



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

**Integral use of immunopeptidomics and immunoinformatics for the characterization of antigen presentation and rational identification of BoLA-DR- presented peptides and epitopes**

**Citation for published version:**

Fisch, A, Reynisson, B, Benedictus, L, Nicastrì, A, Vasoya, D, Morrison, I, Buus, S, Ferreira, BR, de Miranda Santo, IKF, Ternette, N, Connelley, T & Nielsen, M 2021, 'Integral use of immunopeptidomics and immunoinformatics for the characterization of antigen presentation and rational identification of BoLA-DR- presented peptides and epitopes', *Journal of Immunology*. <https://doi.org/10.4049/jimmunol.2001409>

**Digital Object Identifier (DOI):**

[10.4049/jimmunol.2001409](https://doi.org/10.4049/jimmunol.2001409)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

Journal of Immunology

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



1 **Integral use of immunopeptidomics and immunoinformatics**  
2 **for the characterization of antigen presentation and rational**  
3 **identification of BoLA-DR-presented peptides and**  
4 **epitopes<sup>1,2,3,4</sup>**

5

6 Andressa Fisch<sup>\*,i</sup>, Birkir Reynisson<sup>†,i</sup>, Lindert Benedictus<sup>‡</sup>, Annalisa Nicastrì<sup>§</sup>, Deepali  
7 Vasoya<sup>‡</sup>, Ivan Morrison<sup>‡</sup>, Søren Buus<sup>¶</sup>, Beatriz Rossetti Ferreira<sup>\*</sup>, Isabel Kinney Ferreira de  
8 Miranda Santos<sup>||</sup>, Nicola Ternette<sup>§</sup>, Tim Connelley<sup>‡</sup>, Morten Nielsen<sup>†,#,^</sup>

9

10 \* Ribeirão Preto College of Nursing, University of São Paulo, Av Bandeirantes 3900, Ribeirão Preto,  
11 Brazil

12 † Department of Health Technology, Technical University of Denmark, DK-2800 Lyngby, Denmark

13 ‡ The Roslin Institute, Edinburgh, Midlothian EH25 9RG, UK

14 § The Jenner Institute, Nuffield Department of Medicine, Oxford, OX37BN, UK

15 ¶ Laboratory of Experimental Immunology, Department of Immunology and Microbiology, Faculty of  
16 Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark

17 || Ribeirão Preto Medical School, University of São Paulo, Av Bandeirantes 3900, Ribeirão Preto, Brazil

18 # Instituto de Investigaciones Biotecnológicas, Universidad Nacional de San Martín, CP1650 San  
19 Martín, Argentina

20 Running Title: **NetBoLAIIPan**

21 ^ Corresponding author: mniel@dtu.dk

22 <sup>i</sup> Shared first authorship

---

<sup>1</sup> This work was supported in part by funds from the National Institute of Allergy and Infectious Diseases, National Institutes of Health, Department of Health and Human Services, under Contract No. HHSN272201200010C and from the Fundação de Amparo à Pesquisa do Estado de São Paulo – FAPESP (2015/09683-9) to BRF.

<sup>2</sup> The Bill & Melinda Gates Foundation and with UK aid from the UK Foreign, Commonwealth and Development Office (Grant Agreement OPP1127286) under the auspices of the Centre for Tropical Livestock Genetics and Health (CTLGH), established jointly by the University of Edinburgh, SRUC (Scotland's Rural College), and the International Livestock Research Institute (the findings and conclusions contained within are those of the authors and do not necessarily reflect positions or policies of the Bill & Melinda Gates Foundation nor the UK Government).

<sup>3</sup> The BBSRC through the ISP award made to The Roslin Institute (BBS/E/D/20002174).

<sup>4</sup> AF was supported by FAPESP scholarships 2014/11010-9, 2017/21401-4 and 2018/23579-8.

## 23 **Abstract**

24 Major histocompatibility complex (MHC) peptide binding and presentation is the most  
25 selective event defining the landscape of T cell epitopes. Consequently, understanding the  
26 diversity of MHC alleles in a given population and the parameters that define the set of ligands  
27 that can be bound and presented by each of these alleles (the immunopeptidome) has an  
28 enormous impact on our capacity to predict and manipulate the potential of protein antigens to  
29 elicit functional T cell responses. Liquid chromatography-mass spectrometry (LC-MS)  
30 analysis of MHC eluted ligands (EL data) has proven to be a powerful technique for identifying  
31 such peptidomes, and methods integrating such data for prediction of antigen presentation have  
32 reached a high level of accuracy for both MHC class I and class II. Here, we demonstrate how  
33 these techniques and prediction methods can be readily extended to the bovine leukocyte  
34 antigen class II DR locus (BoLA-DR). BoLA-DR binding motifs were characterized by EL  
35 data derived from bovine cell lines expressing a range of DRB3 alleles prevalent in Holstein-  
36 Friesian populations. The model generated (NetBoLAIIpan - available as a web-server at  
37 [www.cbs.dtu.dk/services/NetBoLAIIpan](http://www.cbs.dtu.dk/services/NetBoLAIIpan)) was shown to have unprecedented predictive power  
38 to identify known BoLA-DR restricted CD4 epitopes. In summary, the results demonstrate the  
39 power of an integrated approach combining advanced MS peptidomics with  
40 immunoinformatics for characterization of the BoLA-DR antigen presentation system and  
41 provide a novel tool that can be utilised to assist in rational evaluation and selection of bovine  
42 CD4 T cell epitopes.

43

44

45

46

47 **Key Points**

48 • MS immunopeptidomics and motif characterization for 7 prevalent BoLA-DRB3  
49 molecules

50 • The first pan-specific predictor, NetBoLAIIpan, for BoLA-DRB3 antigen presentation

51 • NetBoLAIIpan demonstrated unprecedented CD4 T cell epitope prediction  
52 performance

53

54

55

56

57

58

59

60

61

62

63

64

65

66

## 67 Introduction

68 Major histocompatibility complex (MHC) genes play a vital role in the regulation of adaptive  
69 immunity. Whilst classical MHC class I genes are expressed on most nucleated cells, MHC  
70 class II (MHCII) molecules show a more restricted expression and are predominantly expressed  
71 on professional antigen-presenting cells such as dendritic cells, B-cells, and macrophages. The  
72 MHCII system enables peptides derived from both extracellular and intracellular proteins that  
73 have been delivered in the endocytic pathway to be loaded into the peptide-binding groove of  
74 MHCII molecules and be displayed as stable peptide-MHCII complexes (pMHCII) on the cell  
75 surface (1). CD4 T cells bearing cognate TCRs capable of binding specific pMHCII complexes  
76 can become activated and perform a range of functions, including supporting other immune  
77 effector cells such as macrophages, B cells and CD8 T cells (2). Thus, pMHCII molecules play  
78 a critical role in initiating and developing both humoral and cell-mediated adaptive immune  
79 responses.

80

81 MHCII molecules are heterodimers composed of an  $\alpha$  and  $\beta$  chain, each consisting of an  
82 extracellular domain, a transmembrane region, and an intracytoplasmic tail. The distal  
83 membrane domains ( $\alpha 1$  and  $\beta 1$ , respectively) form an open peptide-binding groove that binds  
84 peptides of variable length, mainly of 13–25 amino acid residues (3). The peptide-binding  
85 groove most often contains four major pockets that interact with the side-chains of anchoring  
86 residues located at positions 1, 4, 6, and 9 of the 9-mer binding-core of the bound ligand. These  
87 pockets thus determine the binding motif of the peptides that can be presented by an MHCII  
88 molecule (4, 5). A key feature of the MHC genes is the high level of polymorphism. For  
89 example in humans, three conventional MHCII heterodimers are expressed – DR, DQ and DP  
90 – and a total of ~2, ~2,500, ~100, ~1,200, ~80 and ~1,000 protein-coding variants of the  $\alpha$  (A)  
91 and  $\beta$  (B) chain genes, DRA, DRB, DQA, DQB, DPA, and DPB respectively, have been

92 identified. Except for DRA, the polymorphism of MHCII genes is focused predominantly  
93 within the  $\alpha 1$  and  $\beta 1$  domains (6), resulting in variations in the residues of the binding groove,  
94 and consequently determining the variable binding motifs and so the capacity of different  
95 MHCII molecules to bind different peptide sets.

