

PEPVAC: a web server for multi-epitope vaccine development based on the prediction of supertypic MHC ligands

Pedro A. Reche^{1,2,*} and Ellis L. Reinherz^{1,2}

¹Laboratory of Immunobiology and Department of Medical Oncology, Dana-Farber Cancer Institute and

²Department of Medicine, Harvard Medical School, 44 Binney Street, Boston, MA 02115, USA

Received December 23, 2004; Revised January 18, 2005; Accepted January 31, 2005

ABSTRACT

Prediction of peptide binding to major histocompatibility complex (MHC) molecules is a basis for anticipating T-cell epitopes, as well as epitope discovery-driven vaccine development. In the human, MHC molecules are known as human leukocyte antigens (HLAs) and are extremely polymorphic. HLA polymorphism is the basis of differential peptide binding, until now limiting the practical use of current epitope-prediction tools for vaccine development. Here, we describe a web server, PEPVAC (Promiscuous EPitope-based VACcine), optimized for the formulation of multi-epitope vaccines with broad population coverage. This optimization is accomplished through the prediction of peptides that bind to several HLA molecules with similar peptide-binding specificity (supertypes). Specifically, we offer the possibility of identifying promiscuous peptide binders to five distinct HLA class I supertypes (A2, A3, B7, A24 and B15). We estimated the phenotypic population frequency of these supertypes to be 95%, regardless of ethnicity. Targeting these supertypes for promiscuous peptide-binding predictions results in a limited number of potential epitopes without compromising the population coverage required for practical vaccine design considerations. PEPVAC can also identify conserved MHC ligands, as well as those with a C-terminus resulting from proteasomal cleavage. The combination of these features with the prediction of promiscuous HLA class I ligands further limits the number of potential epitopes. The PEPVAC server is hosted

by the Dana-Farber Cancer Institute at the site <http://immunax.dfci.harvard.edu/PEPVAC/>.

INTRODUCTION

T-cells are the key component of the adaptive immune system, playing a pivotal role fighting both infectious agents and cancer cells (1). T-cell-based immune responses are driven by antigenic peptides (epitopes), presented in the context of major histocompatibility complex (MHC) molecules (2). Therefore, the prediction of peptides that can bind to MHC molecules has become the basis for the anticipation of T-cell epitopes (3). MHC molecules fall into two major classes, namely MHC class I (MHCI) and MHC class II (MHCII). Antigens presented by MHCI and MHCII are recognized by two distinct sets of T-cells, CD8⁺ T and CD4⁺ T-cells, respectively. Identification of T-cell epitopes is important for both understanding disease pathogenesis and vaccine design. Thus, the availability of computational methods that can readily identify potential epitopes from primary protein sequences has fueled a new paradigm in vaccine development that is driven by this epitope discovery.

A major complication to this vaccine development approach is the extreme polymorphism of the MHC molecules. In the human, MHC molecules are known as human leukocyte antigens (HLAs), and there are hundreds of allelic variants of the class I (HLA I) and the class II (HLA II) molecules. These HLA allelic variants bind distinct sets of peptides as MHC polymorphism is the basis for peptide-binding specificity (4), and are expressed at vastly variable frequencies in different ethnic groups (5). This complexity suggests that a large number of HLA molecules will have to be targeted for peptide-binding predictions, requiring so many peptides to elicit a broadly protective multi-epitope vaccine as to be

*To whom correspondence should be addressed. Tel: +1 617 632 3824; Fax: +1 617 632 3351; Email: reche@research.dfci.harvard.edu
Correspondence may also be addressed to Ellis L. Reinherz. Tel: +1 617 632 3412; Fax: +1 617 632 3351; Email: ellis_reinherz@dfci.harvard.edu

impractical. Interestingly, groups of several HLA molecules (supertypes) can bind largely overlapping sets of peptides (6,7). The identification of these HLA supertypes facilitates the epitope-based vaccine development for the following two reasons: first, targeting of representative HLA alleles from distinct supertypes allows the immune response to be stimulated in a variety of genetic backgrounds; second, the selection of promiscuous peptide binders to those alleles included within a given supertype limits the number of peptides to be considered without decreasing the spectrum of the immune response.

