

# Effective assessment of CD4<sup>+</sup> T cell Immunodominance patterns: impact of antigen processing and HLA restriction

Miguel Álvaro-Benito<sup>1,2,†,\*</sup>, Esam T Abualrous<sup>1,3,4,†</sup>, Holger Lingel<sup>5</sup>, Stefan Meltendorf<sup>5</sup>, Jakob Holzapfel<sup>1</sup>, Jana Sticht<sup>1</sup>, Benno Kuroпка<sup>6</sup>, Cecilia Clementi<sup>7</sup>, Frank Kuppler<sup>1</sup>, Monika C Brunner-Weinzierl<sup>5</sup>, Christian Freund<sup>1</sup>.

<sup>1</sup>Laboratory of Protein Biochemistry. Department of Biology, Chemistry and Pharmacy. Freie Universität Berlin, Thielallee 63, 14195 Berlin, Germany

<sup>2</sup>Department of Immunology, Ophthalmology and ENT, Universidad Complutense, School of Medicine and 12 de Octubre Health Research Institute (imas12), Madrid, Spain

<sup>3</sup>Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin, Germany

<sup>4</sup>Department of Physics, Faculty of Science, Ain Shams University, Cairo, Egypt

<sup>5</sup>Department of Experimental Pediatrics, Medical Faculty, Otto-von-Gericke-University, Leipziger Str. 44, 39120 Magdeburg, Germany

<sup>6</sup>Mass Spectrometry Core Facility (BioSupraMol), Freie Universität Berlin, Thielallee 63, 14195 Berlin, Germany

<sup>7</sup>Theoretical & Computational Biophysics, Institute for Physics, Freie Universität Berlin, Arnimallee 12, 14195 Berlin

<sup>†</sup> Contributed equally

\* Correspondence to: [m.alvaro@fu-berlin.de](mailto:m.alvaro@fu-berlin.de)

## Abstract

Identifying T cell epitopes is essential for studying and potentially tuning immune responses to pathogens. The polymorphic nature of major histocompatibility complex of class II (MHCII)-genes, and the complexity of the antigen processing mechanisms hinders the effective prediction of immunodominant patterns in humans, specially at the population level. Here, we combined the output of a reconstituted antigen processing system and of *in silico* prediction tools for SARS-CoV-2 antigens considering a broad-population coverage DRB1\* panel to gain insights on immunodominance patterns. The two methods complement each other, and the resulting model improves upon single positive predictive values (PPV) from each of them to explain known epitopes. This model was used to design a minimalistic peptide pool (59 peptides) matching the performance reported for large overlapping peptide pools (> 500 peptides). Furthermore, almost 70 % of the candidates (23 peptides) selected for a frequent HLA background (DRB1\*03:01/\*07:01) feature immunodominant responses *ex vivo*, validating our platform for accessing T cell epitopes at the population level. The analysis of the impact of processing constraints reveals distinct impact of proteolysis and solvent accessible surface area on epitope selection depending on the antigen. Thus, considering these properties for antigens in question should improve available epitope prediction tools.

**Keywords:** Immunodominance, HLA class II, peptide, CD4<sup>+</sup> T cell, Antigen Processing, *in vitro*

## Introduction

T cell responses to any given antigen are focused on immunogenic peptides restricted by individual's Major Histocompatibility Complex (MHC) molecules. MHC-bound peptides triggering T cell activation are considered T cell epitopes and thus represent the switch for accessing T cell function. Peptides featuring T cell activation are extremely important in basic and clinical immunology research, and when they are recurrently selected by one or several MHC allotypes they are considered immunodominant<sup>1</sup>. Immunodominant T cell epitopes may be more or less easily accessed on inbred and syngeneic animals with defined MHCII haplotypes. However, addressing and rationalizing immunodominance patterns in humans poses major challenges. There are important differences at the individual and population level between human beings exposed to pathogens, and experimental animal models immunized under laboratory conditions. Humans express at least three highly polymorphic sets of MHCII molecules (DRB1, DQ and DP), and it is therefore extremely difficult to assign a clear restriction for peptides that yield positive hits for T cell activation at a large scale. Furthermore, certain peptides feature promiscuous binding, hence they are restricted by more than one MHCII-allotype. Additionally, inter-individual differences on antigen processing mechanisms and distinct T cell Receptor (TcR) genes and repertoires between individuals affect the preferential recognition of peptide-MHCII combinations, hence the observed responses.

Peptide binding to available MHCII, more precisely the kinetic stability of the complex, is regarded as the key feature correlated with immunogenicity<sup>2</sup>. Sequence agreement of antigenic determinants with those of peptides eluted from MHC molecules is typically considered as a proxy for kinetic stability, and thus it is used to predict immunodominance. Importantly, not all antigenic peptides with a predicted or measured binding affinity become immunogenic. Under physiological conditions antigens are proteolytically degraded, mainly in late endosomal compartments where they are loaded onto MHCII molecules in the presence of proteases and HLA-DM<sup>3</sup>. However, besides the canonical processing of antigens in endosomes, work on murine models of infection revealed the relevance of alternative processing pathways to immunodominance<sup>4</sup>. Resistance to proteolytic degradation of antigenic regions, along with their location in folded or unfolded regions amenable to bind to MHCII have been proven to impact immunogenicity. In this context, novel platforms based on artificial intelligence (AI) have sought to implement additional constraints to improve epitope discovery with encouraging, yet suboptimal results.

The interest on accessing immunodominant determinants for understanding and manipulating immune responses to pathogens has motivated different strategies and conceptual frameworks over the last decades. Bottom-up strategies prioritize fundamentally peptide binding to, or presentation by single MHCII molecules. Recent examples of this "MHC-centric" view include investigations in the context of influenza<sup>5</sup>, HIV<sup>6</sup> and more recently, SARS-CoV-2<sup>7,8</sup>. Although these studies are usually limited to a handful of MHCII allotypes, they contribute relevant mechanistic insights for epitope selection. Top-down approaches on the other hand make use of peptide libraries spanning regions with high potential to become immunogenic, or in the best cases entire proteomes. This "MHCII agnostic" framework theoretically provides an unbiased overview of immunogenicity when spanning the entire proteome and has also been considered for SARS-CoV-2 infection<sup>9,10</sup>. However, this approach is limited to small pathogens and may suffer from library design considerations that neglect relevant binders. Nevertheless, "MHC-agnostic" studies can provide insights on immunodominance patterns when combined with elaborated epitope mapping strategies.

The great impact of SARS-CoV-2 on our society led to unprecedented research efforts to score the role of the immune system to fight the virus. Viral control was associated with CD4<sup>+</sup> T cell

function<sup>11</sup>, and disease severity was soon correlated with poor, delayed or inefficient CD4<sup>+</sup> T cell responses<sup>12,13</sup>. Given their relevance, the presence and function of SARS-CoV-2-reactive CD4<sup>+</sup> T cells in healthy, infected, and vaccinated individuals has been thoroughly investigated<sup>9,14–23</sup>. Thus, the average response to SARS-CoV-2 in an individual has been estimated to result from 19 epitopes, as measured with peptide pools of broad viral-antigen and MHCII coverage<sup>9</sup>. These works, together with others dedicated efforts on single MHCII allotypes<sup>8,24</sup> yield an extremely attractive background for scoring the potential antigen processing constraints for immunodominance.