96

97 In cattle, there are only two categories of conventional MHCII molecules, BoLA-DR and  
98 BoLA-DQ (7). The DRB, DQA, and DQB genes are highly polymorphic, whilst, as in other  
99 species, the DRA gene is essentially monomorphic (8). Although there are three DRB loci,  
100 only DRB3 is considered to be functionally expressed since DRB1 is a pseudogene and DRB2  
101 is expressed at very low levels if at all (9). Consequently, the variability of expressed BoLA-  
102 DR molecules can be characterized by sequencing of the DRB3 gene (10). The ability to  
103 perform rapid sequence-based typing of DRB3 using Sanger technology has resulted in DRB3  
104 being the most intensely studied bovine MHC gene (11–19), with 357 alleles registered in the  
105 IPD-MHC database (November 2020: <https://www.ebi.ac.uk/ipd/mhc/group/BoLA/>).

106

107 Characterisation of the peptide repertoires presented by different MHCII molecules can enable  
108 the development of algorithms that predict potential MHC binding peptides within proteins  
109 rapidly. Integration of large data sets of peptides directly eluted off MHC molecules and  
110 sequenced by mass-spectrometry (MS), so-called eluted ligand (EL) data, have facilitated the  
111 generation of accurate MHC-binding prediction algorithms (20–27). Such *in silico* tools can  
112 accelerate antigen selection for vaccine development and are of particular relevance to vaccines  
113 against pathogens with large proteomes (e.g. eukaryotic parasites), where screening and  
114 selection of candidate antigens from a large number of expressed proteins would be a major  
115 obstacle.

116

117 Analysis and interpretation of EL data are made challenging by ambiguous ligand MHC  
118 assignment resulting from the multiple MHC molecules expressed on the surface of most cells.  
119 Several approaches have been proposed to address this, spanning from the engineering of cell  
120 lines and/or expressed MHC molecules to allow for analysis of ligands of single MHC  
121 specificities (single allele (SA) ligands) (23, 24, 28) to computational motif deconvolution  
122 techniques (21, 22, 29) handling more complex multi-allele (MA) datasets. Within the latter  
123 category, the machine learning framework NNAlign\_MA (30) has been demonstrated to  
124 efficiently deconvolute MA ligand data obtained from samples expressing multiple MHC  
125 alleles, enabling the construction of improved pan-specific predictors for antigen presentation  
126 for both the MHC class I and class II systems (30–32). NNAlign\_MA achieves this by  
127 annotating the MA data during training in a semi-supervised manner based on MHC co-  
128 occurrence, MHC exclusion, and pan-specific binding prediction (30). This deconvolution  
129 expands the potential training data beyond binding affinity (BA) peptides and SA ligands to  
130 include the more complex and numerous MA ligands.

131

132 EL data differs from BA data in the sense that it not only captures peptide-MHC binding but  
133 also signals related to antigen processing. Recent MHCII prediction models (20, 21, 32) have  
134 leveraged these kinds of data and improved the prediction of MHCII antigen presentation.

135

136 Although most peptidome studies have focused on human and murine models, the technique  
137 can be equally applied to other species. In the context of livestock, we have earlier published  
138 studies demonstrating the ability to use mass spectrometry data to generate highly accurate  
139 prediction algorithms for BoLA-I molecules (33) which have been integrated into the  
140 NetMHCpan-4.1 server (31). Currently, there is no equivalent algorithm that can be used to  
141 predict peptide binding to BoLA-II molecules.

142

143 In this study, we have used mass-spectrometry to generate peptide elution data for BoLA-DR  
144 molecules and use the derived data to provide the first characterization of binding motifs of  
145 bovine MHCII and to demonstrate the development of the first available *in silico* method for  
146 accurate analysis of BoLA-DR ligands for rational CD4 T cell epitope prediction.

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

## 165 **Materials and Methods**

166 **Animal and cell samples.**



167 Brazilian Holstein-Friesian PBMC samples were obtained from frozen archived materials from  
168 animals within the herd at the University of Sao Paulo that had been included in previous  
169 experiments completed under approval from the Committee on the Ethics of Animals Research  
170 at the Nowavet Veterinary Clinical Studies CRO, Viçosa/MG, certificate numbers 56/2016  
171 (approved on 03 August 2016) and 36/2017 (approved on 09 June 2017). PBMC used for the  
172 characterization of BoLA-DR presented peptides from ovalbumin were isolated from a  
173 Holstein-Friesian animal from the University of Edinburgh herd with sampling conducted  
174 under a license granted under the UK Animal (Scientific Procedures) Act 1986. The *Theileria*  
175 *annulata*- and *Theileria parva*-infected cell lines used in this study had been established and  
176 characterised as part of previous studies and were maintained using routine and well-  
177 established protocols (34). Briefly, cattle PBMC from animals of interest, expressing the  
178 relevant BoLA alleles, were isolated by Ficoll and co-cultured with suspensions of *T. annulata*  
179 or *T. parva*-infected ticks to allow *in vitro* infection. The generated cell lines are transformed  
180 by the parasite and so proliferate indefinitely *in vitro*, while endogenously expressing high  
181 levels of MHCII (35). While it is known that *T. annulata* infects mostly B-cells and myeloid  
182 cells and *T. parva* infects T-cells, no further characterization was performed in the cell lines  
183 used in this study. The optimisation and final protocol used to assess the capacity of PBMC  
184 and *Theileria annulata*-infected cell lines to take up ovalbumin and present peptides on BoLA-  
185 DR molecules are described in Supplementary Figure 1.

186

187

188

189

190 **PBMC isolation, RNA extraction and cDNA synthesis.**

191 Bovine PBMC were isolated by density gradient centrifugation using Ficoll Paque Plus (GE  
192 Healthcare Bio-Sciences, Amersham. UK) according to manufacturers' instructions. RNA was  
193 extracted from PBMC using TRIzol (Thermo Scientific, Renfrew, UK) and cDNA synthesised  
194 using the GOscript Kit (Promega, Southampton, UK), both according to the manufacturers'  
195 instructions.

196

### 197 **BoLA-DRB3 sequencing.**

198 For BoLA-DRB3 amplification, primers (For - CCAGGGAGATCCAACCACATTTCC; Rev  
199 - TCGCCGCTGCACAGTGAAACTCTC) incorporating Illumina adaptors and multiplex  
200 identifier tags were obtained from IDT (Leuven, Belgium). PCR was performed using Phusion  
201 High Fidelity PCR kit (New England Biolabs), and the reaction was carried out in a final  
202 volume of 40  $\mu$ L containing 2  $\mu$ L of cDNA, 5X Phusion HF Buffer, 0.8 U  $\mu$ L of Phusion DNA  
203 Polymerase, 3% DMSO, 0.4 mM of dNTP and 0.5  $\mu$ M of each primer. The reaction was  
204 performed in a G-Storm Thermal Cycle System (G-Storm) programmed for one cycle at 98  $^{\circ}$ C  
205 for 30 s, followed by 30 cycles at 98  $^{\circ}$ C for 10 s, 61  $^{\circ}$ C for 30 s, and 72  $^{\circ}$ C for 45 s, with a final  
206 extension period at 72  $^{\circ}$ C for 10 min. 5  $\mu$ L of PCR product from each sample were pooled  
207 together, run on a 1.5% agarose gel, and the band of the appropriate size was extracted and  
208 purified using the QIAquick PCR Purification Kit (Qiagen). A final purification using  
209 Agencourt AMPure XP Beads (Beckman Coulter) at a ratio of 1:1 beads to PCR product was  
210 conducted prior to quantification of the sample and submission to Edinburgh Genomics for  
211 sequencing on the Illumina MiSeq V.3 platform. Analysis of the data was conducted using a  
212 bespoke bioinformatics pipeline (Vasoya *et al.* in preparation).

213

214

### 215 **pBoLA-DR complexes purification.**

216 Cultured cells ( $1 \times 10^9$ ) were washed twice with ice-cold PBS and then lysed in buffer (1%  
217 IGEPAL, 15mM TRIS pH 8.0, 300 mM NaCl and cOmplete protease inhibitor (Roche)) at a  
218 density of  $2 \times 10^8$  cells/mL for 1 min, diluted with PBS 1:1 and solubilized for 45 min at 4 °C.  
219 Lysates were cleared by two-step centrifugation at 500g for 15 min at 4 °C and then at 15,000g  
220 for 45 min at 4 °C. For initial samples pBoLA-DR complexes were directly captured from the  
221 cleared lysates using 5 mg anti-BoLA-DR antibody (ILA21), immobilized in 1 mL of protein  
222 A resin (Amintra, Expedeon, Cambridge, UK). For later samples, pBoLA-DR complexes were  
223 captured from cleared lysates that had been depleted of peptide-BoLA-I (pBoLA-I) complexes  
224 by prior immunoprecipitation with 5 mg anti-BoLA-I antibody (ILA88), immobilized in 1 mL  
225 protein A resin. Captured pBoLA-DR complexes were washed, and peptides eluted from  
226 BoLA-DR molecules using 10% acetic acid and the resulting proteins dried as described in  
227 (36).