In this paper, we describe a web server, PEPVAC (Promiscuous EPitope-based VACCine), that allows the prediction of promiscuous epitopes to five HLA I supertypes: A2 (A*0201-07, A*0209 and A*6802), A3 (A*0301, A*1101, A*3101, A*3301, A*6801 and A*6601), A24 (A*2402 and B*3801), B7 (B*0702, B*3501, B*5101-02, B*5301 and B*5401) and B15 (A*0101, B*1501_B62 and B1502). These supertypes were defined using a method based on the clustering of the predicted peptide-binding repertoire of MHC molecules (8). The combined phenotypic frequency of these supertypes is >95% for five major American ethnicities (Black, Caucasian, Hispanic, Native American and Asian). Thus, targeting these supertypes with epitope predictions would potentially provide a population coverage $\geq 95\%$, regardless of ethnicity.

Peptides binding to HLA I molecules are potential CD8⁺ T-cell epitopes. *In vivo*, the C-terminus of these antigenic epitopes results from the selective proteolysis of cytosolic proteins mediated by the proteasome (9). The proteasome is thus important for determining these epitopes. Therefore, PEPVAC has also been implemented with an algorithm for the identification of those peptides containing a C-terminus that is likely to be the result of proteasomal cleavage. Finally, PEPVAC also allows the prediction of conserved epitopes from sequences with variability masked. The combination of these two features serves in both refining the predictions of T-cell epitopes and limiting the number of potential epitopes.

Prediction of peptide-MHCI binding

The peptide-binding mode of MHCII molecules differs from that of MHCI (10–12), and as result, the prediction of peptide-MHCII binding is less reliable than that of peptide-MHCI binding. Thereby, we have focused here in the prediction of MHCI ligands, a class that is specifically recognized by CD8⁺ cytotoxic T lymphocytes. Peptides binding to a specific MHCII molecule are related by sequence similarity, and thus we use position-specific scoring matrix (PSSM) from aligned MHCII ligands as the predictors of peptide-MHCI binding in combination with a dynamic algorithm. PSSMs are also known as profiles and weight matrices and have previously been shown to be adequate tools for the prediction of peptide-MHC binding (13–16). PSSMs are derived from block alignments of MHCII ligands that are of the same length. Such a restriction guarantees proper structural alignment of ligands and subsequent accuracy of the peptide-binding predictions (13,14). Given that MHCII-ligands are usually of nine residues in length, PSSMs used in this study are for the prediction of ligands of that same size (nine residues). Accuracy of the prediction of peptide-MHCI binding using PSSMs varies depending on threshold and the targeted MHCII molecule.

On average, however, ROC analyses of the predictions at different thresholds result in *AUC* values (Area Under ROC Curve) above 0.8, indicating that these PSSMs are very good for predictors of peptide-MHCI binding. Furthermore, >80% of known CD8⁺ T-cell epitopes can be predicted at a 2% threshold from their protein sources.

Supertypes: identification and population coverage analysis

We defined HLA I supertypes through clustering of predicted MHC peptide-binding repertoires (8). In brief, the core of the method consists of the generation of a distance matrix whose coefficients are inversely proportional to the peptide binders shared by any two HLA molecules (Figure 1). Subsequently, this distance matrix is fed to a phylogenetic clustering algorithm to establish the kinship among the distinct HLA peptide-binding repertoires. Figure 2 shows a phylogenetic tree built upon the peptide-binding repertoire of 55 HLA I molecules, using a Fitch and Margoliash clustering algorithm (17). We defined supertypes (Figure 2) as groups of HLA I alleles with $\geq 20\%$ peptide-binding overlap (pairwise between any pair of alleles). The supertypes identified in this study include the A2, A3, B7, B27 and B44 supertypes previously identified by Sidney *et al.* (16). Furthermore, we have also identified three new supertypes, BX, B15 and B57 (Figure 2). The cumulative phenotypic frequency (CPF) of these supertypes is shown in Table 1. CPF was calculated using the gene and haplotype frequencies reported for five distinct American ethnic groups including Blacks, Caucasians, Hispanic, North American

