotherapy is a promising route towards personalized cancer treatment with a variety of possible realizations, see [1–4] for recent reviews. One path is the administration of nanoparticle vaccines customized with neoantigens. The major histocompatibility complex plays a central role in this process as it is supposed to bind to peptides derived from proteins of the cell or from pathogens in order to display them on the surface of



© The Author(s). 2020 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

the cell for recognition by T-cells. There are three classes of MHC molecules, where MHC class I and II are most important due to their involvement in the targeted immune response. Due to the special nature of the MHC protein, it can bind to peptides that are potentially structurally very different from each other. Therefore, the prediction if a MHC molecule binds to certain peptide is a very challenging task that is, however, a crucial sub task for neoantigen identification for practical realizations of personalized immunotherapy [1].

MHC binding prediction is a well-established problem in bioinformatics with a large number of existing algorithmic solutions. Although many of these algorithms show an excellent performance, they typically rely on complicated training procedures to achieve this performance, such as pretraining on prediction tasks for related alleles or training with artificial negative peptides. In addition, existing solutions use complicated model selection procedures to select a small number of well-performing models from potentially hundreds of trained models to eventually construct an ensemble classifier. Most of the existing approaches are restricted to peptides of fixed length, where shorter sequences are padded or longer sequences are trimmed to an appropriate length by well-motivated but still heuristic rules to identify so-called binding regions. The most prominent MHC I prediction algorithms are summarized in Table 1. We refer to dedicated reviews for more detailed comparisons [5, 6].

**Table 1** Comparison of MHC I prediction tools

<i>Architecture</i>	
SMPMBEC [7]	One-hot encoding, linear model (scoring matrix)
consensus [8]	Linear model (scoring matrix), median rank as prediction
NetMHC4 [9]	Input: 9mer fixed length blocks substitution matrix (BLOSUM) encoding plus additional features; multilayer perceptron with one hidden layer
NetMHCpan4 [10]	Input: 9mer fixed length BLOSUM encoding for peptide, pseudo-sequence for MHC molecule plus additional features; multilayer perceptron with one hidden layer
MHCFlurry [11]	Input: 15mer fixed length BLOSUM62 encoding, missing residues filled with wildcard amino acid (AA); feedforward neural network (NN) with 0 to 2 locally connected and one fully connected hidden layer
<b>USMPep</b> (this work)	Learned embedding layer; AWD LSTM with one hidden layer
<i>Training procedure</i>	
SMPMBEC	Ridge regression with modified regularization, peptide MHC binding energy covariance (PMBEC) similarity matrix as Bayesian prior
consensus	Four scoring matrices from existing algorithms
NetMHC4	Training on non 9mer peptides by insertion of wildcard AA or deletion at all possible positions; augmented training set with natural peptides for each length assumed to be negative
NetMHCpan4	Same insertion/ deletion procedure as NetMHC4; augmented training set with random artificial negatives
MHCFlurry	Pretraining on BLOSUM62 similar allele for alleles with little training data; augmented training set with artificial negative peptides
<b>USMPep</b>	Optional: language model pretraining on unlabeled sequences
<i>Model selection</i>	
SMPMBEC	Single model
consensus	Single model
NetMHC4	Ensemble of 4 NNs
NetMHCpan4	Ensemble of 100 NNs
MHCFlurry	Ensemble of 8-16 NNs selected from 320 models on a validation set
<b>USMPep</b>	Optional: ensemble of 10 NNs with identical architectures and hyperparameters

Finally, not all binding prediction tools are evaluated on standard benchmark datasets, which reduces the comparability, and, even where this is the case, it is often hard to disentangle algorithmic advancements from improvements due to larger amounts of training data. In addition, statements about the generalization in the sense of the algorithm's performance when applied to unseen data are often difficult due to potential overlaps between train and test sets, in particular as training sets often remain undisclosed. This urges for the creation of benchmark repositories, where the existing data are processed in a standardized fashion and split into training, validation and test sets.

In this manuscript, we demonstrate that state-of-the-art performance can be reached with a straightforward approach: We use a single-layer recurrent neural network that is *trained end-to-end* on a regression task *without any task-specific prior knowledge* such as fixed embeddings in the form of amino acid similarity matrices. By construction, this model is able to incorporate *input of variable length without the need for heuristics*, such as for the identification of binding regions. The model is trained using *standard training procedures* without any artificial data or pretraining on related classification tasks. Even *single models are very competitive*. Ensembling or language model pretraining only slightly improve this performance. We fix hyperparameters only once and use standard benchmark datasets to assess the model performance. We provide, amongst others, evaluation results on the recently published HPV dataset [12], demonstrating an excellent performance, which strongly suggests that the measured model performance generalizes to unseen peptide data.

Recurrent architectures have already been used previously for MHC binding prediction [13, 14] and we discuss in more detail how *USMPep* stands out from these approaches. MHCnuggets [13] is rather similar to the proposed approach (apart from the use of fixed embeddings), but relies on a complex transfer learning protocol to achieve its performance. Only limited benchmarking results are available, which makes it difficult to realistically assess its prediction performance. The very recent MHCSeqNet [14] also uses a recurrent architecture, again with pretrained rather than learned embeddings, incorporating both peptide and allele sequence to train a single prediction model for all alleles. However, the paper frames the prediction task as a classification task, which makes it difficult to align the results with the large number of existing benchmark datasets that are predominantly targeted at regression tasks. Nevertheless, the inclusion of the allele sequence represents an exciting opportunity for MHC binding affinity prediction in particular in the light of recent advances in natural language processing on tasks involving two input sequences such as question answering tasks.

## Implementation

### USMPep: universal sequence models for peptide binding prediction

The approach builds on the *UDSMProt*-framework [15] and related work in natural language processing [16]. We distinguish two variants of our approach, either train the regression from scratch or employ language model pretraining. A language model tries to predict the next token given the sequence up to this token, on unlabeled sequence data, here: of simulated proteasome-cleaved peptides. The architecture of the language model is, at its core, a recurrent neural network regularized by different kinds of dropout,

and more specifically an averaged stochastic gradient descent weight-dropped long short-term memory (AWD LSTM) model [17]. After the language model pretraining step, the model is finetuned on the regression task of MHC binding prediction by replacing the output layer with a concat pooling layer and two fully connected layers, see Fig. 1 for a schematic representation. The setup closely follows that used in [15], where protein properties were predicted. The smaller dataset sizes and shorter sequence lengths in the peptide setting (in comparison to protein classification) do not allow for building up large contexts and were accounted for by the reduction of the number of layers from 3 to 1, of the number of hidden units from 1150 to 64 and of the embedding size from 400 to 50.

Similar to [15], the training procedure included 1-cycle learning rate scheduling [18] and discriminative learning rates [16] during finetuning. Target variables for the regression model were log-transformed half-maximal inhibitory concentration ( $IC_{50}$ )-values and a modified mean-squared error loss function [11] that allows to incorporate qualitative data.