228

#### 229 **HPLC.**

230 The dried pBoLA-DRB3 complexes were resuspended in 150  $\mu$ L of loading buffer (0.1%  
231 formic acid, 1% acetonitrile) and loaded onto a  $4.6 \times 50$  mm ProSwift<sup>TM</sup> RP-1S column  
232 (Thermo Scientific) for reverse-phase chromatography on an Ultimate 3000 HPLC system  
233 (Thermo Scientific). Elution was performed using a 0.5 mL/min flow rate over 5 min on a  
234 gradient of 2 to 35% buffer B (0.1% formic acid in acetonitrile) in buffer A (0.1% formic acid).  
235 Eluted fractions were collected from 1 to 8.5 min, for 30 s each. Protein detection was  
236 performed at 280 nm. Even and odd eluted fractions were pooled together, vacuum dried and  
237 stored at -80 °C until use.

238

239

#### 240 **LC-MS/MS.**

241 Dried samples were resuspended in 20  $\mu$ L of loading buffer and analyzed in an Ultimate 3000  
242 nano UPLC system online coupled to an Orbitrap Fusion™ Lumos™ Tribrid™ Mass  
243 Spectrometer (Lumos) (Thermo Scientific) or Q Exactive™ HF Hybrid Quadrupole-Orbitrap™  
244 Mass Spectrometer (HFX). Peptides were separated in a 75  $\mu$ m  $\times$  50 cm PepMap C18 column  
245 using a 1 h linear gradient from 2 to 30% buffer B in buffer A at a flow rate of 250 nL/min  
246 (~600 bar). Peptides were introduced into the mass spectrometer using a nano Easy Spray  
247 source (Thermo Scientific) at 2000 V. Subsequent isolation and higher energy C-trap  
248 dissociation (HCD) was induced in the 20 most abundant ions per full MS scan with an  
249 accumulation time of 120 ms and an isolation width of 1.2 Da (Lumos), or 1.6 Da (HFX). All  
250 fragmented precursor ions were actively excluded from repeated selection for 30 s. The mass  
251 spectrometry proteomics data have been deposited to the ProteomeXchange Consortium via  
252 the PRIDE (37) partner repository with the data set identifier PXD024053.

253

#### 254 **Mass spectrometry data analysis.**

255 The sequence interpretations of mass spectrometry spectra were performed using a database  
256 containing all bovine UniProt entries combined with entry P01012 for chicken ovalbumin (total  
257 of 41610 entries) and 4084 entries for *Theileria parva* Muguga proteome (38). The spectral  
258 interpretation was performed using *de novo*-assisted database search with PEAKS 10  
259 (Bioinformatics Solutions), in 'no enzyme' mode, with mass tolerances of 5 ppm for precursor  
260 ions and 0.03 Da for fragment ions. The data was further searched against 313 inbuild peptide  
261 modifications.

262

263

#### 264 **Filtering of MS-identified peptides.**

265 Previous to all analyses, the lists of peptides identified were filtered to remove: 1) peptides  
266 presenting post-translational modifications; 2) peptides with a peptide-spectrum matching  
267 score  $-\text{Log}_{10}(P) < 15$ ; 3) any peptides derived from *T. parva* Muguga, including the ones  
268 identified in both bovine and *T. parva* Muguga entries; and 4) peptides that shared a 9-mer  
269 overlap with the CD4 T-cell epitope benchmark.

270

### 271 **Model Training.**

272 All ligand data were filtered to include only peptides containing 13-21 residues, to exclude any  
273 residual potentially co-eluted MHCI peptides. Negative peptides were added as described  
274 earlier (32) by sampling random natural peptides from the bovine proteome (described below).  
275 Models were trained in a 5-fold cross-validation manner with partitions constructed from 9-  
276 mer common-motif clustering, ensuring no overlap between test- and training-data. Three  
277 model architectures were used (20, 40, and 60 hidden neurons), each trained with ten random  
278 weight initialization, resulting in an ensemble of 150 networks. Models were evaluated in a  
279 percentile rank fashion, meaning that prediction scores are normalized against a distribution of  
280 prediction scores from random natural peptides. Rank scores are more interpretable than raw  
281 prediction scores and allow for fairer comparison across alleles.

282

283 Two models were trained in this project, both using the NNAlign\_MA machine learning  
284 framework (30). The first model (BoLA) was trained on the novel BoLA SA and MA EL data  
285 combined with the BA data from NetMHCIIpan-4.0 with an added set of BoLA BA data  
286 (roughly 250 measurements for each BoLA-DR molecules incorporating the three different  
287 BoLA-DRB3 alleles - generated *in house*). For the second model (All Data), the BoLA EL data  
288 were combined with all the EL data from the NetMHCIIpan-4.0 data set (human and murine

289 EL data) and the same BA data as the BoLA model. The BoLA and All Data models share  
290 partitions.

291

292 Explicit encoding of ligand context was leveraged to capture antigen processing signatures, as  
293 previously described (20). Briefly, in context encoding 12 residues of the ligand and antigen  
294 are fed as input to the model, 6 are from the N-terminal region of the ligand (3 residues  
295 upstream of the ligand in the antigen and 3 N-terminal Peptide Flanking Regions (PFRs)), and  
296 6 are from the C-terminal region (3 C-terminal PFRs and 3 downstream of the ligand).

297

298 Peptide lists resulting from BoLA-DR eluted ligand data are by nature only positive examples  
299 of ligands that interact with MHCII (excepting co-eluting peptide noise from assay). To train a  
300 peptide-MHCII interaction model, the training data must include examples of non-interacting  
301 peptides sampled from the same background as positive data. To achieve this, peptides (and  
302 their context, see above) were randomly sampled from the bovine proteome. Random negative  
303 peptides were made to follow a uniform length distribution of 13-21 residues, sampling for  
304 each length five times the number of peptides in the most commonly observed ligand length  
305 for a dataset. Negatives were sampled independently for each bovine dataset with a uniform  
306 length distribution so the model can learn the length distribution of ligands (27, 39).

307

308

309

310

## 311 **Results**

### 312 **Analysis of the BoLA-DRB3 repertoire in an experimental cohort of Brazilian** 313 **Holstein-Friesians.**

314 The IPD-MHC database includes over 300 BoLA-DRB3 alleles, of which only a small subset  
315 could be included in this study. To identify the alleles that would be most relevant to ongoing  
316 experiments, a novel high-throughput MiSeq BoLA-DRB3 sequencing approach (Vasoya *et*  
317 *al.*, in preparation) was used to examine the frequency of DRB3 alleles in a representative  
318 cohort of 30 Holstein-Friesian animals from the experimental herd at the University of São  
319 Paulo, Brazil. A total of 22 DRB3 alleles were identified, including a novel allele that had not  
320 been previously described (nDRB3.1). Typical of MHC allele distribution in most cattle  
321 populations, there was a small number of dominant alleles, DRB3\*15:01, DRB3\*01:01,  
322 DRB3\*11:01, DRB3\*14:01:01, and DRB3\*12:01, which were present at a frequency of  $\geq 5\%$ ,  
323 whilst the remaining 17 alleles were present at lower frequencies (Figure 1).

324

### 325 **Generation and analysis of MS data for BoLA-DR eluted peptides.**

326 Initial experiments to establish a BoLA-DR elution technique used O11 and 2229 *Theileria*  
327 *annulata* (TA) cell lines which had previously been confirmed to be homozygous for  
328 DRB3\*10:01 and DRB3\*11:01, respectively (Table I). The length distribution of the peptides  
329 obtained from the 2229TA and both replicates (n1 and n2) of O11TA cell lines was bi-modal.  
330 One peak, centred around 14-15mers was the size anticipated for MHCII ligands; the second  
331 peak, centred around 8-10mer peptides, was more consistent with the length distribution of  
332 MHCI ligands (Figure 2A), and it was speculated that this represented a substantial level of co-  
333 purification of BoLA-I molecules during BoLA-DR immunoprecipitation. To investigate this,  
334 NetMHCpan-4.1 (31) was used to predict the binding potential of all 8-13-mer peptides in each  
335 of the MS data sets for each of the BoLA-I molecules expressed in the given cell line (Table

336 I). The sequence logos of these peptide sets (Supplementary Figure 2) showed remarkable  
337 similarity to the motifs previously described for the BoLA-I alleles in these haplotypes (30)  
338 and between 56.8-70.9% of the 8-13-mer peptides in each sample were predicted to be BoLA-  
339 I binders (defined using a binding threshold of 5% rank). This corroborated the hypothesis that  
340 the majority of these peptides originated from co-precipitated BoLA-I ligands and their  
341 removal resulted in a substantial diminution of the 8-10mer peak (Figure 2B).