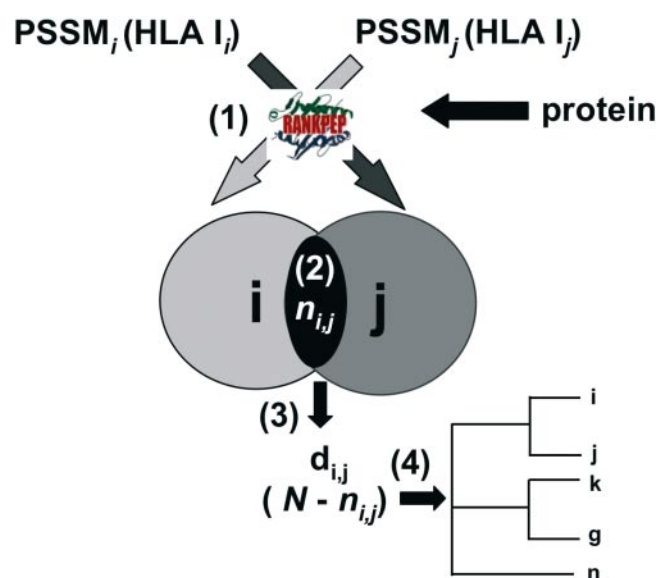


Figure 1. Strategy to define HLA I supertypes. HLA I supertypes are identified by clustering their peptide-binding repertoire (8). The method consists of four basic steps. (i) Predict the peptide-binding repertoire (i, j sets in figure) of each HLA I molecule from the same random protein using the relevant PSSMs in combination with the RANKPEP scoring algorithm (13). (ii) Compute the number of common peptides between the binding repertoire of any two HLA I molecules. (iii) Build a distance matrix whose coefficients are inversely proportional to the peptide-binding overlap between any pair of HLA I molecules. (iv) Use a phylogenetic clustering algorithm to compute and visualize HLA I supertypes (clusters of HLA I molecules with overlapping peptide-binding repertoires).

A

PEPVAC

GENOME WIDE PREDICTION OF PROMISCUOUS EPITOPES FOR VACCINE DESIGN

T-cell vaccines
HLA-peptide binding
HLA-coverage
Proteasome Cleavage

Function:
PEPVAC is a tool aimed to the development of fully covering multi-epitope vaccines against pathogenic organisms based on genome wide predictions of promiscuous MHC I-restricted epitopes

Description:
T-cell epitopes are first anticipated based on their binding to HLA I molecules using profile-matrices, and then filtered for immunoproteasomal cleavage using a probabilistic model. HLA I molecules present many allelic variants with distinct peptides specificities, and thus peptide binding predictions are given for a set of HLA I alleles with a combined phenotypic frequency ~ 95% in 5 distinct ethnic groups. These HLA I alleles are grouped in sets (supertypes) whose predicted binding peptides are largely overlapping. Only the peptides predicted to bind to all HLA I alleles included in each supertype are returned as potential T-cell epitopes. Identification of these promiscuous peptide binders allows to minimize the total number of predicted epitopes without compromising the population coverage required in the design of multi-epitope vaccines.

E-MAIL [\[help\]](#)

Enter your e-mail

GENOMES [\[help\]](#)

Select a genome

- Severe Acute Respiratory Syndrome (SARS) virus
- Influenza Virus A (PR/8)
- Variola Major virus (Strain India)
- Human Immunodeficiency Virus 1 (HIV1) (B clade)

Upload genome (File with translated ORFs in FASTA format)

no file selected

SUPERTYPES [\[help\]](#)

Select HLA-Supertype/s

A2: A*0201, A*0202, A*0203, A*0205, A*0206, A*0207, A6802

A3: A*0301, A*1101, A*3101, A*3301, A*6801, A*6601

A24: A*2402, B*3801

B7: B*0702, B*3501, B*5101, B*5102, B*5301, B*5401

B15: A*0101, B*1501_B62, B1502

Population Coverage with selected supertypes [\[help\]](#)

32.18%

PROTEASOMAL CLEAVAGE [\[help\]](#)