In this work we rationalized the use of a panel of MHCII allotypes with broad population coverage to study SARS-CoV-2 immunodominance at both, allele and population level. We contextualize the candidate epitopes selected by each individual allotype as well as by the entire MHCII panel on the basis of experimental evidence from a reconstituted antigen processing system<sup>25,26</sup>, *in silico* predictions, and information on CD4<sup>+</sup> T cell responses to SARS-CoV-2 available from the IEDB<sup>27</sup>. Integrating the resulting information with known antigen- and MHCII-related features allows us to identify distinct epitope selection patterns for two model antigens. Importantly, the combination of both *in vitro* and *in silico* hits, facilitates defining a limited peptide-pool that triggers T cell responses on a broad panel of MCHII allotypes. Finally, making use of MHCII-tailored peptide pools we could evaluate the extent of immunodominant responses on individuals bearing the same MHCII haplotypes.

## Results

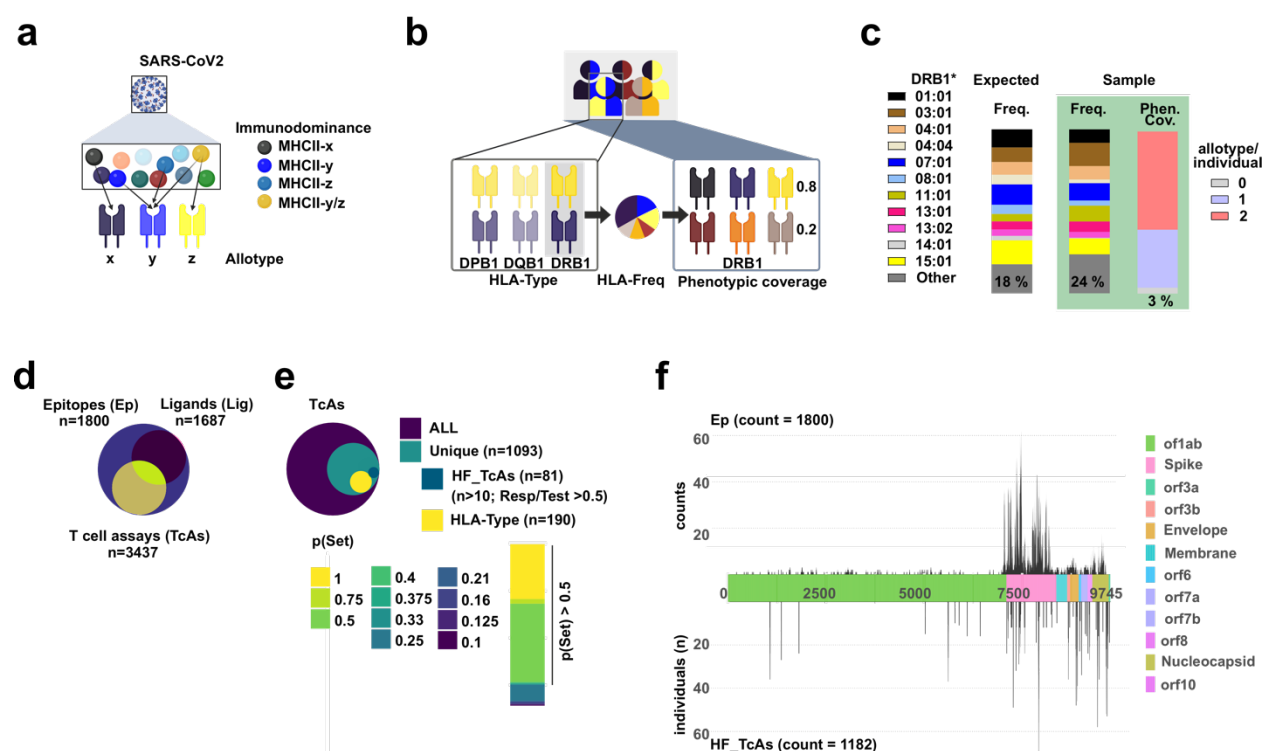
### Contextualization of DRB1 allotypes with a broad population coverage for scoring immunodominance patterns of CD4<sup>+</sup> T cell responses to SARS-CoV-2

Peptides derived from SARS-CoV-2 proteins recurrently selected for their display by one (exclusive restriction) or more than one allotype (promiscuous binder), and giving rise to recurrent T cell responses are considered immunodominant (Figure 1a). We conceived combining information of a reconstituted antigen processing system and *in silico* predictions to query the impact of antigen processing constraints on the selection of immunodominant epitopes for representative allotypes. We reasoned that beyond any allotype-specific immunodominance pattern revealed by our work, an integrative approach considering a representative set of MHCII allotypes will allow us drawing conclusions at the population level, especially when leveraging on the extensive available information.

We focused our efforts on DRB1\* heterodimers as main drivers of MHCII-specific immune responses to viruses<sup>28</sup>. We defined a panel of frequent allotypes with high phenotypic coverage (presence of selected allotypes at the individual level, Figure 1b). Selecting 11 common DRB1\* allotypes, based on available frequencies for European Caucasian populations<sup>29</sup>, we anticipated a phenotypic coverage higher than 90 %. HLA-typing of a pool of 109 donors recruited between November 2019 and February 2020 revealed minor differences between the expected frequencies, and those achieved by our sample. However, the phenotypic coverage achieved by the selected panel was in the expected range, with only 4 out of 109 donors (3 %) not expressing at least one of the allotypes included in our set (Figure 1c, Table S1).

We then evaluated the potential use of available information on immune responses to SARS-CoV-2 to contextualize our prospective findings. Thus, we queried the broadest and most comprehensive repository of human immune responses, the Immune Epitope Data Base (IEDB)<sup>27</sup>. We retrieved all data on CD4<sup>+</sup> T cell responses to SARS-CoV-2 and subset the relevant information for our goal (details in Figure S1a-e). As of September 2022, there were similar total numbers of entries annotated as CD4<sup>+</sup> T cell epitopes (Ep, n=1800) and MHCII Ligands (Lig,

n=1687), and almost double the number for T cell Assays (TcAs, n=3437) but only a handful of multimer-identified epitopes (Tet\_TcAs, n=149) (Figure 1d). There is a clear heterogeneity in terms of size of the studies, of the applied experimental approaches as well as in the source of the antigens tested (Table S2 and Table S3). Information on restriction could be attained from ligand data (Lig entries, with measured binding affinities), peptides from mono-allelic models or defined by display methodologies), or TcA data with associated MHC restrictions (e.g. Multimers, Tet\_TcAs). The subset information is directly applicable for targeted questions, e.g. validating predicted or experimentally defined restrictions. However, we reasoned that the 1093 unique TcAs entries may contain as well relevant information for our work. We therefore evaluated the potential coverage of our panel of 11 DRB1\* allotypes for the studies where typing information was available. This analysis reveals a considerably good representation throughout all datasets (50-100 % phenotypic coverage). Exemplarily, the p(Set) reporting the probability for each TcA record within the subset HLA-Type for being restricted by one of the allotypes of our panel, is higher than 0.5 in 80 % of the cases (Figure 1e).