Dropout rate, the number of training epochs, hidden layers, hidden units and embedding dimensions, were set based on selected alleles of a particular MHC class I dataset (*Kim14* [19], see the detailed description below) by using the score on one of the provided cross-validation folds. The learning rate was determined based on range tests [18]. After this step, the aforementioned hyperparameters were kept fixed for all datasets and alleles both for MHC class I and class II prediction. In particular, neither hyperparameters nor models were selected based on test set scores.

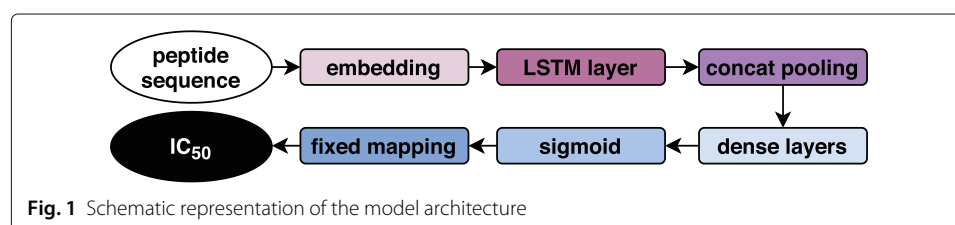
For later convenience, the following acronyms refer to the prediction tools introduced in this work:

- **USMPep\_FS\_sng** single prediction model trained from scratch
- **USMPep\_FS\_ens** ensemble of ten prediction models trained from scratch
- **USMPep\_LM\_sng** single prediction model with language model pretraining
- **USMPep\_LM\_ens** ensemble of ten prediction models with language model pretraining

For simplicity, we consider ensembles of models with identical architectures and hyperparameters and average the final individual predictions.

### MHC binding prediction datasets

For the downstream task of peptide MHC binding prediction, we benchmarked our model on three MHC class I and one MHC class II binding affinity datasets (details listed in Table 2). These datasets comprise peptide sequences and the corresponding binding affinities to specific MHC alleles.



**Table 2** Details of training and test datasets

Dataset	Usage	Total size	Share of binders	# alleles	Median size	Share of quant. meas.	Sequence length
MHC class I							
<i>BD2009</i>	train	117326	0.25	53	1971	0.58	8–11
<i>Blind</i>	test	27680	0.33	53	470	0.58	8–11
<i>MHCFlurry18</i>	train	120720	0.25	32	3659	0.68	8–15
<i>IEDB16_I</i>	test	2827	0.54	32	73	1.0	9
<i>MHCFlurry18</i>	train	68117	0.26	7	6884	0.64	8–15
<i>HPV</i>	test	743	0.34	7	125	0.37	8–11
MHC class II							
<i>Wang10</i>	train	23203	0.37	24	999	1.0	15–37
<i>IEDB16_II</i>	test	15691	0.33	24	641	1.0	15

The threshold for MHC class I binders is 500nM, except for the HPV dataset, where the threshold is 100 000nM. For MHC class II binders, the threshold is 1000nM

*Kim14* is a commonly used binding affinity dataset compiled by [19], available on the Immune Epitope Database (IEDB)<sup>1</sup> [20], and is split into a non-overlapping training (*BD2009*) and test set (*Blind*). Similar peptides (of same length with at least 80% sequence identity) shared by training and test set were removed from *Blind*. For *BD2009*, we selected the provided cross-validation split without similar peptides between the subsamples ("*cv\_gs*"). There are 53 class I alleles (human and mouse/ macaque alleles) with respectively 117326 and 27680 affinity measurements in *BD2009* and *Blind*. For comparability with recently developed systematical benchmarks [5, 12], we tested *USMPep* on two further MHC I datasets, which we refer to as *HPV* and *IEDB\_16*. The training data of the tools reported in the literature vary in size and compilation.

We trained our models on data provided by [11] and refer to this dataset as *MHCFlurry18*. It is assembled from an IEDB snapshot of December 2017 and the *Kim14* dataset.

*HPV* is a recently published dataset of 743 affinity measurements of peptides derived from two human papillomavirus 16 (HPV16) proteins binding to seven human leukocyte antigen (HLA) class I alleles [12]. Peptides were considered as binders if they had  $IC_{50}$ -values below 100000nM. For peptides classified as non-binders, quantitative measurements are not available.

*IEDB16\_I* is made up of an IEDB snapshot of October 2016 [5]. It was filtered for quantitative measurements with  $IC_{50} \leq 50000$ nM and 9mer peptides. Training sequences of other tools were removed from the dataset. It consists of 2827 affinity measurements across 32 class I alleles. We removed any sequences occurring in the test dataset from our training data *MHCFlurry18*.

In addition, we trained and tested *USMPep* on MHC class II binding data: *Wang10* is an experimental binding affinity dataset from the IEDB site based on the dataset by [21]. We used it to train our prediction tools.

*IEDB16\_II* is a MHC II test dataset provided by [5] from the same IEDB snapshot as the MHC I *IEDB16\_I* test set above, filtered for quantitative measurements with  $IC_{50} \leq 50000$ nM and 15mer peptides. After removing sequences present in the training data, 15034 affinity measurements covering 24 alleles remained in the test dataset. We benchmarked our models on this dataset.

<sup>1</sup><http://tools.iedb.org/main/datasets/>

### Evaluation metrics

For performance evaluation, we consider two evaluation metrics that are most frequently considered in the literature [5, 11]: The area under the receiver operating characteristic curve (AUC ROC) measures the performance of binary classifying binders and non-binders. While AUC ROC is straightforward to evaluate, it comes with the disadvantage of having to specify a threshold to turn the targets into binary labels, which discards valuable label information during the evaluation procedure. Commonly applied threshold values exist for the datasets under consideration, as discussed in the previous section, but the simplicity of this procedure neglects a possible allele dependence of these threshold values [22]. Ranking metrics such as Spearman  $r$  evaluate the correlation between the rankings of measured and predicted affinities and circumvent this issue. Spearman  $r$  can only be evaluated for quantitative measurements, which discards information on test sets that contain also qualitative measurements. For both metrics, we calculated error bars based on 95% empirical bootstrap confidence intervals. For single models, we report the mean performance across 10 runs and the maximal deviation of the point estimate compared to the lower and upper bounds provided by the respective confidence intervals as a conservative error estimate. An alternative approach, which is taken for example in [14], is to frame the problem as a classification rather than a regression problem i.e. to predict the probability of binding of a peptide to a given allele and to use AUC ROC as performance metric. In principle, the raw prediction of our models before transforming back to  $IC_{50}$ -values could also be interpreted as probabilities and are also available from our code repository to ensure straightforward comparability, even though the results between regression and classification models are only partially comparable due to different training objectives.