Filter: OFF Model: 1

B

RESULTS

PEPVAC: Peptide based vaccines using promiscuous MHC I-restricted epitopes

SUMMARY

Genome: Influenza Virus A (PR/8)

- Genome Size: 4617 amino acids
- ORFs: 11

Immunoproteasome Filter: OFF
Selected SuperAntigens: A3
Predicted Peptides: 37 (0.80% of all 9mer peptides from selected GENOME)
Minimum Population Coverage: 32.18%

HLA SUPERTYPE

HLA alleles included: A*3301, A*1101, A*3101, A*0301, A*6801, A*6601
Minimum population coverage: 32.18% [\[help\]](#)

A3 promiscuous peptides: 37

RANK	SOURCE GI	POS.	N	SEQUENCE	C	MW (Da)	SCORE	% OPT.
1	P03452	217	AYV	SVVTSNYNR	RFT	1021.09	101.0	82.11 %
2	P03485	179	ENR	MVLASTTAK	AME	903.09	96.0	78.05 %
3	P03431	342	IAP	IMFSNKMAR	LGK	1079.34	93.0	75.61 %
4	P03431	221	LJR	ALTLNTMTK	DAE	974.16	91.0	73.98 %
5	P03431	713	GIS	SMVEAMVSR	ARI	991.19	88.0	71.54 %
6	P03433	601	AES	SVKEDDMTK	EFF	1047.22	87.0	70.73 %
7	P03428	372	RAT	AILRKATRR	LIQ	1066.32	85.0	69.11 %
9	P03466	142	WHS	NLNDATYQR	TRA	1076.13	85.0	69.11 %
10	P03428	113	HYP	KIYKTYPER	VER	1229.45	84.0	68.29 %
11	P03428	9	ELR	NLMSQSRTR	EIL	1074.22	83.0	67.48 %
12	P03431	471	RTC	KLLGINMSK	KKS	985.24	79.0	64.23 %
13	P03431	355	GKG	YMFESKSMK	LRT	1132.36	79.0	64.23 %
14	P03452	177	SYP	KLKNSYVNK	KGK	1075.26	79.0	64.23 %
15	P03433	177	RLF	TIRQEMASR	GLW	1073.24	79.0	64.23 %
18	P03431	199	TKK	MITQRTIGK	RKQ	1029.25	78.0	63.41 %
22	P03433	627	SPK	GVEESSIGK	VCR	886.96	75.0	60.98 %
24	P03431	190	RKR	RVRDNTMFK	MIT	1129.33	75.0	60.98 %
25	75188	111	ALS	NIRVSSRSK	TTL	1028.18	75.0	60.98 %
27	P03431	578	EIK	KLWEQTRSK	AGL	1134.33	74.0	60.16 %
28	P03433	162	ADY	TLDEFSSAR	IKT	1058.13	74.0	60.16 %

Figure 3. The PEPVAC web server. (A) PEPVAC input page. The page is divided into several sections. E-MAIL, for obtaining the results via e-mail (optional). GENOMES, where a selection of genomes from pathogenic organisms is available, as well as the possibility of uploading a user-provided genome. SUPERTYPES, the supertypes A2, A3, B7, A24, and B15 are available for selection. Alleles targeted for peptide-binding predictions in each supertype are indicated. The minimum population coverage of the selected supertypes is calculated on the fly and shown on the relevant window. PROTEASOMAL CLEAVAGE, prediction of proteasomal cleavage using three optimal language models is carried out in parallel to the peptide-binding predictions. (B) PEPVAC result page. An example result page where the A3 supertype was selected for peptide-binding predictions from the genome of *Influenza A virus (A/PR/8/34)* is shown. The result page first displays a summary of the predictions, followed by the predicted peptide binders to each of the selected supertypes (only A3 in the shown example). Peptides highlighted in violet contain a C-terminal residue that is predicted to be the result of proteasomal cleavage. If the proteasomal cleavage filter is checked ON in the input page, only violet peptides will be shown.

site <http://immunax.dfc.harvard.edu/PEPVAC/> hosted by the Molecular Immunology Foundation/Dana-Farber Cancer Institute. The web interface to PEPVAC is divided into several sections that facilitate intuitive use (Figure 3A). Main features of the web server are discussed below.