342 To address the observed co-enrichment of pBoLA-I in pBoLA-DR immunoprecipitations, it  
343 was decided to apply a sequential immunoprecipitation protocol, starting with pBoLA-I  
344 complex depletion using an anti-BoLA-I monoclonal antibody (IL-A88), followed by pBoLA-  
345 DR precipitation. This two-step protocol was applied to samples from a series of seven *T.*  
346 *parva*-infected cell lines (Table I) which expressed a range of DRB3 alleles present in our  
347 experimental cohort (\*11:01, \*10:01, \*1501, \*1201) or which were of interest because of  
348 ongoing *T. parva* CD4 T cell epitope identification studies that included these alleles (\*16:01  
349 and \*20:01). The total numbers of peptides identified in these samples ranged between 1280  
350 and 8335 (Table I), and the distribution of the peptide lengths is shown in Figure 2C. The  
351 results in this figure show a substantially lower representation of 8-10mer peptides, indicating  
352 successful reduction but not complete depletion of BoLA-I eluted peptides (Figure 2C).

353 Analysis of the binding potential of the peptides in the 8-10mer peak confirmed that the  
354 majority were, in fact, still BoLA-I binders (Table I and Figure 2D); indicating that although  
355 the preliminary BoLA-I depletion had a profound effect on reducing peptides from co-eluted  
356 pBoLA-I, it did not eliminate them completely. Removal of predicted MHCI binders from the  
357 datasets (ranging in frequency from 0.9-8.9%, Table I) effectively abolished the 8-10mer peak  
358 (Figure 2D), establishing that i) combined BoLA-I depletion by prior immunoprecipitation and  
359 bioinformatic removal of predicted MHCI-binders provided the optimal results and ii)  
360 consistent with other MHCII molecules, BoLA-DRB3 molecules have a preference for binding



361 peptides of length 13-21 amino acids (after the combined filtering, 80.7% of the peptides fall  
362 in this length range).

363

364 **Motif deconvolution and prediction model generation from MS data sets of**  
365 **BoLA-DR eluted ligands.**

366 Using the MS BoLA-DR EL data sets, alternative models for BoLA-DRB3 motif  
367 deconvolution were assessed and a prediction model for BoLA-DRB3 ligands was developed.  
368 Details for the model training and model parameters are described in the materials and methods.  
369 In short, bovine ligand data was filtered only to include peptides of 13-21 residues and were  
370 used as positive data points, with negative data points added as previously described (32). Two  
371 models were trained: a 'BoLA' model using the novel BoLA-DR elution data combined with  
372 the BA (binding affinity) data from NetMHCIIpan-4.0 and a set of BA data covering three  
373 different BoLA-DRB3 alleles; and an 'All Data' model, which includes the BA and EL data of  
374 the BoLA model with added murine and human EL data from the NetMHCIIpan-4.0 data set.  
375 Both models were trained with and without assessing the 'context' of the peptide within the  
376 parent protein (MAC- and MA-models, respectively). Here, ligand context refers to including  
377 residues near the ligand termini, inside and outside the ligand, to capture signals of antigen  
378 processing. Further details on data partitioning, model training and context definition are  
379 provided in materials and methods.

380

381 The results of the cross-validation evaluation measured in terms of the AUC are shown in  
382 Figure 3 and show clear differences in the performance of the models used. Firstly, for both  
383 the 'BoLA' and the 'All Data' models, every cell line data set displayed a higher AUC for the  
384 MAC-model than the MA-Model (p-value: 0.00097 in a binomial test counting number of cell  
385 lines with higher AUC for MAC-models versus MA-models). This agrees with earlier studies

386 for the human and mouse MHCII system (20, 32, 40), showing the value of incorporating  
387 encoding context into the prediction models. Secondly, the 'BoLA' MAC-model has  
388 significantly higher median AUC compared to the 'All Data' MAC-Model (p-value: 0.00195  
389 in a binomial test counting cell lines where 'BoLA' MAC-model has higher AUC compared to  
390 'All Data' MAC-model, excluding ties), indicating that inclusion of the human and murine  
391 training data had no benefit in the generation of a model for BoLA-DR binding prediction. This  
392 comparative evaluation clearly demonstrated the 'BoLA-MAC' model exhibited the best  
393 performance and so was selected for subsequent use.

394

395 Examples of BoLA-DRB3 allele motif deconvolution from EL data-sets as performed by the  
396 BoLA-MAC model are shown in Figure 4. The motif deconvolution results for each sample  
397 included in this study are displayed in Supplementary Figure 3A, and the motifs for each of the  
398 seven BoLA-DRB3 alleles covered by the EL data (combining the data from all samples) are  
399 shown in Supplementary Figure 3B. As can be seen in Figure 4, the deconvolution results in  
400 well-defined motifs, with the anticipated preference for residues at positions 1, 4, 6 and 9 of  
401 the binding core and limited exclusion of non-conforming peptides (average of 8.6% of ligands  
402 assigned as contaminants in samples included in Figure 4). The data presented here also shows  
403 the ability of the deconvolution to discriminate the motifs of both BoLA-DRB3 alleles in  
404 heterozygous samples (495TP and 2123TP) as well as the consistency in the motifs for the  
405 same BoLA-DRB3 molecule obtained from different EL data-sets (e.g. BoLA-DRB3\*10:01 in  
406 495TP and 5072TP). These observations are consistent across all of the samples included in  
407 this study, with non-conforming (trash) peptides constituting only ~12.5%, a high average  
408 Pearson correlation between motifs for the same BoLA-DRB3 molecule (0.92 for BoLA-  
409 DRB3\*10:01 and 0.908 for BoLA-DRB3\*11:01, Supplementary Figure 4A), and a very high  
410 specificity being demonstrated for individual motifs (PPV values in the range 0.751-0.868,

411 across the different deconvolutions, Supplementary Figure 4B). As such, the data confirms that  
412 the BoLA-MAC model permitted the generation of high resolution and reproducible BoLA-  
413 DRB3 binding motifs from EL data. This model, renamed as NetBoLAIpan, has been made  
414 publicly available at [www.cbs.dtu.dk/services/NetBoLAIpan](http://www.cbs.dtu.dk/services/NetBoLAIpan).

415

416 **NetBoLAIpan can be used to predict BoLA-DRB3 presented peptides derived**  
417 **from exogenous proteins.**

418 To extend our studies on the utility of the NetBoLAIpan method developed above, the model's  
419 ability to predict which peptides would be presented by BoLA-DR molecules from an  
420 exogenous protein was examined. Here, both PBMC (BoLA-DRB3\*01:01 and \*11:01) and the  
421 O11TA\_n2 cell line (BoLA-DRB3\*10:01) described above were pulsed with soluble  
422 ovalbumin (OVA, see materials and methods and Supplementary Figure 1 for details) before  
423 performing pBoLA-DR elution. Only one OVA-derived peptide ("SSANLSGISSAESLK")  
424 was identified in the O11TA sample, which demonstrated very poor predicted binding to  
425 BoLA-DRB3\*10:01 with a predicted percentile rank value of 29.2%, strongly suggesting this  
426 peptide to be a contaminant co-purified during the BoLA IP enrichment and hence not a  
427 genuine BoLA presented peptide. In contrast, seven OVA-derived peptides were identified in  
428 the PBMC sample. Mapping the seven peptides onto the OVA protein sequence (Figure 5 -  
429 Inserted panel) shows that all the peptides clustered around the 9-mer core "INKVVRFDK",  
430 located at OVA<sub>54-62</sub>, with a common motif IxxVxRxxK – matching the motif described in  
431 Supplementary Figure 3B for BoLA-DRB3\*01:01. Also of interest is that six out of the seven  
432 ligands observed had proline in the C-2 position, which is a common feature in context motifs  
433 (20). The NetBoLAIpan model was applied to predict potential DRB3\*01:01 and  
434 DRB3\*10:01 ligands in the OVA protein sequence. To achieve this, the OVA protein was *in*  
435 *silico* digested into overlapping 13-21-mer peptides, and binding to DRB3\*01:01 and

436 DRB3\*10:01 was predicted for each peptide with predicted ligands identified using a 1% rank  
437 score threshold; this resulted in the identification of 48 predicted ligands covering binding to  
438 both BoLA-DRB3 molecules. The MS identified and *in silico* predicted ligands were then  
439 stacked onto the OVA protein sequence, and a profile was calculated showing the relative  
440 number of measured and predicted ligands mapped to each amino acid position within the  
441 protein. The MS identified and *in silico* predicted ligand profiles demonstrated a striking  
442 concordance, with the MS identified peptides overlapping with the dominant peak of *in silico*  
443 predicted peptides (38 overlapping peptides located at positions 45-71) (Figure 5) (similar data  
444 were obtained using rank threshold values in the range 0.5-2.0%, results not shown), indicating  
445 that NetBoLAIIpan can accurately predict ligands derived from defined proteins that are  
446 experimentally shown by MS to be presented by BoLA-DR.