**Figure 1. Overview of immunodominance from a human population perspective.** **a.** From all potential peptides in the viral proteome of SARS-CoV-2 (shown as spheres) every MHCII allotype considered (MHCII-x-z) would preferentially select a limited pool, which will define the corresponding immunodominance pattern. There will be peptides that become immunogenic for only one allotype (black sphere) whereas others may be restricted by more than one allotype (beige sphere). Peptides being immunodominant for several allotypes feature promiscuous binding to the corresponding MHCs. **b.** Every individual in a population carries specific MHCII haplotypes consisting of DP, DQ and DR molecules (HLA-type, color coded). From all MHCII molecules, DRB1 dimers play a key role in responses to viral pathogens. Based on their frequency it is possible to design a restricted panel of allotypes with a high phenotypic coverage representative for a given population. **c.** A panel of 11 HLA class II molecules (color coded) with an expected cumulative frequency of up to 82 % of European populations (left bar) is conceived to determine the CD4<sup>+</sup> T cell specific responses over a panel of voluntarily recruited donors, and evaluate immunodominant and/or prevalent responses in published data. The coverage of the sampled population reaches up to 76 % (right bar) and a phenotypic coverage of up to 97 % (in 109 donors). **d.** Overview of the number of entries in the Immune Epitope Data Base (IEDB) for all available MHCII-related SARS-CoV-2 epitopes, ligands and T cell Assays (as of September 2022). **e.** The TcAs subset account for 3437 entries, comprising 1093 unique hits that could be further subdivided into different categories. We considered subset entries that provide responses in a high proportion of individuals (HF\_TcAs) and those that have a



broad HLA-typing information associated (HLA-type). In case of HF\_TcAs there is often an inferred restriction assigned to the epitope. Many of those entries are reported to be restricted by some of the DRB1\* allotypes of our panel. However, the heterogeneity on the methods to report restrictions (e.g. predictions or exclusion of alleles from non-reacting donors) led us to omit this information for our analysis. HLA-Type entries on the other hand usually report entries that had been tested in many individuals and only a handful respond to them. In this case the HLA-types of the individuals responding is also provided, and permits the estimate of the probability of phenotypic coverage of the designed panel. Overall, if  $p(\text{Set}) > 0.5$  the individuals tested would bear at least one of the allotypes selected in our panel. If more than one individual was tested, and only one individual would carry any of the HLA-types considered the  $p(\text{Set})$  value ranges between 0.1 and 0.4). f. Summary of the counts of CD4<sup>+</sup> T cell epitopes (upper panel), and number of individuals responding to prevalent immunogenic regions (lower panel) retrieved from the IEDB, and plotted over a scheme of the SARS-CoV-2 proteome (color coded by orfs as stated in the legend). Prevalent or High Frequency (HF\_TcAs) responses are considered those observed in experiments including a minimum of 10 individuals, in which at least 5 individuals responded to the corresponding peptide.

Sampling and/or reporting biases across the viral proteome may impact the contextualization of our planned work. Therefore, we checked for antigen-sampling skewing as a final step towards validating the usage of the IEDB data for our goal. There is a clear over-representation of the three structural proteins: Spike, Nucleocapsid and Membrane when considering entries tagged as Ep (Figure 1f upper part). Such picture is mirrored as well in Lig entries and probably reflects a pronounced interest on accessing the immunogenic determinants of structural proteins (Figure S1c). Despite this bias, we identified a set of peptides tested over a large number of individuals (not selected on the basis of specific allotypes) and triggering prevalent responses. We subset these entries as High Frequency TcAs (HF\_TcAs), if the Response Frequency (RF)<sup>30</sup> was higher than 0.5 and the minimum number of individuals in which it was tested was 10. Interestingly, these entries are spread in most cases over the entire viral proteome (Table S2 and Figure S1e). Given the good phenotypic coverage of our panel we reasoned that these highly immunogenic regions represent relevant targets that may be preferentially restricted by allotypes between those considered by us (Figure 1d lower part). If this is the case, these regions may be as well selected by our approach.

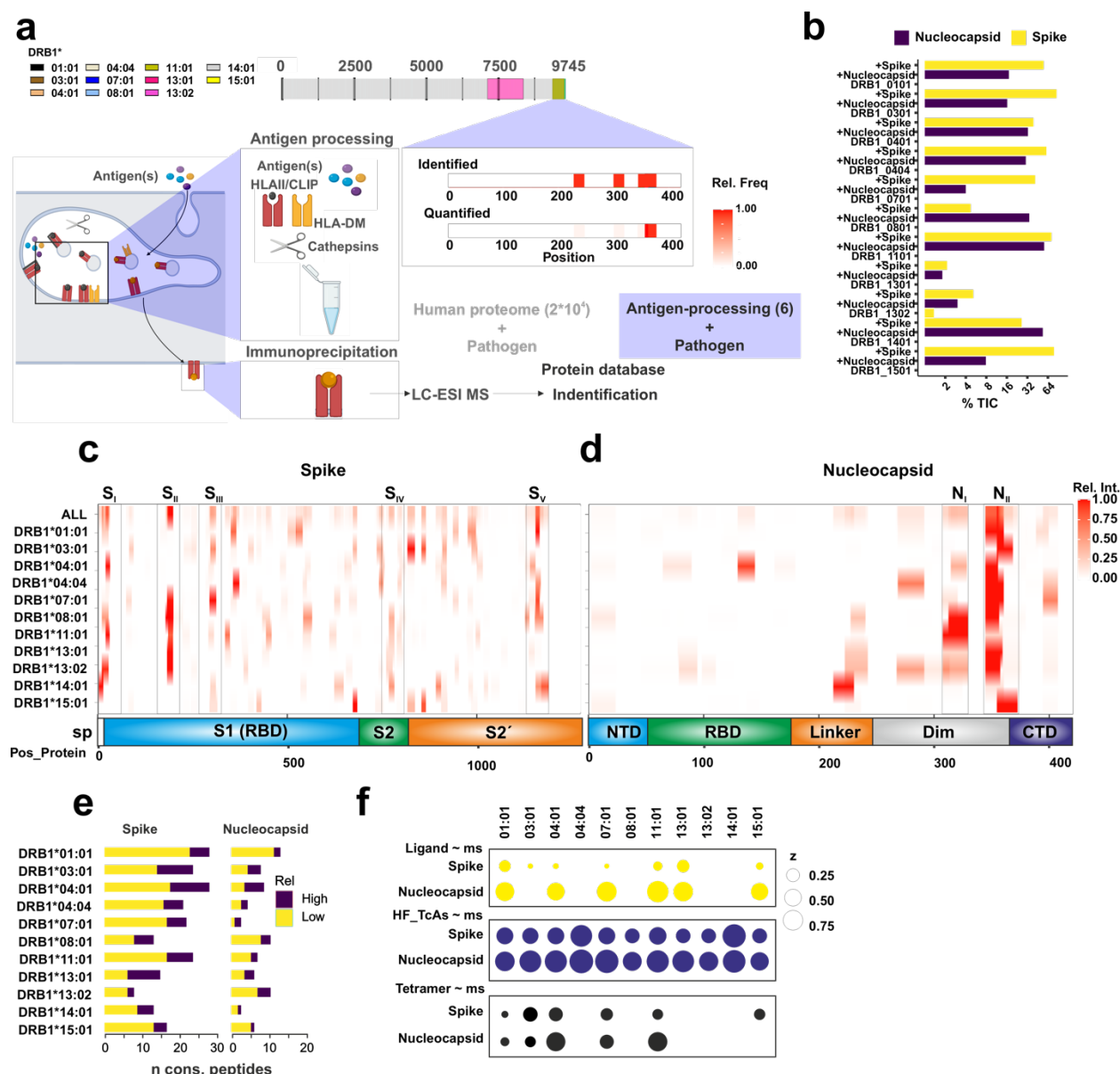
## **A reconstituted antigen processing system points out different pathways for peptide selection from the Spike and Nucleocapsid proteins**