The prediction performance across different alleles that make up a single MHC benchmark dataset can be quantified in different ways. *Overall* performance measures can be calculated across multiple alleles by concatenating all target and prediction results and evaluating the respective metrics on this set. This predominantly used but rarely discussed method has to be contrasted with reporting the *mean* or the median of the respective performance measures across all alleles, which is the default evaluation metric for related tasks such as remote homology detection [23] or transcription factor binding site prediction [24]. The difference between both evaluation approaches is related to the discussion about micro vs. macro averages for the evaluation of multi-class classification problems [25]. In particular, there are two fundamental differences between both evaluation approaches: First, the datasets enter the *overall* score with different weights determined by the size of the respective test sets, which is a weighting based on the experimental availability of binding affinities whereas the *mean* score assigns equal weight to all test sets. Second, the *overall* performance measure implicitly assumes that prediction scores are directly comparable across different alleles, which seems slightly questionable in the light of the discussion of allele-dependent binding thresholds [22]. To give the reader a complete picture of the prediction performance, we will report *overall* as well as *mean* scores. In any case, we advocate to provide individual prediction for all peptides, which allows to possibly redo the analysis using a different performance metric at a later point in time. To this end, the peptide-wise binding affinity predictions for our tools are provided in the accompanying code repository.

## Results

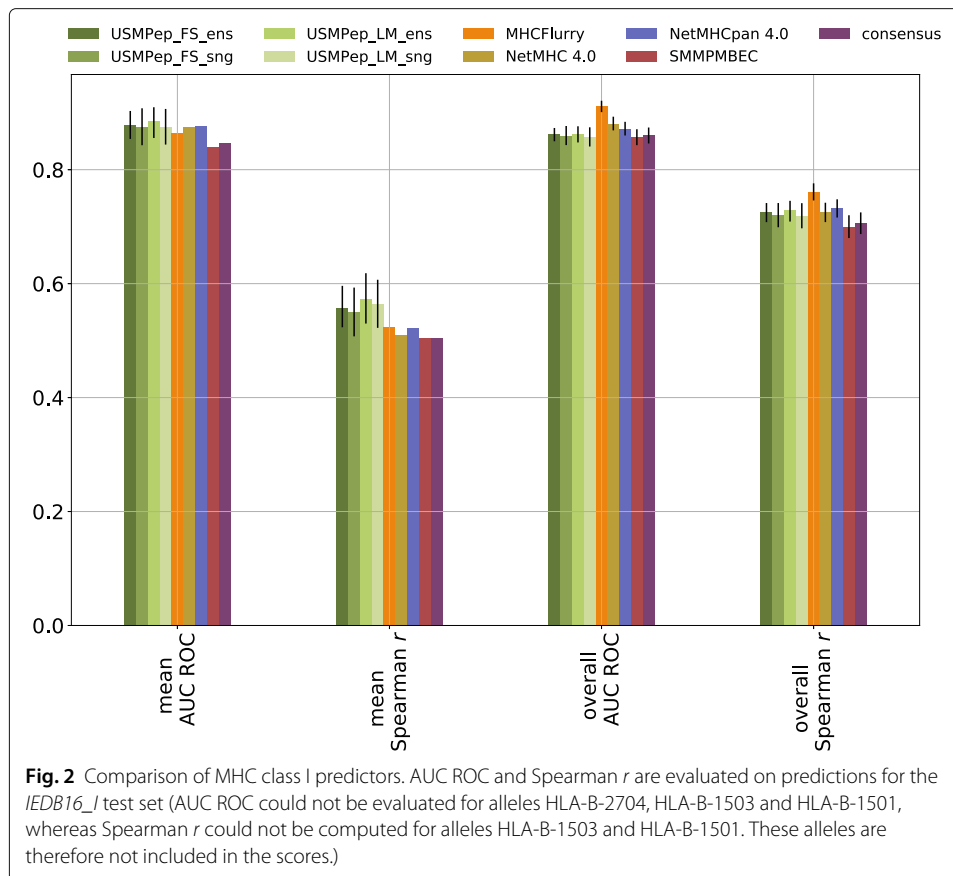
The results section is organized as follows: In “[MHC class I binding prediction](#)” section, we present a detailed evaluation of the performance of *USMPep* for MHC class I binding affinity prediction. Three different benchmark datasets highlight different performance characteristics. In “[MHC class II binding prediction](#)” section, we investigate the applicability of our methods for MHC class II binding affinity prediction. Finally, we discuss language modeling on peptide data and its impact on downstream performance in “[Language modeling on peptide data and its impact on downstream performance](#)” section.

### MHC class I binding prediction

#### *IEDB16 dataset*

We open the assessment of MHC class I binding prediction with results on the *IEDB16* dataset that showcases the excellent predictive performance of *USMPep*. We compare to literature results that were evaluated in a recent comprehensive benchmark [5] on this dataset. This benchmark includes evaluation metrics testing not only accuracy of binder classification, but also accuracy of binding affinity ranking and of direct binding affinity prediction accuracy. Covering 32 HLA alleles, the *IEDB16* dataset reflects a broad spectrum of MHC molecules.

In Fig. 2, we show *overall* AUC ROC and *overall* Spearman  $r$  as reported by [5] for the latest versions of the NetMHC tools, MHCFlurry, SMMPMBEC and consensus and our scores for the different versions of *USMPep*. This is supplemented by *mean* AUC ROC and *mean* Spearman  $r$  compared to results provided in the data repository accompanying [5]. For *mean* AUC ROC and Spearman  $r$  error bars could not be calculated for the literature approaches due to the fact that only allele-wise scores but no peptide-wise predictions were provided. In the light of the issues discussed in “[Evaluation metrics](#)” section, we advocate the use of *mean* scores rather than *overall* scores. For easy comparability, we also provide *overall* scores as they are used predominantly in the literature. It turns out that an ensemble of ten predictors with language model pretraining (*USMPep\_LM\_ens*), reaches the highest scores in both *mean* evaluation metrics. In this respect, the results of all four *USMPep*-variants are consistent with each other and similar (within error bars) to the result of MHCFlurry, the best-performing method in the benchmark [5]. This result stresses the claims of excellent prediction performance even for a single model trained from scratch. Interestingly, the performance of all proposed prediction tools is slightly worse when considering *overall* scores. The error boundaries of our tools barely touch those of MHCFlurry with regard to *overall* Spearman  $r$ . In particular, in terms of *overall* AUC ROC none of our predictors is consistent with MHCFlurry within error bars. We further investigated the origin of this performance deficiency and found that it could be traced back to a single allele, HLA-B-3801, which is peculiar in the sense that 172 of the 176 test set samples fall into a single Hobohl cluster [19] of sequences with more than 80% sequence similarity, i.e. show a particularly high sequence identity that is not seen in other test datasets. These 172 samples constitute a sizable amount of the overall 2827 test samples and strongly influence the predictive performance when using *overall* performance metrics. With the exceptions in terms of the *overall* metrics for the *IEDB16* dataset, our proposed methods are consistent with the best-performing methods for all MHC I benchmark datasets both for *overall* and *mean* performance metrics.