447

448 **Validation of the BoLA model for BoLA-DRB3 presented CD4 T cell epitope**  
449 **prediction.**

450 Next, the performance of NetBoLAIIpan was validated using a set of 25 BoLA-DR restricted  
451 *T. parva* CD4 T cell epitopes experimentally validated using *T. parva*-specific CD4 T-cell lines  
452 generated from immunised animals in a IFN $\gamma$  secretion T cell assay (Morrison et al., manuscript  
453 in preparation, refer to Supplementary Table 1). Here, NetMHCIIpan-4.0 was included as a  
454 reference model to test the extent to which peptide presentation rules learned from human and  
455 murine data extrapolate to bovine epitopes. Each epitope source protein was *in silico* digested  
456 into peptide strings matching the length of the epitopes, and each peptide was then assigned  
457 the lowest predicted rank score from the set of 13-19-mers whose binding core overlapped with  
458 the peptide string. Next, the epitope's F-rank value was calculated as the percentage of peptides  
459 with a greater prediction score than the epitope. Hence, a perfect prediction has an F-rank value  
460 of 0, and a random prediction presents a value of 50. Comparison of F-rank values obtained by

461 the different models for the set of *T. parva* epitopes (Figure 6), shows that the NetBoLAIIpan  
462 models with or without context achieved equivalent prediction performance both achieving a  
463 median F-rank value of 0.697% and median prediction percentile rank score for the epitopes  
464 of 0.2. In practical terms, these results translate into 12 out of 25 epitopes being ranked as the  
465 top predicted peptide within the given source protein. Both NetBoLAIIpan models achieved  
466 significantly better F-ranks compared to NetMHCIIpan-4.0 (p-values: <0.001 comparing the  
467 two NetBoLAIIpan models to NetMHCIIpan-4.0). The large difference in the performance of  
468 the NetMHCIIpan-4.0 and NetBoLAIIpan models clearly demonstrates the power of  
469 combining BoLA-DR EL data and advanced immunoinformatics to generate novel tools for  
470 characterizing antigen presentation epitope identification in the BoLA-DR system.

471

472

473

474

475

476

477

478

479

480

481

482

483

484 **Discussion**

485 A pre-requisite for the development of next-generation subunit vaccines is the identification of  
486 antigens containing epitopes that can be recognised by B cells, CD8 T cells and CD4 T cells,  
487 as appropriate for the immune response required. Several bioinformatic tools that enable the  
488 prediction of CD4 T cell epitopes in humans have been developed and the recent integration of  
489 large-scale MHC-eluted peptide data have led to a dramatic improvement in their performance  
490 (21, 24, 32). In contrast, there is a lack of equivalent bioinformatics tools designed specifically  
491 for bovine MHCII molecules, and since the currently available tools have not incorporated  
492 bovine MHCII EL data during their development, they perform with limited accuracy when  
493 applied to bovine data (as demonstrated in this study - Figure 6). In previous studies, we have  
494 shown how the use of high-quality EL mass spectrometry data combined with advanced  
495 immunoinformatics and machine-learning techniques can further our understanding of the rules  
496 underlying MHC antigen processing and presentation, allowing the development of improved  
497 prediction methods for MHC ligands and T cell epitopes (30–32). Here, we have extended this  
498 work to cover, for the first time, BoLA-DR molecules.

499 Results from our initial experiments indicated that the peptides isolated following pBoLA-DR  
500 immunoprecipitation were heavily contaminated with co-eluted pBoLA-I-presented peptides.  
501 This phenomenon has been reported previously in other studies using equivalent protocols for  
502 immunoprecipitation of MHCII molecules from human cell lines and has been hypothesised to  
503 reflect that the protocol for lysing the cells results in the immunoprecipitation of membrane  
504 fractions, which contain both MHCI and MHCII molecules (41, 42). In this study neither prior  
505 depletion of pBoLA-I (by immunoprecipitation) nor bioinformatic prediction and removal of  
506 BoLA-I contaminant ligands were completely effective in eliminating the BoLA-I-binding  
507 contamination when applied alone - both left a remnant peak of 8-10-mer peptides. However,  
508 the combined use of these two approaches was successful in removing the 8-10-mer peptide  
509 peak, resulting in 13-21-mer dominated profiles characteristic of MHCII presented peptides.

510 On this basis, we would propose that future studies for BoLA-II immuno-peptidomics should  
511 routinely make use of both preliminary depletion of pBoLA-I complexes by use of an initial  
512 pBoLA-I immunoprecipitation step (consistent with recently developed approaches for human  
513 MHCII immuno-peptidomic studies (21, 43)), and *in silico* immunoinformatic BoLA-I peptide-  
514 binding depletion using currently available prediction methods (31, 33) (or if working with cell  
515 lines expressing alternative BoLA-I haplotypes by generating BoLA-I peptide-binding motifs  
516 by subjecting the product of the preliminary pBoLA-I immunoprecipitation to elution, mass-  
517 spectrometric analysis and subsequent motif deconvolution).

518 In this study, we compared two models for developing the BoLA-DR prediction algorithm. The  
519 first of these was trained using EL data only from BoLA-DR, whilst the second was trained on  
520 the same data augmented by an exhaustive human (HLA) and murine (H-2) MHCII-eluted  
521 peptide dataset (both models also incorporated human, murine and a small amount of bovine  
522 BA data). A cross-validation evaluation demonstrated that the former model had superior  
523 performance, suggesting that integration of cross-species EL datasets was not beneficial to the  
524 accuracy of the results generated by this model. However, this evaluation was restricted to the  
525 limited set of BoLA-DRB3 alleles covered by the EL data generated in the current study, and  
526 it remains to be seen whether a model integrating cross-species EL data would allow improved  
527 prediction when extrapolated to data generated from samples expressing other BoLA-DRB3  
528 alleles. As over 300 BoLA-DRB3 alleles have been described at present, further evaluation of  
529 how best to incorporate inter- and intra-species data to improve the algorithm's performance is  
530 warranted as it will not be feasible for BoLA-DR EL data to be generated for more than a subset  
531 of these alleles. The seven BoLA-DRB3 alleles included in this study were selected  
532 predominantly based on their frequency in the experimental herd of Holstein-Friesian cattle at  
533 the University of São Paulo (USP) (in combination with the availability of DRB3-genotyped  
534 TA/TP cell lines and validated BoLA-DRB3 presented epitope data). The cumulative total

535 frequency of these seven alleles in the samples of animals from the USP herd was ~48% and  
536 retrospective analysis of the University of Edinburgh herd shows that these alleles have an even  
537 higher representation (~67.9%). This is broadly in line with the frequencies observed in  
538 Holstein-Friesian herds across South America and other parts of the world (51.2-73%) (18).  
539 Analysis of the BoLA-DRB3 molecules in Holstein-Friesian animals is attractive for several  
540 reasons: i) due to the high levels of inbreeding, characterisation of a small number of DRB3  
541 alleles will allow comprehensive coverage of the breed (e.g. inclusion of another five DRB3  
542 alleles would give 77-98% coverage of Holstein-Friesian populations (18) and ii) as high-value  
543 dairy animals there is great interest in introducing Holstein-Friesians into low-income countries  
544 (frequently tropical) as part of the process of increasing agricultural productivity and food  
545 security; a major limitation to this process is the Holstein-Friesian susceptibility to many of the  
546 pathogens prevalent in regions of the world. Consequently, there is a particular interest in  
547 finding interventions, such as vaccination, that can be used to protect Holstein-Friesian animals  
548 in tropical environments.