We used an experimental system to probe immunodominance patterns of human responses to SARS-CoV-2 using our DRB1 panel. We focused on the Spike and Nucleocapsid proteins since these two antigens have a great contribution to the total immune response to SARS-CoV-2, and there is a large pool of information available on the curated IEDB entries (Table S2 and S3). We applied a reconstituted antigen processing system previously described<sup>25</sup> and tuned by us for identifying T cell epitopes from complex antigenic mixtures<sup>26</sup>. While this experimental model system will not entirely reflect the cellular environment in its complexity, it features the main steps of the MHCII-antigen processing and presentation pathway. In particular, this model incorporates several key components such as proteases and the peptide loading catalyst HLA-DM<sup>3</sup>, as well as the place-holder peptide CLIP and a reducing environment (Figure 2a). We performed the experiments following previously described protocols facilitating the interaction of the antigen with CLIP-loaded MHCII molecules in the presence of DM prior to adding proteases<sup>31</sup>. After immunoprecipitation of MHCII-peptide complexes peptides are eluted and cleaned up. Subsequently liquid chromatography mass spectrometry (LC-ESI-MS) facilitates the identification of bound peptides.

We detect a significant exchange of CLIP for peptides derived from the SARS-CoV-2 antigens applied. On average the total ion current (TIC) for SARS-CoV-2-derived peptides from these antigens reached more than 50 % with only two deviations from this behavior (DRB1\*13:01 and DRB1\*13:02) with TICs for non-CLIP peptides lower than 10 % for both antigens (Figure 2b). This

low exchange rate of CLIP for these allotypes was observed with two batches of the corresponding proteins assayed in triplicates validating the outcome of the experiments (see summary in Table S4 and Figure S2). Our analysis pipeline combines the information on the series of nested peptides identified by MS for each allotype into consensus peptides as previously described<sup>32</sup>. Consensus peptides, also referred as experimentally predicted epitopes (epEp), were then mapped to their corresponding antigen, Spike (Figure 2c) or Nucleocapsid protein (Figure 2d) along with their relative intensity for each allotype. Additionally, we combined all intensity values for each residue in all allotypes assayed (MS\_ALL). This analysis revealed five regions of more than 15 amino acids from the Spike protein ( $S_I$  to  $S_V$ ) that are preferentially selected by more than 4 different allotypes (relative intensity above the median of the candidates for each allotype). Additionally, we observe a clear bias towards the selection of peptides from two regions of the Nucleocapsid protein ( $N_I$  and  $N_{II}$ ) (details in Figure S3). Interestingly, the  $N_{II}$  region, selected by all DRB1\* allotypes, map to a specific segment of the dimerization domain. Thus, while we detect an allotype-specific pattern of peptide selection for the Spike there is a clear bias to a defined region ( $N_{II}$ ) in case of the Nucleocapsid protein (Figure 2c and d).

We noted that for all the experiments, a handful of peptides accumulate most of the TIC measured. Thus, we hypothesized that peptides preferentially selected represent regions with a higher likelihood to be displayed to T cells in a cellular environment and thus become epitopes (Figure 2e). To test this hypothesis, we evaluated to what extent the intensity values of the MS data explain residues that have been described as: i) Ligands for each allotype, Lig (subset of those with measured affinities, and including only as positive hits those with Binding affinities lower than 1000 nM), ii) entries yielding frequent/prevalent responses, HF\_TcAs, and iii) those identified/confirmed via multimer staining, Tet\_TcAs. We considered independent binary logistic regression models for each of these entries at the residue level (positive and negative hits coded as 1 and 0, respectively). The MS intensity values also coded at the residue level, and supplied as continuous variable, were considered the independent variable to fit the corresponding models. This analysis concludes that the measured MS data has a decent performance to identify residues from peptides reported as HF\_TcAs, with PPV that go up to 0.75 (Figure 2g). PPV for Lig or Tet\_TcAs on the other hand are relatively small, indicating a certain disagreement between the experimental data and that previously reported. The smaller number of entries considered in these cases could account, at least partly, for these lower PPVs. Together, we conclude that the reconstituted antigen processing system selects HF\_TcAs which have not been previously described as binders or in Tetramer stains. However, since these peptides are eluted from MHCII molecules they should be considered ligands, and may represent not yet probed peptides in multimer stains.



**Figure 2. Overview of the reconstituted *in vitro* antigen processing system.** **a.** Rationale and experimental overview of the antigen processing system. The protein database used for the searches consists of 311 entries including common MS-contaminants and abundant proteins from the host used for recombinant protein expression (Exp. Host, *S. frugiperda*), those used for protein manipulation and antigen processing proteins as well as all reference SARS-CoV-2 proteins. **b.** Summary of the performance of the antigen processing system for the selection of candidate antigenic peptides. The sum of the MS1 intensities for the SARS-CoV-2 identified peptides is represented for each of the 3 sets of experimental conditions tested. Note that the average of  $n = 2$  experiments with 3 technical replicates measured for each has been considered. **c.** Summary of the antigenic regions selected by each allotype (stated in the left) over the Spike protein and depicted according to their relative MS1 intensity. In the lower part the main domains defined for this protein are indicated. The relative frequency of the overlap of candidates selected by all allotypes considered is shown in upper row, referred as ALL. **d.** The same as in c. but for the Nucleocapsid protein. **e.** Summary of the number of consensus peptides (maximum of the overlap of the series of nested peptides) identified by MS. The color code refers to the number of peptides defined as those with higher MS1 intensity than the mean (High), or lower (Low). **f.** Summary of the Positive Predictive Value (PPV) for the identification of known ligands (Binding affinity measured lower than 1000 nM, shown in the upper panel, yellow), epitopes that give rise to prevalent responses (middle panel, blue, HF\_TcAs), and those validated by tetramer stains (lower panel, black), based on MS-data for the two antigens used. Independent models were attained and tested for each allotype considered by our DR-panel. The size of the dot refers to the PPV as shown in the legend.