### HPV dataset

As the training data is not publicly available for some MHC I prediction tools, a possible overlap between training and test datasets and correspondingly an overestimation of the predictive performance cannot be excluded. The same applies to the most common procedure of reducing the overlap between training and test set by merely removing sequences from the test set that are also contained in the training set in identical form rather than using more elaborate measures for sequence similarity. These issues can be circumvented by a performance evaluation on a dataset of different origin that has so far not been used to train MHC prediction tools. This applies to the recently released HPV binding affinity data [12]. However, in this benchmark, it is not possible to disentangle superior prediction performance due to larger amounts of training data from algorithmic advances since size and compilation of the training set of the algorithms vary.

There are only quantitative measurements available for the peptides considered as binders and we therefore chose to evaluate the predictive performance with AUC ROC. We report the performance of all models considered in [12] and our tools measured by AUC ROC in Table 3, where we used the predictions provided by [12]. Our *USMPep* tools show an excellent prediction performance. For three out of seven alleles, an *USMPep*-model even reaches the highest AUC ROC. All neural-network-based predictors show a similar AUC ROC evaluated across all measurements in the dataset, while the ensemble with language model pretraining (*USMPep\_LM\_ens*) shows the highest *mean* and



**Table 3** Benchmarking MHC class I predictors on recently published binding affinity data (*HPV16*), see also Table 4 for allele-wise scores

allele	mean AUC ROC	overall AUC ROC
USMPep_FS_ens	0.824(3)	0.814(3)
USMPep_FS_sng	0.818(4)	0.808(4)
USMPep_LM_ens	<b>0.831(4)</b>	<b>0.815(3)</b>
USMPep_LM_sng	0.813(5)	0.802(4)
MHCFlurry	0.817(4)	0.809(4)
NetMHC 3.4	<b>0.831(3)</b>	0.794(3)
NetMHC 4.0	0.803(4)	0.780(3)
NetMHCpan 2.8	0.818(4)	0.792(3)
NetMHCpan 3.0	0.815(4)	0.787(3)
NetMHCpan 4.0	0.820(3)	0.792(4)
SMM	0.684(5)	0.695(4)
SMMPMBEC	0.722(6)	0.723(4)
Pickpocket 1.1	0.760(5)	0.708(4)
consensus	0.751(5)	0.766(4)
IEDB recommended	0.756(5)	0.772(4)
NetMHCcons 1.1	0.827(4)	0.799(3)

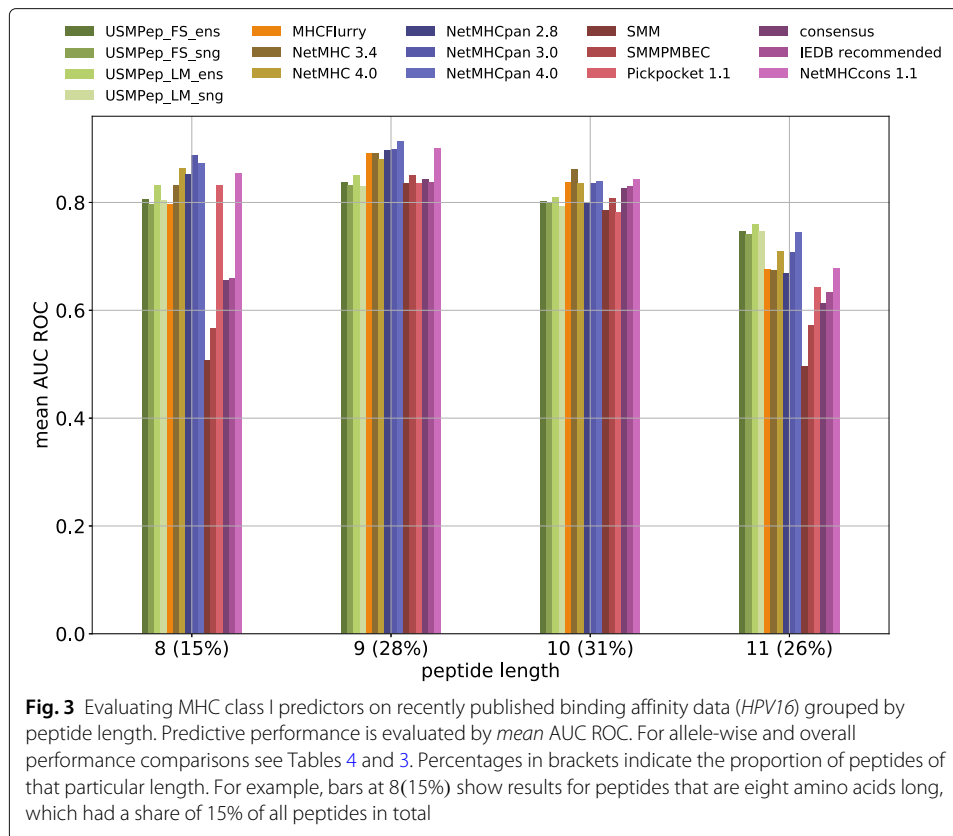
Predictive performance is evaluated by AUC ROC (threshold for binders < 100 000nM) on single alleles and across all alleles (*mean* and *overall*). The scores for literature approaches were calculated based on peptide-wise predictions provided in [12]. Numbers in brackets in the table concisely denote the corresponding bootstrap confidence intervals. For instance, 0.824(3) stands for a mean AUC ROC of 0.824 ± 0.003

*overall* scores among all prediction tools. As for the *IEDB16* dataset, even the single model *USMPep*-tools are very competitive.