549 A critical and general issue for rational vaccine development is the identification of relevant  
550 antigens. Approaches dependent on conventional antigen-screening techniques have  
551 limitations, especially when applied to complex pathogens (e.g. eukaryotic pathogens), where  
552 the size of the proteomes makes a comprehensive analysis of the full potential antigen  
553 repertoire prohibitively expensive and laborious. For such pathogens, bioinformatic tools that  
554 can help rationalise antigen screening assays and/or selection are of particular value and have  
555 a significant potential for accelerating vaccine development. A potential approach would be to  
556 use bioinformatics tools to predict which peptides from a candidate antigen would be present  
557 by BoLA molecules when delivered as a vaccine. To directly evaluate this, we examined  
558 NetBoLAIIpan's ability to correctly identify the peptides from ovalbumin that had been pre-  
559 loaded onto cell's then subjected to MHC-elution analysis. A comparison of the set of eluted



560 peptides from a PBMC sample and the *in silico* predicted BoLA-DRB3 binding peptides  
561 demonstrated an exceptionally high level of concordance. However, in the experiment  
562 performed with the O11 cell line, only a single OVA peptide was identified. Subsequent  
563 evaluation using NetBoLAIIpan was not able to identify this peptide as a predicted ligand,  
564 suggesting that it was not sourced directly from O11 BoLA-DRB3 molecule DRB3\*10:01, but  
565 rather could represent an OVA degradation peptide product co-purified during the BoLA IP  
566 enrichment. As the *in-silico* analysis further identified several DRB3\*10:01 restricted strong  
567 binding peptide in the OVA sequence, the failure to discover OVA ligand in the O11 MS  
568 experiment strongly suggests that the uptake and presentation of OVA protein in this model  
569 was unsuccessful, supporting the idea of the previously identified peptide as a false ligand.  
570 These results illustrate the integral power of combining *in-silico* modelling and MS elution  
571 studies both for the exact stratification of false-positive sequences identified in such IP  
572 experiments due to co-purification, and to confirm the extent of true-positive peptide ligands.  
573 Further, this analysis suggests that the ability of NetBoLAIIpan to accurately model the  
574 peptides derived from an exogenously administered protein could be exploited to provide an  
575 efficient and inexpensive *in silico* preliminary evaluation of the potential immunogenicity of  
576 candidate antigens and so contribute to the rational selection of antigens (44) prior to  
577 undertaking expensive and laborious *in vivo/in vitro* experiments. In particular, such an  
578 analysis could be used to assess the MHC coverage of individual antigens, and thus inform the  
579 construction of optimal vaccine designs. An example of how such *in silico* analysis could be  
580 employed is given in Supplementary Figure 5.

581 During the development of the prediction model, it was clear that the integration of signals  
582 relating to antigen-processing was beneficial. That is, the inclusion of information regarding  
583 the ‘context’ of the peptides (i.e. both the amino acid residues in the protein flanking the  
584 peptides and the amino acids at the termini of the peptide) significantly improved the power of

585 the models for predicting ligands. The NetBoLAIIpan model exhibited an unprecedented high  
586 performance when evaluated using a set of validated BoLA-DRB3 presented epitopes from *T.*  
587 *parva*, achieving a median F-rank score of 0.697% (corresponding to 12 out of 25 of the defined  
588 epitopes being the highest predicted peptides within the source protein). This performance was  
589 significantly higher than the 19.23% achieved by the previously available NetMHCIIPan model  
590 which had not been trained on the BoLA-DRB3 elution peptide data, demonstrating the utility  
591 of generating and incorporating these data sets. In line with earlier work, context did not impart  
592 the same benefit in the task of ranking CD4 epitopes as was found for ligand data. Here, the  
593 context model was found to perform equivalent to the non-context model. These results align  
594 with earlier work using the mouse and human MHC class II systems (20, 32, 40). Interestingly,  
595 however further improvements in epitope prediction could be obtained by ranking antigen  
596 peptides based on the number of binders within overlapping 13-19-mers. This method of  
597 assigning epitope ranks is based on the intuitive assumption that protein regions with multiple  
598 predicted binders have a greater chance of being presented by BoLA-DRB3 molecules. Using  
599 this approach, the median F-rank score was 0.362%, suggesting a non-trivial improvement in  
600 the prediction. However, further benchmarking on larger epitope sets to systematically evaluate  
601 the comparative performance of this methodology is needed before the recommendation that it  
602 is routinely adopted can be made.

603 In conclusion, this study has proven the high value and important synergistic effect of  
604 combining peptide-MHC elution MS data and advanced immunoinformatics to characterize  
605 antigen presentation and perform ligand/epitope identification in the BoLA-DR system.

606

607

608

609

610

611

612

613

614

615

616 **References**

617 1. Neefjes, J., M. L. M. Jongstra, P. Paul, and O. Bakke. 2011. Towards a systems understanding of  
618 MHC class I and MHC class II antigen presentation. *Nat. Rev. Immunol.* 11: 823–836.

619 2. Luckheeram, R. V., R. Zhou, A. D. Verma, and B. Xia. 2012. CD4<sup>+</sup>T cells: differentiation and  
620 functions. *Clin. Dev. Immunol.* 2012: 925135.

621 3. Chicz, R. M., R. G. Urban, W. S. Lane, J. C. Gorga, L. J. Stern, D. A. Vignali, and J. L. Strominger.  
622 1992. Predominant naturally processed peptides bound to HLA-DR1 are derived from MHC-related  
623 molecules and are heterogeneous in size. *Nature* 358: 764–768.

624 4. Schueler-Furman, O., Y. Altuvia, and H. Margalit. 2001. Examination of possible structural  
625 constraints of MHC-binding peptides by assessment of their native structure within their source  
626 proteins. *Proteins* 45: 47–54.

627 5. Jones, E. Y., L. Fugger, J. L. Strominger, and C. Siebold. 2006. MHC class II proteins and disease: a  
628 structural perspective. *Nat. Rev. Immunol.* 6: 271–282.

629 6. Beck, S., and J. Trowsdale. 2000. The human major histocompatibility complex: lessons from the  
630 DNA sequence. *Annu. Rev. Genomics Hum. Genet.* 1: 117–137.

631 7. Andersson, L., and L. Rask. 1988. Characterization of the MHC class II region in cattle. The number  
632 of DQ genes varies between haplotypes. *Immunogenetics* 27: 110–120.

633 8. Zhou, H., J. G. H. Hickford, Q. Fang, and S. O. Byun. 2007. Short communication: Identification of  
634 allelic variation at the bovine DRA locus by polymerase chain reaction-single strand conformational  
635 polymorphism. *J. Dairy Sci.* 90: 1943–1946.

636 9. Burke, M. G., R. T. Stone, and N. E. Muggli-Cockett. 1991. Nucleotide sequence and northern  
637 analysis of a bovine major histocompatibility class II DR beta-like cDNA. *Anim. Genet.* 22: 343–352.

638 10. Behl, J. D., N. K. Verma, N. Tyagi, P. Mishra, R. Behl, and B. K. Joshi. 2012. The major  
639 histocompatibility complex in bovines: a review. *ISRN Vet. Sci.* 2012: 872710.

640 11. Takeshima, S., M. Ikegami, M. Morita, Y. Nakai, and Y. Aida. 2001. Identification of new cattle  
641 BoLA-DRB3 alleles by sequence-based typing. *Immunogenetics* 53: 74–81.

642 12. Miltiadou, D., A. S. Law, and G. C. Russell. 2003. Establishment of a sequence-based typing  
643 system for BoLA-DRB3 exon 2. *Tissue Antigens* 62: 55–65.

644 13. Baxter, R., N. Hastings, A. Law, and E. J. Glass. 2008. A rapid and robust sequence-based  
645 genotyping method for BoLA-DRB3 alleles in large numbers of heterozygous cattle. *Anim. Genet.* 39:  
646 561–563.

647 14. Takeshima, S.-N., Y. Matsumoto, and Y. Aida. 2009. Short communication: Establishment of a  
648 new polymerase chain reaction-sequence-based typing method for genotyping cattle major  
649 histocompatibility complex class II DRB3. *J. Dairy Sci.* 92: 2965–2970.

650 15. Miyasaka, T., S.-N. Takeshima, H. Sentsui, and Y. Aida. 2012. Identification and diversity of bovine  
651 major histocompatibility complex class II haplotypes in Japanese Black and Holstein cattle in Japan. *J.*  
652 *Dairy Sci.* 95: 420–431.

653 16. Giovambattista, G., S. Takeshima, M. V. Ripoli, Y. Matsumoto, L. A. A. Franco, H. Saito, M.  
654 Onuma, and Y. Aida. 2013. Characterization of bovine MHC DRB3 diversity in Latin American Creole  
655 cattle breeds. *Gene* 519: 150–158.

656 17. Takeshima, S. N., T. Miyasaka, M. Polat, M. Kikuya, Y. Matsumoto, C. N. Mingala, M. A.  
657 Villanueva, A. J. Salces, M. Onuma, and Y. Aida. 2014. The great diversity of major histocompatibility  
658 complex class II genes in Philippine native cattle. *Meta Gene* 2: 176–190.

659 18. Takeshima, S.-N., G. Giovambattista, N. Okimoto, Y. Matsumoto, A. Rogberg-Muñoz, T. J. Acosta,  
660 M. Onuma, and Y. Aida. 2015. Characterization of bovine MHC class II DRB3 diversity in South  
661 American Holstein cattle populations. *Tissue Antigens* 86: 419–430.

662 19. Takeshima, S.-N., C. Corbi-Botto, G. Giovambattista, and Y. Aida. 2018. Genetic diversity of BoLA-  
663 DRB3 in South American Zebu cattle populations. *BMC Genet.* 19: 33.