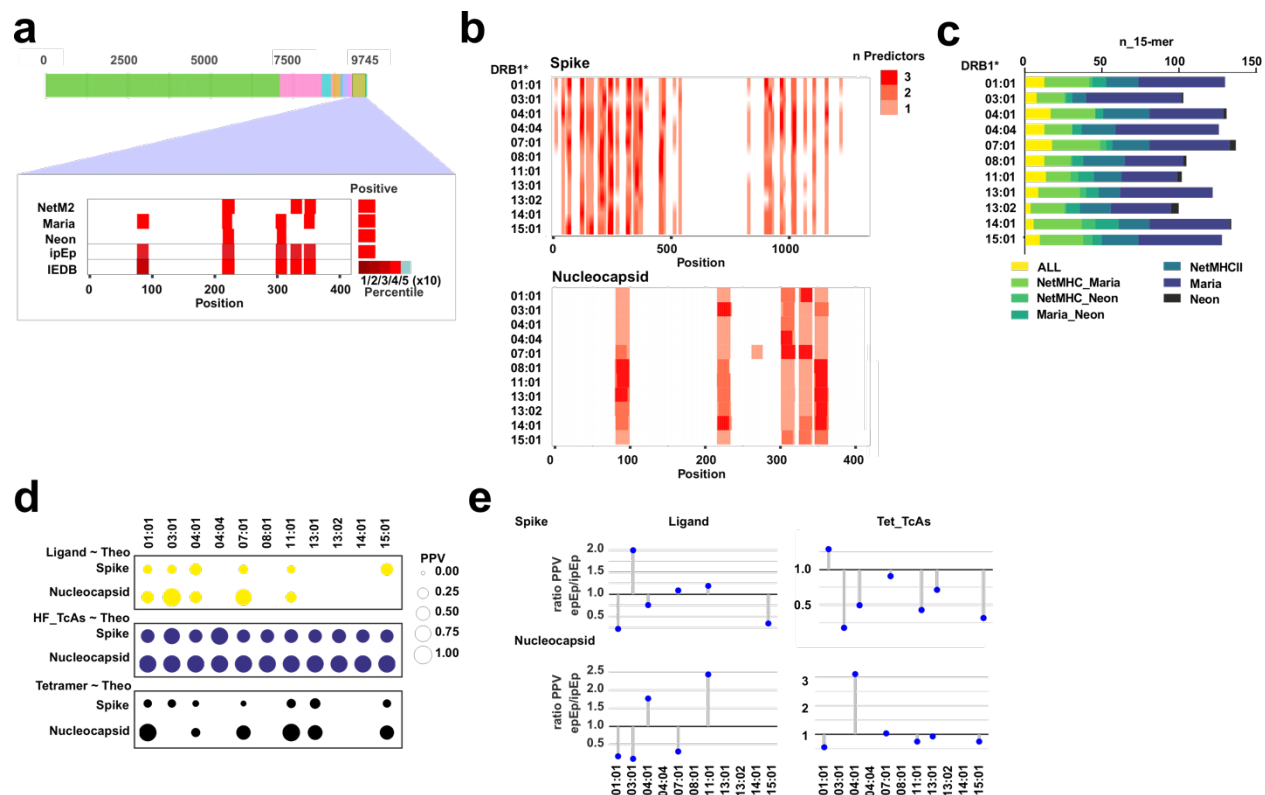


## A systematic *in silico* analysis reveals potentially immunogenic regions not overlapping with those selected experimentally

We applied *in silico* prediction tools to gain further insights on immunodominance patterns of human responses to the Nucleocapsid and Spike proteins of SARS-CoV-2. We considered validating the results of these predictions on the IEDB curated entries as we did for candidates selected by the reconstituted *in vitro* antigen processing system. We performed the analysis over the entire viral proteome and subset, where necessary the analysis for the Spike and Nucleocapsid proteins. All potential 15-mers in the viral proteome (average size of MHCII-epitopes) were queried with each of the DRB1\* allotypes from our panel using NetMHCIIpan4.0<sup>33</sup> (binding), MARIA<sup>34</sup> and NeonMHCII<sup>35</sup> (presentation predictions). We also considered the outcome of the IEDB CD4<sup>+</sup> T cell immunogenicity prediction tool<sup>36</sup> as a flagging criterion. This tool ranks peptides according to their similarity to previously described epitopes for a panel of seven common MHCII allotypes (DRB1\*03:01, 07:01 and 15:01 as well as DRB3\*0101, \*02:02, DRB4\*01:01 and DRB5\*01:01), and its combined score brings together binding motifs and immunogenicity information. Interestingly, each of the three tools considered selects specific antigenic regions as depicted in Figure 3a for the Nucleocapsid protein as model antigen, and DRB1\*03:01 as allotype. Furthermore, most of the peptides from all these regions are flagged by the combined score of the IEDB immunogenicity prediction tool. Therefore, we considered the sum of all predictors as *in silico* predicted Epitopes (ipEP).

The combined output of these tools selects recurrently around 5 (average 22 peptides) and 25 (100 peptides) antigenic regions (consecutive and non-interrupted appearance of 15-mers), for both the Nucleocapsid and Spike protein, respectively (Figure 3b) for our DRB1\*-panel. A similar picture is observed for the viral proteome with 152 regions and 623 peptides selected by all allotypes considered (Figure S4a). Noteworthy, most ipEP are selected only by one or two of the predictors considered for the Spike and Nucleocapsid protein (Figure 3c), as well as for the entire viral proteome (Figure S4b), and no ipEP is tagged by all predictors at the same time for the 11 DRB1\* allotypes (Figure S4c, ALL). Surprisingly, the two AI-based tools predicting “presentation” rather than “binding”, namely Maria and Neon, show the lowest congruence level with each other (in the range of 2-5 %). The level of agreement between these two predictors drops dramatically for DRB1\*13:02 to less than 1 % and seems, in general, lower for allotypes with charged residues in their binding motifs (Figure S4d).

We reasoned that the appearance of ipEP clusters over antigenic regions, but especially their recurrent selection by the different predictors could represent higher confidence candidates. Thus, we defined a Score prediction (Sp) for every residue consisting on the number of times that it has been identified by the three tools used. Under these premises we tested the predictive potential of Sp to classify IEDB curated entries as we did for the MS1 intensity in case of the epEP candidates. Sp is a discrete and continuous variable and we considered as well binary logistic regression models, with: i) Ligands for each allotype, Lig, ii) entries yielding frequent/prevalent responses, HF\_TcAs, and iii) those confirmed via multimer staining, Tet\_TcAs, as dependent variable. The models attained for the Spike and Nucleocapsid protein (Figure S5a) yield positive predictive values (PPV) similar to those attained for the epPe (Figure 3d). Interestingly, if we considered the entire proteome these values drop considerably up to 3-fold (Figure S5b). Thus, either a low density of positive hits, a high number of false positives, or most likely a combination of both should be responsible for these values. Interestingly, the ratio between the PPV for the experimental model vs. the *in silico* predictions (ratio epEP/ipEP) for Lig and Tet\_TcAs identifications (Figure 3e) reveals that none of the methods outperforms clearly the other one when explaining entries from the IEDB curated data for these two antigens.