It is instructive to investigate the performance of the different MHC prediction tools restricted to peptides of a certain length, which is only possible for the HPV dataset, where peptide-wise predictions for all literature approaches are provided. The result of such an analysis is shown in Fig. 3. Our tools outperform the other models on only the 11mer peptides. This observation can be explained by the fact that the internal state of the recurrent neural network has to build up over the sequence. The longer the peptide,

**Table 4** Allele-wise results on (*HPV16*), see also Table 3 for *mean* and *overall* scores

Allele	HLAA1	HLAA11	HLAA2	HLAA24	HLAA3	HLAB15	HLAB7
USMPep_FS_ens	0.793	<b>0.885</b>	0.830	0.807	0.768	0.803	0.884
USMPep_FS_sng	0.785	0.883	0.822	0.798	0.764	0.799	0.883
USMPep_LM_ens	<b>0.848</b>	0.880	0.809	<b>0.821</b>	0.766	0.824	0.871
USMPep_LM_sng	0.813	0.869	0.805	0.802	0.755	0.805	0.854
MHCFlurry	0.816	0.850	<b>0.833</b>	0.755	0.793	0.797	0.867
NetMHC 3.4	0.841	0.867	0.793	0.765	<b>0.840</b>	0.825	0.884
NetMHC 4.0	0.823	0.855	0.792	0.730	0.779	0.825	0.801
NetMHCpan 2.8	0.756	0.863	0.787	0.778	0.794	0.857	0.880
NetMHCpan 3.0	0.841	0.848	0.781	0.739	0.778	0.876	0.825
NetMHCpan 4.0	0.839	0.854	0.805	0.742	0.784	<b>0.891</b>	0.836
SMM	0.476	0.828	0.730	0.643	0.788	0.704	0.646
SMMPMBEC	0.593	0.846	0.777	0.639	0.799	0.716	0.670
Pickpocket 1.1	0.744	0.773	0.757	0.709	0.731	0.808	0.802
consensus	0.570	0.870	0.772	0.687	0.767	0.832	0.756
IEDB recommended	0.566	0.877	0.769	0.702	0.772	0.852	0.755
NetMHCcons 1.1	0.807	0.872	0.797	0.777	0.819	0.847	<b>0.889</b>

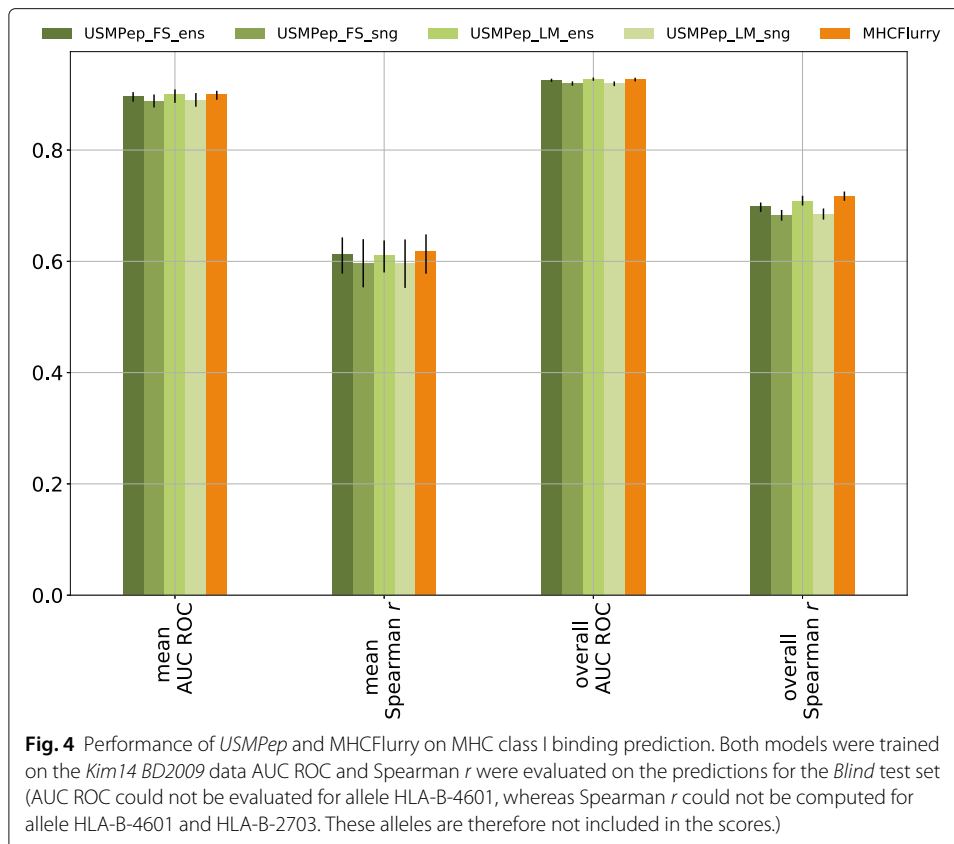


the more context is available, which is why *USMPep* generates comparably more accurate predictions for long sequences than for shorter ones.

#### Kim14 dataset

As final benchmark dataset for MHC class I prediction, we consider the *Kim14* dataset that is interesting for a number of reasons. In order to investigate how the predictive power of our approach depends on the size of the training data set, we trained and tested our model on the *Kim14 BD2009* and *Blind* data. The authors of [11] kindly provided us with the *Blind* predictions of their tool trained on *BD2009*, which allow for a direct comparison with a state-of-the-art tool. Corresponding training routines are by now also available in the code repository accompanying [11].

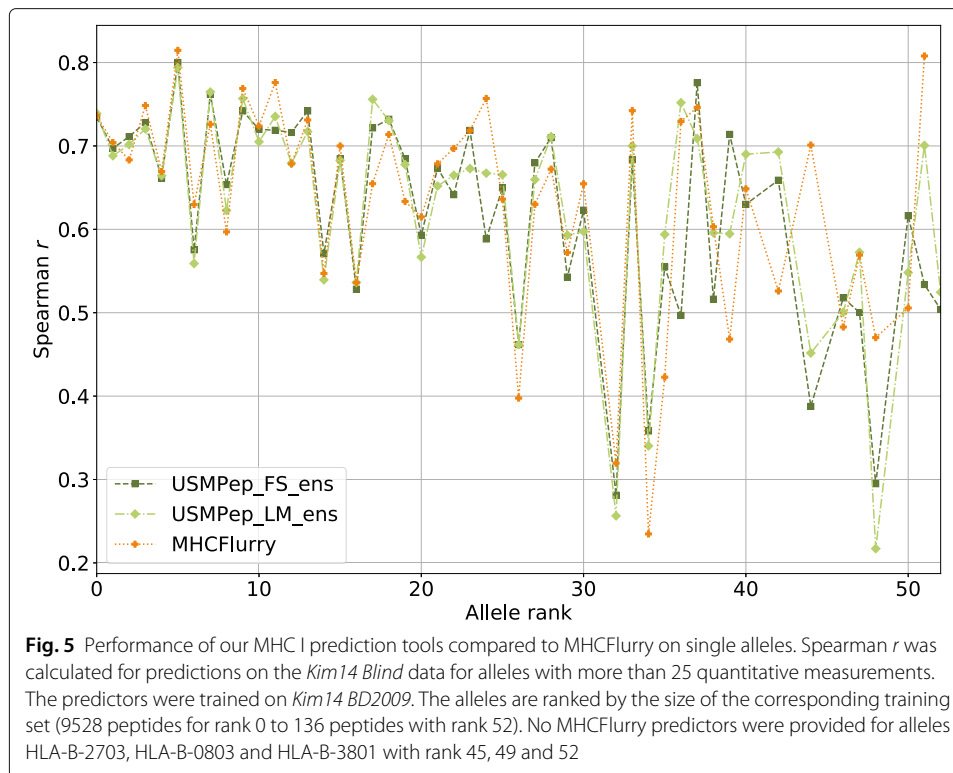
First, we compare the prediction success measured by AUC ROC and Spearman  $r$  computed across all alleles (Fig. 4). No MHCFlurry predictors exist for alleles HLA-B-2703, HLA-B-0803 and HLA-B3801 with rank 45, 49 and 52 due to insufficient training data. These alleles were therefore also excluded for the scores of our tools. The predictors perform very similarly with regard to all metrics. Due to the different scales of the metrics, the minor performance variations appear to be more pronounced for Spearman  $r$  in comparison to AUC ROC. Nevertheless, the ranking of the almost equally performing predictors remains consistent. Our pretrained tool *USMPep\_LM\_ens* performs only slightly better than *USMPep\_FS\_ens* trained from scratch. This also holds for the single model versions. Both *USMPep* ensemble predictors are compatible with MHCFlurry.