664 20. Barra, C., B. Alvarez, S. Paul, A. Sette, B. Peters, M. Andreatta, S. Buus, and M. Nielsen. 2018.

665 Footprints of antigen processing boost MHC class II natural ligand predictions. *Genome Med.* 10: 84.  
666 21. Racle, J., J. Michaux, G. A. Rockinger, M. Arnaud, S. Bobisse, C. Chong, P. Guillaume, G. Coukos, A.  
667 Harari, C. Jandus, M. Bassani-Sternberg, and D. Gfeller. 2019. Robust prediction of HLA class II  
668 epitopes by deep motif deconvolution of immunopeptidomes. *Nat. Biotechnol.* 37: 1283–1286.  
669 22. Bassani-Sternberg, M., and D. Gfeller. 2016. Unsupervised HLA Peptidome Deconvolution  
670 Improves Ligand Prediction Accuracy and Predicts Cooperative Effects in Peptide-HLA Interactions. *J.*  
671 *Immunol. Baltim. Md 1950* 197: 2492–2499.  
672 23. Abelin, J. G., D. B. Keskin, S. Sarkizova, C. R. Hartigan, W. Zhang, J. Sidney, J. Stevens, W. Lane, G.  
673 L. Zhang, T. M. Eisenhaure, K. R. Clauser, N. Hacohen, M. S. Rooney, S. A. Carr, and C. J. Wu. 2017.  
674 Mass Spectrometry Profiling of HLA-Associated Peptidomes in Mono-allelic Cells Enables More  
675 Accurate Epitope Prediction. *Immunity* 46: 315–326.  
676 24. Abelin, J. G., D. Harjanto, M. Malloy, P. Suri, T. Colson, S. P. Goulding, A. L. Creech, L. R. Serrano,  
677 G. Nasir, Y. Nasrullah, C. D. McGann, D. Velez, Y. S. Ting, A. Poran, D. A. Rothenberg, S. Chhangawala,  
678 A. Rubinsteyn, J. Hammerbacher, R. B. Gaynor, E. F. Fritsch, J. Greshock, R. C. Oslund, D. Barthelme,  
679 T. A. Addona, C. M. Arieta, and M. S. Rooney. 2019. Defining HLA-II Ligand Processing and Binding  
680 Rules with Mass Spectrometry Enhances Cancer Epitope Prediction. *Immunity* 51: 766-779.e17.  
681 25. Bulik-Sullivan, B., J. Busby, C. D. Palmer, M. J. Davis, T. Murphy, A. Clark, M. Busby, F. Duke, A.  
682 Yang, L. Young, N. C. Ojo, K. Caldwell, J. Abhyankar, T. Boucher, M. G. Hart, V. Makarov, V. T. D.  
683 Montpreville, O. Mercier, T. A. Chan, G. Scagliotti, P. Bironzo, S. Novello, N. Karachaliou, R. Rosell, I.  
684 Anderson, N. Gabrail, J. Hrom, C. Limvarapuss, K. Choquette, A. Spira, R. Rousseau, C. Voong, N. A.  
685 Rizvi, E. Fadel, M. Frattini, K. Jooss, M. Skoberne, J. Francis, and R. Yelensky. 2018. Deep learning  
686 using tumor HLA peptide mass spectrometry datasets improves neoantigen identification. *Nat.*  
687 *Biotechnol.* .  
688 26. Sarkizova, S., S. Klaeger, P. M. Le, L. W. Li, G. Oliveira, H. Keshishian, C. R. Hartigan, W. Zhang, D.  
689 A. Braun, K. L. Ligon, P. Bachireddy, I. K. Zervantonakis, J. M. Rosenbluth, T. Ouspenskaia, T. Law, S.  
690 Justesen, J. Stevens, W. J. Lane, T. Eisenhaure, G. Lan Zhang, K. R. Clauser, N. Hacohen, S. A. Carr, C.  
691 J. Wu, and D. B. Keskin. 2020. A large peptidome dataset improves HLA class I epitope prediction  
692 across most of the human population. *Nat. Biotechnol.* 38: 199–209.  
693 27. Garde, C., S. H. Ramarathinam, E. C. Jappe, M. Nielsen, J. V. Kringelum, T. Trolle, and A. W.  
694 Purcell. 2019. Improved peptide-MHC class II interaction prediction through integration of eluted  
695 ligand and peptide affinity data. *Immunogenetics* 71: 445–454.  
696 28. Prilliman, K., M. Lindsey, Y. Zuo, K. W. Jackson, Y. Zhang, and W. Hildebrand. 1997. Large-scale  
697 production of class I bound peptides: assigning a signature to HLA-B\*1501. *Immunogenetics* 45: 379–  
698 385.  
699 29. Andreatta, M., B. Alvarez, and M. Nielsen. 2017. GibbsCluster: unsupervised clustering and  
700 alignment of peptide sequences. *Nucleic Acids Res.* 45: W458–W463.  
701 30. Alvarez, B., B. Reynisson, C. Barra, S. Buus, N. Ternette, T. Connelley, M. Andreatta, and M.  
702 Nielsen. 2019. NNAlign\_MA; MHC peptidome deconvolution for accurate MHC binding motif  
703 characterization and improved T cell epitope predictions. *Mol. Cell. Proteomics MCP* .  
704 31. Reynisson, B., B. Alvarez, S. Paul, B. Peters, and M. Nielsen. 2020. NetMHCpan-4.1 and  
705 NetMHCIpan-4.0: improved predictions of MHC antigen presentation by concurrent motif  
706 deconvolution and integration of MS MHC eluted ligand data. *Nucleic Acids Res.* 48: W449–W454.  
707 32. Reynisson, B., C. Barra, S. Kaabinejadian, W. H. Hildebrand, B. Peters, and M. Nielsen. 2020.  
708 Improved Prediction of MHC II Antigen Presentation through Integration and Motif Deconvolution of  
709 Mass Spectrometry MHC Eluted Ligand Data. *J. Proteome Res.* 19: 2304–2315.  
710 33. Nielsen, M., T. Connelley, and N. Ternette. 2018. Improved Prediction of Bovine Leucocyte  
711 Antigens (BoLA) Presented Ligands by Use of Mass-Spectrometry-Determined Ligand and in Vitro  
712 Binding Data. *J. Proteome Res.* 17: 559–567.  
713 34. Goddeeris, B. M., and W. I. Morrison. 1988. Techniques for the generation, cloning, and  
714 characterization of bovine cytotoxic T cells specific for the protozoan *Theileria parva*. *J. Tissue Cult.*  
715 *Methods* 11: 101–110.  
716 35. DeMartini, J. C., N. D. MacHugh, J. Naessens, and A. J. Teale. 1993. Differential in vitro and in vivo

717 expression of MHC class II antigens in bovine lymphocytes infected by *Theileria parva*. *Vet. Immunol.*  
718 *Immunopathol.* 35: 253–273.

719 36. Purcell, A. W., S. H. Ramarathinam, and N. Ternette. 2019. Mass spectrometry-based  
720 identification of MHC-bound peptides for immunopeptidomics. *Nat. Protoc.* 14: 1687–1707.

721 37. Perez-Riverol, Y., A. Csordas, J. Bai, M. Bernal-Llinares, S. Hewapathirana, D. J. Kundu, A. Inuganti,  
722 J. Griss, G. Mayer, M. Eisenacher, E. Pérez, J. Uszkoreit, J. Pfeuffer, T. Sachsenberg, S. Yilmaz, S.  
723 Tiwary, J. Cox, E. Audain, M. Walzer, A. F. Jarnuczak, T. Ternent, A. Brazma, and J. A. Vizcaíno. 2019.  
724 The PRIDE database and related tools and resources in 2019: improving support for quantification  
725 data. *Nucleic Acids Res.* 47: D442–D450.

726 38. Tretina, K., R. Pelle, J. Orvis, H. T. Gotia, O. O. Ifeonu, P. Kumari, N. C. Palmateer, S. B. A. Iqbal, L.  
727 M. Fry, V. M. Nene, C. A. Daubenberger, R. P. Bishop, and J. C. Silva. 2020. Re-annotation of the  
728 *Theileria parva* genome refines 53% of the proteome and uncovers essential components of N-  
729 glycosylation, a conserved pathway in many organisms. *BMC Genomics* 21: 279.

730 39. Jurtz, V., S. Paul, M. Andreatta, P. Marcatili, B. Peters, and M. Nielsen. 2017. NetMHCpan-4.0:  
731 Improved Peptide-MHC Class I Interaction Predictions Integrating Eluted Ligand and Peptide Binding  
732 Affinity Data. *J. Immunol. Baltim. Md 1950* 199: 3360–3368.