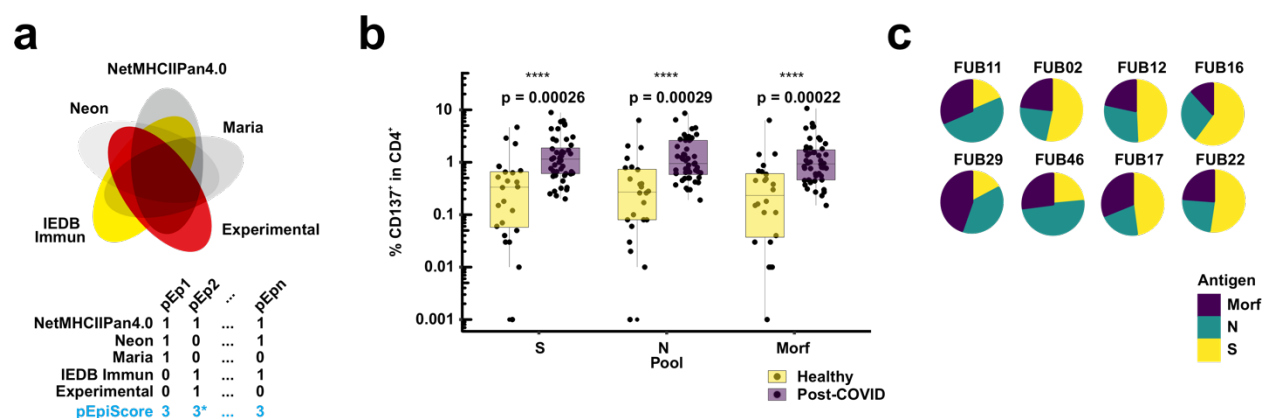
**Figure 3. Overview of the *in silico* predictions considered.** **a.** All potential 15-mers (sliding window of 1 amino acid) of the entire SARS-CoV-2 proteome (shown as concatenated proteins color-coded by orfs) were queried *in silico* as MHCII-binders/presented peptides using three different tools (NetMHCIIpan, Maria and Neon). The IEDB combined score of the CD4<sup>+</sup> T cell immunogenicity prediction tool was used to flag potential candidates. Candidates identified by each of these tools (exemplified by DRB1\*03:01 using the nucleocapsid protein), are mapped to each antigen (highlighted in red for the example). The IEDB immunogenicity prediction tool combined score ranges between 0 and 100, and candidates above the 50 % percentile are considered very unlikely to become immunogenic (see legend). Candidates identified by any of these tools are considered *in silico* predicted Epitopes (ipEp). **b.** Correlation between predictors for each DR-allotype, indicated as number of predictors identifying each residue (overlap) for each allotype for the Spike (up) and Nucleocapsid (down) proteins. **c.** Overlap between the different predictors used for each allotype for the Spike and Nucleocapsid proteins shown as a bar-chart. **d.** Positive Predictive Value of the logistic regression models generated based on predictions to identify residues as: Ligands for individual allotypes where there is available information (Lig); Prevalent T cell epitopes (HF\_TcAs, identified in pools of more than 10 individuals as eliciting responses in more than 5); Tetramers including peptide and MCHII allotype information (Tet\_TcAs). **e.** PPV ratios for MS-derived and *in silico* models for Ligand and Tetramer identifications (see legend). A ratio lower than 1 refers to a better PPV of the *in silico* model whereas a ratio higher than 1 refers to a higher PPV of the MS-model for each allotype considered.

### Minimal peptide pools recapitulate CD4<sup>+</sup> T cell responses observed with larger peptide pools *ex vivo*

We have previously considered combining both experimental and *in silico* information to enable efficient epitope discovery from complex antigenic mixtures<sup>26,31</sup>. However, we did not consider using this information to gain mechanistic insights on epitope selection. We explored the potential of the combination of these two approaches to define a minimalistic set of candidates with a broad population coverage, validating their immunogenicity and addressing antigen processing mechanistic questions. We considered averaged MS1 intensity from the reconstituted antigen processing system for all DRB1\* allotypes, and a cumulative Sp for all allotypes according to the *in silico* tools as explanatory variables. The performance of these variables, independently and in combination, to classify residues from the Spike and Nucleocapsid proteins as known HF\_TcAs

IEDB reveals a certain improvement of their combination (Figure S6a). Owing to the different sizes, we considered whether any 15-mer identified as ipEp lies within epEp regions and used the number of tools selecting them suggest a prioritizing scheme (Figure 4a). This approach is limited two structural proteins Spike and Nucleocapsid, main targets of immune responses in natural infections<sup>9,10</sup>. However, since Membrane protein, as well as orf3a and orf8 are also relevant targets of immune responses upon infection<sup>9,10</sup>, we incorporated the main candidates from these antigens identified by the *in silico* approach for testing their immunogenicity. Our final selection of candidates comprised 59 peptides with 31 entries for the Spike (S pool), 15 for the Nucleocapsid (N pool), and 13 for the combined set of Membrane (6), orf3 (5) and orf8 (2) (Morf pool, Table S5, Figure S6b). To minimize the number of peptides we considered padding on either the N- or the C- termini to cover different restrictions where register shifts may occur between allotypes.

We assayed the ability of each of the peptide pools to trigger T cell activation on Peripheral Blood Mononuclear Cells (PBMCs) from our donors by Intracellular Cytokine Staining (ICS). The panel of donors includes healthy individuals, who had not been exposed to SARS-CoV-2, and were thus considered as controls, and individuals who had been diagnosed and recovered from infection with SARS-CoV-2. Peptides from each antigen (Spike and Nucleocapsid, S and N Pools, respectively) as well as from the Membrane, orf3 and orf8 (Morf Pool) were used to trigger T cell activation. Overall, each of the pools yields a considerably higher activation level of CD4<sup>+</sup> T cell on samples from post-infected individuals than it does for the control group, (CD4/CD137 shown in Figure 4b). The percent of CD4<sup>+</sup> T cells activated in these individuals ranges between 2-15 %, reaching a very similar frequency of responders as previously shown by others<sup>9</sup>. There is a trend, not reaching statistical significance, on the extent of the increased response in individuals carrying a higher dosage of allotypes from our DRB1\* panel. Of note, we detect similar activation levels in 4 pre-pandemic samples (FUB55, FUB57, FUB58 and FUB59). Overall, these results along, with the increased levels of all cytokines, activation and exhaustion markers tested for all post-COVID individuals when compared to those of healthy individuals validate our minimalistic peptide pools designed for the broad coverage DR-panel (Figure S7a).



**Figure 4. Overview of the design and validation of the broad coverage peptide pools based on the experimental information and the *in silico* prediction tools.** **a.** Schematic representation of the scoring system considered. **b.** Frequency of activated (CD137<sup>+</sup>) CD4<sup>+</sup> T cells upon stimulation with each of the pools indicated on the x-axis as determined by flow cytometry. Two groups were tested for each peptide pool, individuals non-previously infected with SARS-CoV-2 (Healthy, yellow, n=24) and individuals that had recovered from SARS-Cov-2 infection (Post-COVID, purple, n=48). The difference between the median of the responses was compared applying a non-parametric Mann-Whitney test, and the significance is reported as follows: \*p < 0.05; \*\*p < 0.01; \*\*\*p < 0.001; \*\*\*\*p < or = 0.0001. Individuals with no response detected were set to 0.001. **c.** Pie charts illustrating the percentage of total cells responding to each of the peptide pools for a subset of Post-COVID individuals showing biased responses towards one of the pools.

Previous studies addressing immune responses to the entire viral proteome have shown the relative frequency of responses to each antigen<sup>9</sup>. We evaluated the extent of the response to each of these pools at the individual level. There is a homogeneous distribution of roughly one third of the total frequency of response measured to each of the three pools in almost every individual (Figure S7b). By exception, a handful of individuals (n = 8) feature an altered behavior showing an increased proportion of T cells responding to either the S or the N pools with only one individual (FUB29) with a relatively large proportion of T cells towards the Morf pool (Figure 4c). Together, these minimalistic peptide pool recapitulate activation levels achieved by larger pools and we could hypothesize that they shall contain immunodominant epitopes restricted by the broad-coverage DR allotypes.