Second, to examine the impact of the training set size, we report allele-wise Spearman *r* scores in Fig. 5 for our predictors and MHCFlurry. The alleles are ranked by the size of the corresponding training set. While 9528 training sequences exist for the rank 0 MHC molecule HLA-A-0201, there are only 136 training peptides for allele HLA-B-3801 with rank 52. Spearman *r* is only shown for alleles with more than 25 quantitative measurements. The variance of the allelwise performances of the different tools becomes more pronounced the less training data are available. However, none of the models outperforms the others for the subset of alleles with less than 1000 training data points (rank range 33 to 52). When averaging over the alleles with rank 33 to 52, the mean Spearman *r* scores of 0.54(5), 0.57(4), 0.57(5) for *USMPep\_FS\_ens*, *USMPep\_LM\_ens* and MHCFlurry, respectively, remain consistent within error bars. This observation is interesting considering the fact that for alleles with fewer than 1000 training measurements, MHCFlurry was pretrained on an augmented training set with measurements from BLOSUM similar alleles, *USMPep\_LM\_ens* was pretrained on a large corpus of unlabeled peptides and *USMPep\_FS\_ens* in contrast only saw the training sequences corresponding to one MHC molecule. These results stress that further efforts might be required to truly leverage the potential of unlabeled peptide data in order to observe similar improvements as seen for proteins [15] in particular for small datasets.

#### MHC class II binding prediction

Turning to MHC Class II binding prediction, we aim to demonstrate the universality of our approach beyond its applicability to different MHC I alleles. Here, we stress again

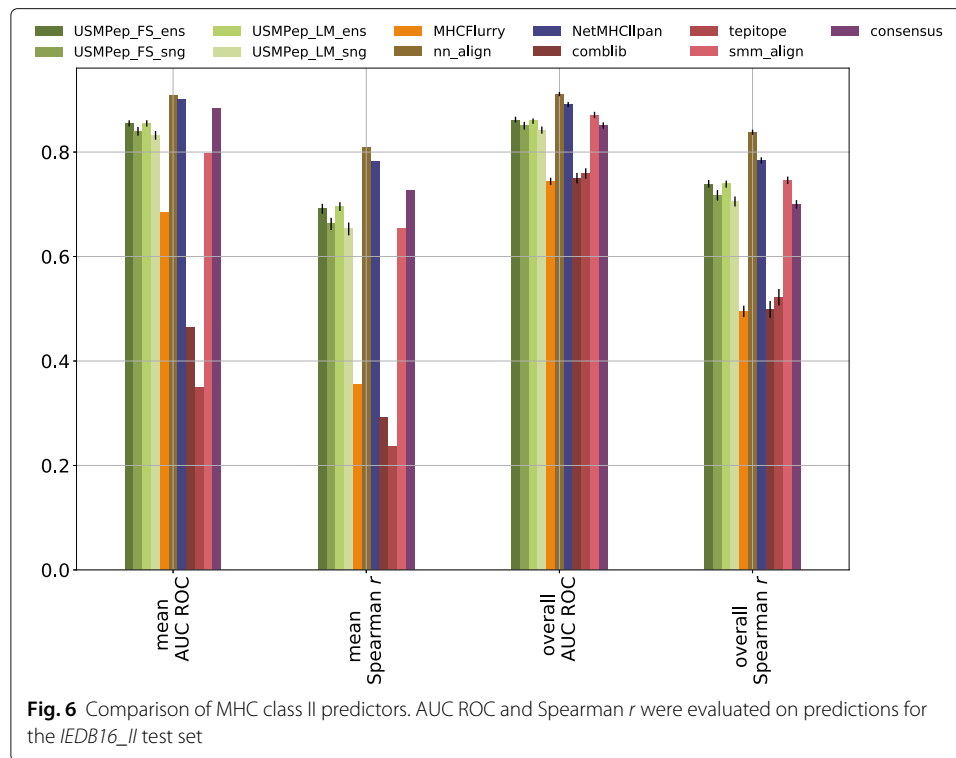


that we use the same model architecture, the same pretrained language model in case of pretraining, and even the same set of hyperparameters for all MHC class I and class II alleles. The main difference between and MHC class I and class II binding prediction is the typically larger length of 15 amino acids for MHC class II compared to at most 11 for MHC class I. The analysis of the prediction performance in dependence of the length of the peptide in the previous section suggests that this setting is particularly suitable for the *USMPep* prediction tools. Unfortunately, the reported literature results vary widely concerning the selection of training data, which makes it difficult to distinguish between algorithmic improvements and improvements due to larger amounts of training data.

The *USMPep* prediction tools, and in particular the ensemble variants, show a solid performance compared to literature results, see Fig. 6. Whereas the *USMPep*-predictors always provided the best-performing method for MHC class I prediction, it is outperformed for MHC class II by *NetMHCIIpan* and *nn\_align*. We deliberately decided to train on *Wang10* instead of a more recent IEDB snapshot to work on a well-defined published dataset. However, this makes it hard to assess if the performance differences between our results and the best-performing methods can be attributed to the fact that the *USMPep*-predictors were trained using IEDB data up to 2010 whereas in particular the best-performing tools were trained on larger amounts and more recent data or if there a particular intricacies inherent to the MHC class II prediction task.

#### Language modeling on peptide data and its impact on downstream performance

As final analysis, we analyze language modeling on peptide data and its impact on MHC binding affinity prediction as downstream task. To this end, we constructed a dataset of simulated proteasome-cleaved peptides to pretrain *USMpep* on a large corpus of



unlabeled sequences. We filtered the SwissProt release 2018\_10 for the human proteome and employed NetChop [26] to obtain proteasome cleavage sites for these proteins. The stochastic process of protein slicing was modeled by cutting with the cleavage probability provided by NetChop. We discarded sequences of less than eight and more than 20 amino acids length and obtained 6547641 peptides. We compare the performance of a peptide language model to that of a language model trained on human protein data using prediction accuracy as metric.