733 40. Paul, S., E. Karosiene, S. K. Dhanda, V. Jurtz, L. Edwards, M. Nielsen, A. Sette, and B. Peters. 2018.  
734 Determination of a Predictive Cleavage Motif for Eluted Major Histocompatibility Complex Class II  
735 Ligands. *Front. Immunol.* 9: 1795.

736 41. Partridge, T., A. Nicastrì, A. E. Kliszczak, L.-M. Yindom, B. M. Kessler, N. Ternette, and P. Borrow.  
737 2018. Discrimination Between Human Leukocyte Antigen Class I-Bound and Co-Purified HIV-Derived  
738 Peptides in Immunopeptidomics Workflows. *Front. Immunol.* 9: 912.

739 42. Bettencourt, P., J. Müller, A. Nicastrì, D. Cantillon, M. Madhavan, P. D. Charles, C. B. Fotso, R.  
740 Wittenberg, N. Bull, N. Pinpathomrat, S. J. Waddell, E. Stylianou, A. V. S. Hill, N. Ternette, and H.  
741 McShane. 2020. Identification of antigens presented by MHC for vaccines against tuberculosis. *NPJ*  
742 *Vaccines* 5: 2.

743 43. Chong, C., F. Marino, H. Pak, J. Racle, R. T. Daniel, M. Müller, D. Gfeller, G. Coukos, and M.  
744 Bassani-Sternberg. 2018. High-throughput and Sensitive Immunopeptidomics Platform Reveals  
745 Profound Interferon-γ-Mediated Remodeling of the Human Leukocyte Antigen (HLA) Ligandome. *Mol.*  
746 *Cell. Proteomics MCP* 17: 533–548.

747 44. Barra, C., C. Ackaert, B. Reynisson, J. Schockaert, L. E. Jessen, M. Watson, A. Jang, S. Comtois-  
748 Marotte, J.-P. Goulet, S. Pattijn, E. Paramithiotis, and M. Nielsen. 2020. Immunopeptidomic Data  
749 Integration to Artificial Neural Networks Enhances Protein-Drug Immunogenicity Prediction. *Front.*  
750 *Immunol.* 11: 1304.

751 45. Grisi, L., R. C. Leite, J. R. de S. Martins, A. T. M. de Barros, R. Andreotti, P. H. D. Cançado, A. A. P.  
752 de León, J. B. Pereira, and H. S. Villela. 2014. Reassessment of the potential economic impact of  
753 cattle parasites in Brazil. *Rev. Bras. Parasitol. Vet. Braz. J. Vet. Parasitol. Orgao Of. Col. Bras.*  
754 *Parasitol. Vet.* 23: 150–156.

755 46. Garcia, G. R., J. M. Chaves Ribeiro, S. R. Maruyama, L. G. Gardinassi, K. Nelson, B. R. Ferreira, T. G.  
756 Andrade, and I. K. F. de Miranda Santos. 2020. A transcriptome and proteome of the tick  
757 *Rhipicephalus microplus* shaped by the genetic composition of its hosts and developmental stage.  
758 *Sci. Rep.* 10: 12857.

759 47. Manzano-Román, R., V. Díaz-Martín, A. Oleaga, P. Obolo-Mvoulouga, and R. Pérez-Sánchez.  
760 2016. TSGP4 from *Ornithodoros moubata*: molecular cloning, phylogenetic analysis and vaccine  
761 efficacy of a new member of the lipocalin clade of cysteinyl leukotriene scavengers. *Vet. Parasitol.*  
762 227: 130–137.

763 48. Andreotti, R., R. C. Cunha, M. A. Soares, F. D. Guerrero, F. P. L. Leite, and A. A. P. de León. 2012.  
764 Protective immunity against tick infestation in cattle vaccinated with recombinant trypsin inhibitor  
765 of *Rhipicephalus microplus*. *Vaccine* 30: 6678–6685.

766 49. Labuda, M., A. R. Trimnell, M. Licková, M. Kazimírová, G. M. Davies, O. Lissina, R. S. Hails, and P.  
767 A. Nuttall. 2006. An antivector vaccine protects against a lethal vector-borne pathogen. *PLoS Pathog.*  
768 2: e27.

769 50. Harnnoi, T., S. Watchabunsook, T. Sakaguchi, X. Xuan, and K. Fujisaki. 2006. Characterization of  
770 Haemaphysalis longicornis recombinant cement-like antigens and preliminary study of their  
771 vaccination effects. *J. Vet. Med. Sci.* 68: 1289–1295.

772 51. Canales, M., J. M. P. de la Lastra, V. Naranjo, A. M. Nijhof, M. Hope, F. Jongejan, and J. de la  
773 Fuente. 2008. Expression of recombinant Rhipicephalus (Boophilus) microplus, R. annulatus and R.  
774 decoloratus Bm86 orthologs as secreted proteins in Pichia pastoris. *BMC Biotechnol.* 8: 14.

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

## 792 **Figures**

793 **Figure 1 - Frequencies of BoLA-DRB3 alleles detected by a MiSeq genotyping approach in a subset of the**  
794 **experimental Holstein-Friesian cattle herd at the University of Sao Paulo (n=30).** The frequency data is shown  
795 as a Pareto plot with the frequency of individual alleles displayed on the left vertical axis and the cumulative  
796 frequencies of the DRB3 alleles shown on the right vertical axis. Allele nDRB3.1 was a novel sequence.

797

798 **Figure 2 - Length distribution of BoLA-DR eluted peptides.** Kernel density estimates comparing length  
799 distributions of BoLA-DR eluted peptides using different strategies for removal of BoLA-I eluted contaminants:  
800 (A) Direct pBoLA-DR elution; (B) Direct pBoLA-DR elution with subsequent removal of BoLA-I binders as  
801 predicted by NetMHCpan-4.1; (C) Initial immunoprecipitation to deplete pBoLA-I complexes. (D) Same as for  
802 panel (C) but with subsequent removal of BoLA-I binders as predicted by NetMHCpan-4.1. Due to failed pBoLA-  
803 I depletion sample 2229TP is not represented in this figure.

804

805 **Figure 3 - Cross-Validation evaluation of bovine EL data.** Models were evaluated on the BoLA-DR ligand  
806 data in a cross-validation manner. The boxplot shows the AUC per cell line sample for the BoLA and All Data  
807 models with and without context encoding (MAC-Model and MA-Model, respectively). Each point in the figure  
808 represents data from a single sample. Of note, the outlier sample with a cross-validated AUC performance below  
809 0.90 for the BoLA-MAC model was 2229TA; this sample had 27% ligands assigned as contaminants  
810 (Supplementary Figure 3A) causing the decrease in the observed AUC.

811

812 **Figure 4 - Examples of deconvoluted motifs derived from EL BoLA-DR datasets.** From each cell line defined  
813 as being heterozygous for DRB3, two peptide-binding motifs were derived. Where cell lines express the same  
814 DRB3 allele, consistent motifs were identified (e.g., both 2123TP and 495TP express DRB3\*11:01 and show a  
815 similar peptide-binding motif). Motifs were generated from ligands with a rank score of <20 for the context-  
816 model. Ligands with a predicted rank >20 are assigned to the Trash cluster. Logos show alignments of predicted  
817 peptide binding cores where numbers in parentheses represent the number of peptides.

818

819 **Figure 5 - Profiles of predicted and measured OVA ligands in the PBMC sample. (Main Figure)** The gray  
820 shaded area shows the relative number of measured EL ligands in the PBMC sample overlapping each position in  
821 the OVA sequence. The dotted line represents the mapping of 13-21-mers from the OVA sequence predicted with



822 a rank score < 1% for the BoLA-DRs expressed in the PBMC sample; the peaks at positions 6-23, 45-71, 196-  
823 210 and 275-291 represent 5, 38, 1 and 4 predicted BoLA-DR binding peptides, each with median predicted rank  
824 scores of 0.64, 0.45, 0.82, and 0.56, respectively. (**Inserted panel**) Mapping of the seven OVA peptides measured  
825 in the PBMC sample. All but one of the peptides shared a binding core “INKVVRFDK” in positions 54-62 of the  
826 OVA sequence.

827

828 **Figure 6 - Comparison of different BoLA-DR prediction models using validated CD4 T cell epitopes.**

829 Distribution of percentage F-rank performance values for defined BoLA-DR presented *T. parva* epitopes using  
830 the NetBoLAIIpan and NetMHCIIpan-4.0 models with (Context) and without context (No Context). Prediction  
831 scores were assigned to each overlapping epitope length-matched peptide in the epitope source protein as  
832 described in the text. The y-axis is shown in log-scale and F-rank values below 0.1 are presented as 0.1005 to  
833 avoid non-defined values.