### **The combined action of *in silico* predictions and *in vitro* approaches highlights promiscuous binders and inter-individual differences in CD4<sup>+</sup> T cell responses**

Our results validated the performance of the three peptide pools over a broad spectrum of restrictions. Next, we evaluated whether the selected peptides feature immunodominance from an allotype-specific point of view. We performed a targeted validation of these candidates over individuals carrying DRB1\*03:01-\*07:01 allotypes. We reasoned that the five individuals bearing this combination (2 uninfected, and 3 Post-COVID) will allow significant testing in a relatively confined HLA background (full HLA class II-type in Figure S8a). A total of 23 candidates out of the total of 59 peptides represented in all peptide pools were selected. These candidates include epEp/ipEp selected preferentially for either of the allotypes, or both (Figure 5a). We defined two pools, each of them consisting of 14 peptides for each allotypes with an overlap of 5 peptides. These peptides rank on the top positions for each DRB1 allotype according to our selection criteria (see Table S5). First, we tested the immunogenicity of these peptides when combined as pools via ELISPOT. We detected IFN $\gamma$ , as predominant pro-inflammatory cytokine in viral infections, and IL-10 as an immunosuppressive cytokine typically secreted by regulatory T cells. Each of the allotype-specific pools achieves a similar response as a positive control antigen for common human pathogens (Figure 5b). Minimal background responses towards the shortlisted peptides in pre-COVID samples is detected, while T cell reactivities in Post-COVID donors was clearly stronger. Notably, an additive effect was observed for the two pools, suggesting a complementary effect of the peptides found in each pool (Figure 5c). Interestingly, we note an individual-specific cytokine secretion pattern varying from predominant pro-inflammatory (FUB32 focused on IFN $\gamma$ ) to mainly suppressive (FUB69 mostly secreting IL-10), with very low numbers of polyfunctional, or dual cytokine producing cells (Figure S8b).

We then evaluated the contribution of the candidates considered in each of the pools to the total response generated and determined the potential restrictions of the peptides used. Owing to sample limitations we applied dual-peptide combinations in the same ELISPOT setup as described above. For each combination we included one candidate from each of the two pools. The total amount of immune response elicited by the dual peptide combinations reaches similar levels in all tested individuals (100 spots eq.  $4 \times 10^3$  SFU/  $10^6$  of CD4<sup>+</sup> T cells). These results are very consistent in terms of frequencies of responders and cytokine profiles (IFN $\gamma$ :IL-10) for the sum of the dual peptide combinations, each allotype-specific pool, or its combination (Figure S8c). We used the binding affinity of each peptide to either of the two allotypes considered to define whether they have any preferential restriction (Figure 5d, Figure S8d). More than three quarters of the peptide candidates have a high (8 out of 23) or medium (10 out of 23) affinity towards at least one of the selected allotypes. Interestingly, some of the combinations tested triggered relatively high and recurrent responses in all three individuals recovered from SARS-CoV-2





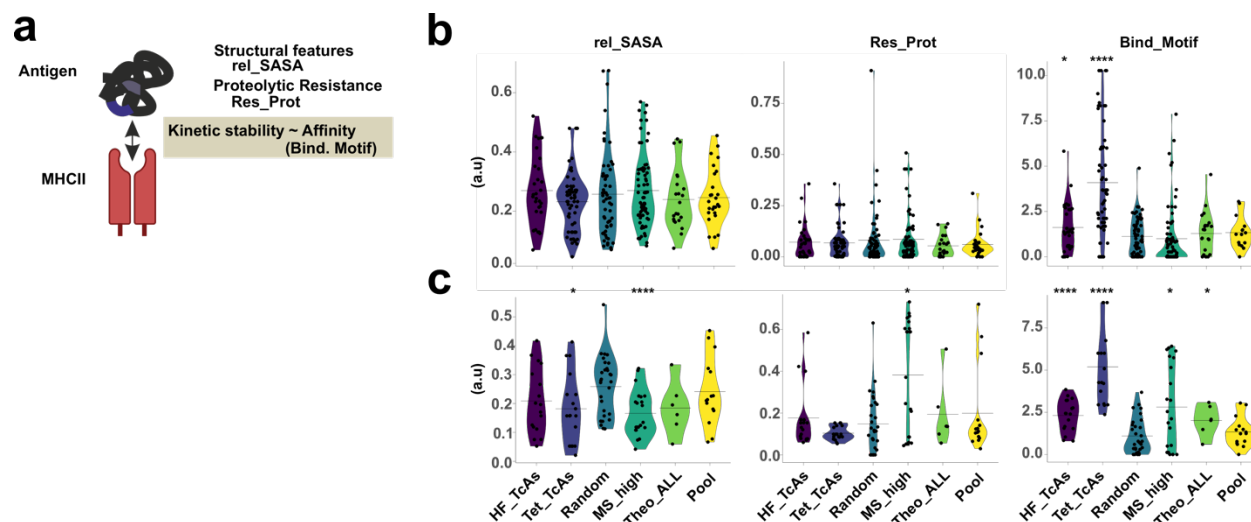


## Molecular signatures of epitope selection for the Spike and Nucleocapsid proteins

All previous results validated the immunodominant character of the selected epitopes. CD4<sup>+</sup> T cell epitopes are considered to be primarily selected on the basis of MHCII binding motifs. However, additional antigenic features are known to play a key role on epitope selection. We wondered whether we could score the relevance of these features based on the available data on the IEDB, as well as on the *in silico* and *in vitro* candidates. We considered at first: i) the binding motifs (Bind\_Motif) available in the antigens for each MHCII allotype considered (defined by NetMHCIIpan4.0), ii) the structural environment in which amino acids are located reflected by the relative surface accessible solvent area (rel\_SASA), and iii) the resistance to proteases as experimentally defined by us (Res\_Prot) (Figure 6a). Each of these antigen-related features were coded at the amino acid level (Figure S9a). Additionally, we considered a Binding-Motif-Summary for the entire DRB1\*-panel referred to as BM\_ALL (for those cases where MHC restriction is unknown), and proteolytic degradation maps that were generated *in vitro* for the two antigens considered (Figure S9b). Together, Nucleocapsid and Spike feature a different pattern for Res\_Prot (Figure S9c), which may or may not have an impact on candidate peptide selection.

We reasoned that it would be beneficial to take into account the context in which the amino acids are settled for scoring the impact of these antigenic features on peptide selection. We applied a sliding-window or context-dependent analysis whereby all residues included in an epitope/candidate are considered to generate the corresponding score (Figure 6a bottom). We attained the scores for each of the three referred features for the hits of our *in silico* analysis (Theo\_ALL) and the reconstituted *in vitro* antigen processing system (MS\_high), as well as for those peptides confirmed in tetramer stains (Tet\_TcAs), and those yielding prevalent responses (HF\_TcAs) in the IEDB. We also obtained these scores for a selection of random peptides of each of the antigens, which were used as controls (Random). This analysis confirms that Tet\_TcAs and HF\_TcAs from both antigens feature increased Bind\_Motif scores. Thus far, only Nucleocapsid entries defined by our experimental approach follow a similar trend in terms of Bind\_Motif (Figure 6b). Furthermore, the reconstituted *in vitro* antigen processing system seems to select a specific subset of entries from the Nucleocapsid protein, with high rel\_SASA and Res\_Prot (Figure 6c). Surprisingly, this behavior is partly followed by HF\_TcAs entries, showing as well lower rel\_SASA values.

These results confirm that, the epitopes described in the IEDB with a clear restriction match the binding motifs of the relevant allotypes (Tet\_TcAs). Moreover, a theoretical binding motif (BM\_ALL) consisting of the combination of those Bind\_Motif from the 11 DRB1\* allotypes used in this work, is proven to be a descriptive feature for peptides yielding recurrent responses in the IEDB (HF\_TcAs). We could differentiate a rather clear pattern for the candidates selected by the reconstituted *in vitro* antigen processing system in case of the Nucleocapsid protein. Thus, low SASA regions, resistant to proteases are preferentially selected, indicating that perhaps regions protected from degradation are made available to binding to MHCII at some point, thus favoring their selection. Interestingly a partial overlap of this behavior is seen as well for the HF\_TcAs from this same antigen.