The results in terms of language model performance along with the corresponding downstream performance (MHC) on the regression task are compiled in Table 5 and allow a number of interesting observations: First, the language model performance increases considerably when training on (proteasome-cleaved) peptide data in accordance with expectations. It is crucial to remark, that the language modeling task on peptide data poses additional difficulties compared to language modeling on protein data as the sequences are comparably short and the model thus cannot build up a lot of context. Additionally, the model does not only have to learn the normal language model task for protein data but implicitly has to learn to stochastically predict cleavage sites. Second, even we evaluated on protein data, the protein language model only reaches an accuracy of 0.137, which is considerably lower than the accuracy of 0.41 reported in the literature [15]. This effect is a direct consequence of the considerably smaller model size (1 instead of 3 layers; 64 instead of 1150 hidden units; embedding size of 50 instead of 400).

The details of the language model pretraining directly impact the downstream performance and show a consistent trend across all experiments described above even though the differences in downstream performance stay small and mostly remain consistent within error bars. In line with the general trend, the most downstream-task-adapted

**Table 5** Language model and MHC class I binding affinity prediction performance

Model	LM		Downstream ( <i>mean</i> )	
	perpl.	acc.	AUC ROC	Spearman <i>r</i>
LM (protein)	39.3	0.083	0.90(2)	0.55(4)
LM (peptide)	13.4	0.206	0.89(2)	0.57(4)
From scratch	–	–	0.89(2)	0.55(3)

Language model metrics perplexity (perpl.) and accuracy (acc.) were in all cases evaluated on peptide data. The downstream performance corresponds to an ensemble of 10 predictors trained on the *MHCFlurry18* and evaluated on the *IEDB16\_I* test set

pretraining on peptide data performs best, generally performing slightly better than the corresponding model trained from scratch. In contrast, pretraining on protein data in general even leads to a loss in performance compared to training from scratch.

## Conclusions

In this work, we put forward *USMPep*, a recurrent neural network that consistently shows excellent performance on three popular MHC class I binding prediction datasets as well as a solid performance on MHC class II binding prediction, see Table 6 for a performance summary. Most remarkably, this is achieved with a standard training procedure without incorporating artificial negative peptides, complicated transfer learning protocols or ensembling strategies and without relying on heuristics.

A central issue that prevents a true comparability of algorithmic approaches to the problem is the fact that the datasets that were used to train the prediction models differ between different literature approaches and are often not publicly available. This entangles the predictive power of a given algorithm with the data it was trained on. This urges for the creation of an appropriate benchmarking repository along with standardized evaluation procedures to allow for a structured benchmarking of MHC binding prediction algorithms. As a first step, we advocate to provide binding affinity predictions for all peptides to allow fine-grained comparisons of the overall predictive performance even at a later stage as opposed to reporting just a single score summarizing the performance across all datasets.

## Availability and requirements

Project name: USMPep

Project home page: <https://github.com/nstrodt/USMPep>

Operating system(s): Platform independent

**Table 6** Performance summary: Rank of *USMPep* compared to competitors across the different datasets

Dataset	<i>mean</i>		<i>overall</i>	
	AUC ROC	Spearman <i>r</i>	AUC ROC	Spearman <i>r</i>
	MHC class I			
IEDB16_I	<b>1<sup>st</sup></b>	<b>1<sup>st</sup></b>	3 <sup>rd</sup>	3 <sup>rd</sup>
HPV	<b>1<sup>st</sup></b>	–	<b>1<sup>st</sup></b>	–
Kim14	<b>2<sup>nd</sup></b>	<b>2<sup>nd</sup></b>	<b>1<sup>st</sup></b>	<b>2<sup>nd</sup></b>
	MHC class II			
IEDB16_II	4 <sup>th</sup>	4 <sup>th</sup>	3 <sup>rd</sup>	4 <sup>th</sup>

Scores marked in bold face are best-performing or consistent with the best-performing result within error bars

Programming language: Python

Other requirements: see project homepage

License: BSD

Any restrictions to use by non-academics: as permitted by BSD License

#### Abbreviations

MHC: Major histocompatibility complex; AWD LSTM: Averaged stochastic gradient descent weight-dropped long short-term memory; IEDB: Immune Epitope Database; AUC ROC: Area under the receiver operating characteristic curve; BLOSUM: Blocks substitution matrix; AA: Amino acid; NN: Neural network; PMBEC: Peptide MHC binding energy covariance

#### Acknowledgements

The authors thank Patrick Wagner for discussions and work on related topics. The authors thank Timothy O'Donnell for correspondence and for kindly providing MHCFlurry predictions on the *Kim14* dataset. *USMPep* was implemented using PyTorch [27] and fast.ai [28].

#### Authors' contributions

Conceptualization: NS, MW; Data Curation: JV, NS; Investigation: JV; Methodology: JV, NS; Software: JV, NS; Supervision: NS, WS; Validation: JV; Visualization: JV; Writing — Original Draft Preparation: JV, NS; Writing — Review & Editing: JV, MW, WS, NS; Final approval of the manuscript: JV, MW, WS, NS

#### Funding

This work was supported by the Bundesministerium für Bildung und Forschung (BMBF) through the Berlin Big Data Center under Grant 01IS14013A and the Berlin Center for Machine Learning under Grant 01IS18037L.

#### Availability of data and materials

The USMPep source code together with a user guide, links to pretrained models and individual predictions of USMPep for the datasets analysed in this study are available at <https://github.com/nstrodt/USMPep>. The training and test datasets are available in the IEDB repository <http://tools.iedb.org> or supplementary data of articles cited in the Methods.

#### Ethics approval and consent to participate

Not applicable.

#### Consent for publication

Not applicable.

#### Competing interests

The authors declare that they have no competing interests.