Input and limitations. In PEPVAC, input query to carry epitope predictions is entered in the GENOME section (Figure 3A). Input consists of a single or various protein sequences in FASTA format. Only the standard 20 amino acid residues are considered. There are several translated genomes from pathogenic organisms that can be selected as inputs. More useful, a user-provided local file containing a set of protein sequences can be uploaded to the server using the choose/browse bottom. PEPVAC can also process files with protein sequences, in which the variable sites have been masked with a dot '.' symbol. In that case, peptide-binding predictions will be carried out only over consecutive stretches of nine or more residues. Sequences with variable positions masked according to the Shannon entropy variability metric (4,19) can be obtained at the site <http://immunax.dfc.harvard.edu/bioinformatics/Tools/sva.html>. Currently, there is a limit

of 200 sequences and 50 000 symbols that can be processed per request. If such limits are exceeded, the server will return an error.

Supertypes and thresholds. The A2, A3, B7, B15 and A24 (Figure 2 and Table 1) supertypes have been chosen for promiscuous peptide-binding predictions in PEPVAC. Only those peptides that are predicted to bind to all the alleles included in the supertypes are returned in the output (Figure 3B). Threshold for the prediction of promiscuous peptide binders in PEPVAC has been fixed to provide a reduced and manageable set of promiscuous peptide binders to each supertype. As an example, predicted promiscuous peptides to the above five supertypes from a genome, such as that of *Influenza virus A* (4160 amino acids) distributed in 10 distinct open reading frames, represent only 5.51% (254 9mer peptides) of all possible peptides (4617 9mer peptides).

Proteasome cleavage. In PEPVAC, predictions of supertypic peptide binders are combined with the prediction of proteasomal cleavage using probabilistic language models derived from HLA I-restricted epitopes (14). Currently, there are three optional models for proteasomal cleavage

that differ in their sensitivity/specificity ratio of the predictions as discussed elsewhere (14). These models are selected within the PROTEASOME CLEAVAGE section. Model 1 has the highest sensitivity (~95%) and the lower specificity (~60%). Conversely, Model 3 has the lowest sensitivity (65%) with the largest specificity (80%). Model 2 has a sensitivity and specificity of ~70%. Promiscuous peptide binders containing a C-terminal end, predicted to be the result of proteasomal cleavage, are shown in violet in the result page (Figure 3B). In the previous example with the *Influenza virus A*, the list of promiscuous peptide binders to the five selected supertypes decreases from 254 down to 170 peptides (3.7% of all 9mer peptides from *Influenza virus A* genome) after considering proteasomal cleavage using Model 1. Furthermore, a combination of the predictions of peptide-MHCI binding and proteasomal cleavage increases the specificity of the epitope predictions by discarding predicted peptide-MHCI binders that are experimentally unable to elicit CD8⁺ T-cell responses (20).

Output. The results page returned by PEPVAC is shown in Figure 3B. This page first displays a summary of the predictions, including the chosen selections, the number of predicted peptides and the minimum population coverage provided by the supertypic selection, followed by the predicted peptide binders to each of the selected supertypes (only A3 in the shown example). Peptides are predicted to bind to all alleles included in the supertype, and appear ranked with regard to the PSSMs of the first allele included in the supertype. Relevant information about each sorted peptide includes its protein source as well as its molecular weight.

ACKNOWLEDGEMENTS

This manuscript was supported by NIH grant AI50900 and the Molecular Immunology Foundation. We wish to acknowledge John-Paul Glutting for programming assistance. Funding to pay the Open Access publication charges for this article was provided by NIH grant AI50900.

Conflict of interest statement. None declared.