Received: 14 November 2019 Accepted: 23 June 2020

Published online: 02 July 2020

#### References

1. Scheetz L, Park KS, Li Q, Lowenstein PR, Castro MG, Schwendeman A, Moon JJ. Engineering patient-specific cancer immunotherapies. *Nat Biomed Eng.* 2019. <https://doi.org/10.1038/s41551-019-0436-x>.
2. Sahin U, Türeci Ö. Personalized vaccines for cancer immunotherapy. *Science.* 2018;359(6382):1355–60. <https://doi.org/10.1126/science.aar7112>.
3. Schumacher TN, Schreiber RD. Neoantigens in cancer immunotherapy. *Science.* 2015;348(6230):69–74. <https://doi.org/10.1126/science.aaa4971>.
4. Hu Z, Ott PA, Wu CJ. Towards personalized, tumour-specific, therapeutic vaccines for cancer. *Nat Rev Immunol.* 2018;18(3):168. <https://doi.org/10.1038/nri.2017.131>.
5. Zhao W, Sher X. Systematically benchmarking peptide-MHC binding predictors: From synthetic to naturally processed epitopes. *PLoS Comput Biol.* 2018;14(11):1006457. <https://doi.org/10.1371/journal.pcbi.1006457>.
6. Mei S, Li F, Leier A, Marquez-Lago TT, Giam K, Croft NP, Akutsu T, Smith AI, Li J, Rossjohn J, Purcell AW, Song J. A comprehensive review and performance evaluation of bioinformatics tools for HLA class I peptide-binding prediction. *Brief Bioinformatics.* 2019. <https://doi.org/10.1093/bib/bbz051>.
7. Kim Y, Sidney J, Pinilla C, Sette A, Peters B. Derivation of an amino acid similarity matrix for peptide:MHC binding and its application as a Bayesian prior. *BMC Bioinformatics.* 2009;10(1):394. <https://doi.org/10.1186/1471-2105-10-394>.
8. Moutaftsi M, Peters B, Pasquetto V, Tschärke DC, Sidney J, Bui H-H, Grey H, Sette A. A consensus epitope prediction approach identifies the breadth of murine TCD8+ cell responses to vaccinia virus. *Nat Biotechnol.* 2006;24(7):817–9. <https://doi.org/10.1038/nbt1215>.
9. Andreatta M, Nielsen M. Gapped sequence alignment using artificial neural networks: application to the MHC class I system. *Bioinformatics.* 2015;32(4):511–7. <https://doi.org/10.1093/bioinformatics/btv639>.
10. Jurtz V, Paul S, Andreatta M, Marcatili P, Peters B, Nielsen M. NetMHCpan-4.0: Improved Peptide-MHC Class I Interaction Predictions Integrating Eluted Ligand and Peptide Binding Affinity Data. *J Immunol.* 2017;199(9):3360–8. <https://doi.org/10.4049/jimmunol.1700893>.
11. O'Donnell TJ, Rubinsteyn A, Bonsack M, Riemer AB, Laserson U, Hammerbacher J. MHCflurry: Open-Source Class I MHC Binding Affinity Prediction. *Cell Syst.* 2018;7(1):129–1324. <https://doi.org/10.1016/j.cels.2018.05.014>.

12. Bonsack M, Hoppe S, Winter J, Tichy D, Zeller C, Küpper MD, Schitter EC, Blatnik R, Riemer AB. Performance Evaluation of MHC Class-I Binding Prediction Tools Based on an Experimentally Validated MHC–Peptide Binding Data Set. *Cancer Immunol Res.* 2019;7(5):719–36. <https://doi.org/10.1158/2326-6066.cir-18-0584>.
13. Bhattacharya R, Sivakumar A, Tokheim C, Guthrie VB, Anagnostou V, Velculescu VE, Karchin R. Evaluation of machine learning methods to predict peptide binding to MHC Class I proteins. *bioRxiv.* 2017. <https://doi.org/10.1101/154757>.
14. Phloyphisut P, Pornputtpong N, Sriswasdi S, Chuangsuwanich E. MHCSeqNet: a deep neural network model for universal MHC binding prediction. *BMC Bioinformatics.* 2019;20(1):. <https://doi.org/10.1186/s12859-019-2892-4>.
15. Strodthoff N, Wagner P, Wenzel M, Samek W. UDSMProt: universal deep sequence models for protein classification. *Bioinformatics.* 2020;36(8):2401–9. <https://doi.org/10.1093/bioinformatics/btaa003>.
16. Howard J, Ruder S. Universal language model fine-tuning for text classification. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Melbourne: Association for Computational Linguistics; 2018. p. 328–339. <https://doi.org/10.18653/v1/P18-1031>.
17. Merity S, Keskar NS, Socher R. Regularizing and optimizing LSTM language models. *arXiv preprint arXiv:1708.02182.* 2017.
18. Smith LN. A disciplined approach to neural network hyper-parameters: Part 1 - learning rate, batch size, momentum, and weight decay. *arXiv preprint arXiv:1803.09820.* 2018.
19. Kim Y, Sidney J, Buus S, Sette A, Nielsen M, Peters B. Dataset size and composition impact the reliability of performance benchmarks for peptide-MHC binding predictions. *BMC Bioinformatics.* 2014;15(1):241. <https://doi.org/10.1186/1471-2105-15-241>.
20. Vita R, Mahajan S, Overton JA, Dhanda SK, Martini S, Cantrell JR, Wheeler DK, Sette A, Peters B. The Immune Epitope Database (IEDB): 2018 update. *Nucleic Acids Res.* 2018;47(D1):339–43. <http://dx.doi.org/10.1093/nar/gky1006>.
21. Wang P, Sidney J, Kim Y, Sette A, Lund O, Nielsen M, Peters B. Peptide binding predictions for HLA DR, DP and DQ molecules. *BMC Bioinformatics.* 2010;11(1):568. <https://doi.org/10.1186/1471-2105-11-568>.
22. Paul S, Weiskopf D, Angelo MA, Sidney J, Peters B, Sette A. HLA Class I Alleles Are Associated with Peptide-Binding Repertoires of Different Size, Affinity, and Immunogenicity. *J Immunol.* 2013;191(12):5831–9. <https://doi.org/10.4049/jimmunol.1302101>.
23. Chen J, Guo M, Wang X, Liu B. A comprehensive review and comparison of different computational methods for protein remote homology detection. *Brief Bioinformatics.* 2016;19(2):231–44. <https://doi.org/10.1093/bib/bbw108>.
24. Alipanahi B, DeLong A, Weirauch MT, Frey BJ. Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning. *Nat Biotechnol.* 2015;33(8):831–8. <https://doi.org/10.1038/nbt.3300>.
25. Manning CD, Raghavan P, Schütze H. Introduction to Information Retrieval. New York: Cambridge University Press; 2008. <https://doi.org/10.1017/cbo9780511809071>.
26. Nielsen M, Lundegaard C, Lund O, Keşmir C. The role of the proteasome in generating cytotoxic T-cell epitopes: insights obtained from improved predictions of proteasomal cleavage. *Immunogenetics.* 2005;57(1):33–41. <https://doi.org/10.1007/s00251-005-0781-7>.
27. Paszke A, Gross S, Chintala S, Chanan G, Yang E, DeVito Z, Lin Z, Desmaison A, Antiga L, Lerer A. Automatic differentiation in PyTorch. In: 31st Conference on Neural Information Processing Systems (NIPS) Workshop Autodiff; 2017.
28. Howard J, et al. fast.ai. GitHub. 2018. <https://github.com/fastai/fastai>. Accessed 26 Apr 2019.

## Publisher's Note

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

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

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

**At BMC, research is always in progress.**

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