REFERENCES

- Paul, W.E. (ed.) (1998) *Fundamental Immunology*. 4th Edn. Raven Press, New York, NY.
- Zinkernagel, R.M. and Doherty, P.C. (1974) Restriction of *in vitro* T cell-mediated cytotoxicity in lymphocytic choriomeningitis within a syngeneic or semiallogeneic system. *Nature*, **248**, 701–702.
- Flower, D.R. and Doytchinova, I.A. (2002) Immunoinformatics and the prediction of immunogenicity. *Appl. Bioinformatics*, **1**, 167–176.
- Reche, P.A. and Reinherz, E.L. (2003) Sequence variability analysis of human class I and class II MHC molecules: functional and structural correlates of amino acid polymorphisms. *J. Mol. Biol.*, **331**, 623–641.
- Gjertson, D.W. and Terasaki, P.I. (1998) *HLA 1998*. ASHI Publications, Lenexa, KS.
- Sette, A. and Sidney, J. (1999) Nine major HLA class I supertypes account for the vast preponderance of HLA-A and -B polymorphism. *Immunogenetics*, **50**, 201–212.
- Sette, A. and Sidney, J. (1998) HLA supertypes and supermotifs: a functional perspective on HLA polymorphism. *Curr. Opin. Immunol.*, **10**, 478–482.
- Reche, P.A. and Reinherz, E.L. (2004) Definition of MHC supertypes through clustering of MHC peptide binding repertoires. *Artificial Immune Systems*. Springer-Verlag Berlin Heidelberg, Catania, Italy, Vol. LNCS 3239, pp. 189–196.
- Craiu, A., Akopian, T., Goldberg, A. and Rock, K.L. (1997) Two distinct proteolytic processes in the generation of a major histocompatibility complex class I-presented peptide. *Proc. Natl Acad. Sci. USA*, **94**, 10850–10855.
- Madden, D.R. (1995) The three-dimensional structure of peptide-MHC complexes. *Annu. Rev. Immunol.*, **13**, 587–622.
- Madden, D.R., Garboczi, D.N. and Wiley, D.C. (1993) The antigenic identity of peptide-MHC complexes: a comparison of the conformations of five viral peptides presented by HLA-A2. *Cell*, **75**, 693–708.
- Stern, L.J. and Wiley, D.C. (1994) Antigen peptide binding by class I and class II histocompatibility proteins. *Structure*, **2**, 245–251.
- Reche, P.A., Glutting, J.P. and Reinherz, E.L. (2002) Prediction of MHC class I binding peptides using profile motifs. *Hum. Immunol.*, **63**, 701–709.
- Reche, P.A., Glutting, J.-P. and Reinherz, E.L. (2004) Enhancement to the RANKPEP resource for the prediction of peptide binding to MHC molecules using profiles. *Immunogenetics*, **56**, 405–419.
- Nielsen, M., Lundegaard, C., Warming, P., Hvid, C.S., Lambert, K., Buus, S., Brunak, S. and Lund, O. (2004) Improved prediction of MHC class I and class II epitopes using a novel Gibbs sampling approach. *Bioinformatics*, **20**, 1388–1397.
- Peters, B., Tong, W., Sidney, J., Sette, A. and Weng, Z. (2003) Examining the independent binding assumption for binding of peptide epitopes to MHC-I molecules. *Bioinformatics*, **19**, 1765–1772.
- Fitch, W.M. and Margoliash, E. (1967) Construction of phylogenetic trees. *Science*, **155**, 279–284.
- Cao, K., Hollenbach, J., Shi, X., Shi, W., Chopek, M. and Fernandez-Vina, M.A. (2001) Analysis of the frequencies of HLA-A, B, and C alleles and haplotypes in the five major ethnic groups of the United States reveals high levels of diversity in these loci and contrasting distribution patterns in these populations. *Hum. Immunol.*, **62**, 1009–1030.
- Stewart, J.J., Lee, C.Y., Ibrahim, S., Watts, P., Shlomchik, M., Weigert, M. and Litwin, S. (1997) A Shannon entropy analysis of immunoglobulin and T cell receptor. *Mol. Immunol.*, **34**, 1067–1082.
- Zhong, W., Reche, P.A., Lai, C.C., Reinhold, B. and Reinherz, E.L. (2003) Genome-wide characterization of a viral cytotoxic T lymphocyte epitope repertoire. *J. Biol. Chem.*, **278**, 45135–45144.
- Dawson, D.V., Ozgur, M., Sari, K., Ghanayem, M. and Kostyu, D.D. (2001) Ramifications of HLA class I polymorphism and population genetics for vaccine development. *Genet. Epidemiol.*, **20**, 87–106.