**Figure 6. Impact of antigenic features on candidate epitope selection.** **a.** Scheme illustrating antigen- and MHCII-dependent features influencing antigenic peptide selection (top), and summary of the context-dependent analysis performed (bottom). Binding Motif (BM) matching for individual allotypes, or a sum of all alleles considered by our DRB1 panel (BM\_ALL), Resistance to Proteases (Res\_Prot) determined experimentally, and relative Solvent Accessible Surface Area (rel\_SASA) are coded at the amino acid level. For each candidate (Pepi) considered an average value for each feature ("y") is retrieved based on the values for these features ("x") of the constituent amino acids(aai). **b.** Distribution of rel\_SASA, Res\_Prot and BM of all hits previously described in the IEDB (HF\_TcAs and Tet\_TcAs), those identified by our reconstituted antigen processing system (MS\_high), the *in silico* approach proposed (Theo\_ALL), and a Random set of peptides for the Spike protein. In case of HF\_TcAs and Theo\_ALL the values from BM\_ALL are considered in the corresponding graph as there is no restriction known for each entry. The violin plots show the density of each distribution, as well as the mean. Saphiro-Wilk's test confirmed that most data do not follow a normal distribution (details in Figure S10). Differences in the average values for each feature considered between the selection of random picked entries and each of the sub-sets indicated were evaluated by Wilcoxon signed-rank test. Significance: \*p < 0.05; \*\*p < 0.01; \*\*\*p < 0.001; \*\*\*\*p < or = 0.0001. **c.** Same as in (b) but for the Nucleocapsid protein.

## Discussion

Our work conceives a pipeline for the identification of immunodominant epitopes to human pathogens and applies it to SARS-CoV-2 antigens. This pipeline relies on experimental information, which is combined and integrated with *in silico* predictions for a broad population coverage panel of DRB1\* allotypes. The output of a reconstituted antigen processing system for the DRB1\* panel considering two major antigens of SARS-CoV-2, and *in silico* predictions for the entire proteome, allow us designing three minimalistic peptide pools consisting of a total of 59 peptides from five viral proteins. The immunogenicity of the resulting candidates is validated *ex vivo* by comparing responses of un-infected and SARS-CoV-2 exposed individuals. The extent of the responses from COVID19-recovered individuals with these pools reaches similar levels as those attained with 10-times larger pools<sup>9</sup>. Moreover, 15 out of a restricted set of 23 peptides considered for fine-testing on a frequent HLA-background (DRB1\*03:01/07:01) feature immunodominant responses, demonstrating the performance of the proposed workflow. Taking advantage of the validated pipeline as well as available knowledge acquired over the pandemics, we provide evidence for immunodominance patterns beyond individual MHCII and binding motifs. The reconstituted antigen processing system recurrently selects candidates from regions with high predicted affinities, low SASA and high resistance to proteases, features previously described to be involved in antigen processing. Interestingly, these features are well represented in epitopes that have been previously validated for this antigen. In case of the Spike glycoprotein, the larger pool of known epitopes and selected by our workflow feature exclusively agreement with the corresponding binding motifs. Our work in terms of antigen processing-related features

should be nurtured with more data to provide a more comprehensive picture. Together, this research provides the basis for conceptualizing minimal peptide pools with broad population coverage that should enable improved peptide vaccination strategies.

Accessing immunodominance patterns to score immunological function would enable maximizing the number of epitopes that provide the broadest coverage while keeping the complexity of the mixture as low as possible. Under these premises it is worth noting that more than 1800 CD4<sup>+</sup> epitope entries and approximately double the number of T cell assays were compiled at the IEDB between November 2019 and September 2022. Despite sampling bias, we reasoned that we could benchmark our epitope discovery scheme with regard to the existing information as it is regularly performed for the *in silico* prediction validation<sup>33–35</sup>, but also to gain insights on the impact of antigen processing constraints and HLA-restriction. We opted for a scheme, where epitopes flagged as tetramers and affinity values below a conventional threshold (Aff < 1000nM) were considered sufficient to define HLA restrictions of a ligand, similar to a recent publication<sup>38</sup>. In contrast to this conservative approach, we considered entries eliciting recurrent responses (RF > 0.5, and n>10, named here HF\_TcAs here) where the restriction element was unknown. Our DRB1\* panel has a decent coverage throughout the individuals considered in these studies, hence these allotypes are potential restricting elements of those HF\_TcAs. Exemplarily, 80 % of the individuals recruited in one of the most comprehensive studies validating immunodominance to SARS-CoV-2 expresses at least one of the allotypes considered by us<sup>9</sup>. Likewise, if we take into account all epitopes described in TcAs data with associated HLA information, the estimated probability that one of our selected allotypes is the restricting element is higher than 0.5 in approximately 80% of the entries. This value reflects the phenotypic coverage for each individual entry of our panel according to those restrictions stated by these studies.

We considered accessing immunodominance patterns at the human population level prioritizing experimental information and refining the searches supported by *in silico* tools. We have previously validated this pipeline for the assessment of immune responses in the context of complex antigenic mixtures for a limited set of MHCII restrictions<sup>26</sup>. Note that these experiments provide direct evidence for peptide binding to the MHCII used, but also degradation under proteolytic conditions and editing by HLA-DM. Important differences with regard to previous works are: the number of restrictions tested, and the use MHCII molecules pre-loaded with CLIP as compared to “empty”-molecules<sup>6,25,37</sup>, mimicking endosomal conditions where epitope selection primarily takes place. Our results highlight a limited number of regions recurrently selected by all allotypes for the Nucleocapsid protein, exemplified by the N<sub>II</sub> region, in contrast to a broader sampling for the Spike protein. Previously, we refined candidate selection from the experimental work, by the identification of candidates fitting binding motifs<sup>26</sup>. Now, we considered the combined output of three state of the art MHCII binding/presenting predictors to streamline this candidate epitope selection. Interestingly, the regions identified by these three predictors lie within the threshold of the IEDB immunogenicity prediction tool to consider them immunogenic. A popular alternative approach for accessing immunodominance as the one presented here considers *in silico* predictions to point out potential immunogenic regions, and then deconvolutes T cell reactivities using peptide pools of candidates spanning these regions experimentally. This framework was implemented for *Mycobacterium tuberculosis* antigens<sup>38</sup> and it has been applied recently to SARS-CoV-2<sup>39</sup>. This strategy limits the initial selection of regions to the 20 % percentile of hits estimated for seven common DRB allotypes by the IEDB immunogenicity prediction tool and it is estimated to reach up to 50 % of individual’s immune response to a pathogen. Together, we can only conclude that each of these approaches has its own advantages and limitations.

An interesting aspect dealt with, is that of the design of peptide pools dedicated to CD4<sup>+</sup> T cell research, which usually consist of more or less complex mixtures of 15-mers overlapping 10-mer<sup>9</sup>

or 11-mer<sup>10</sup>. These pools have proven to be extremely useful but neither all binding registers nor the impact of N- and C-terminal extensions for CD4<sup>+</sup> T cell responses are usually taken into account<sup>40</sup>. Exemplarily, we focused on five highly immunogenic antigens at the population level<sup>9,10</sup> whose coverage for fine epitope-mapping using 15-mers will require either 2000 (to cover all potential 15-mers), 550 (to reach all binding registers), or 454 (to reach the estimated 50 % of the total response considered by the abovementioned report) peptides. Applying our workflow, we considerably reduced these numbers to attain similar activation levels of T cells as it has been reported with peptide pools that are considered to reach 50 % of the total immune responses under very similar experimental conditions (measured as % of CD137<sup>+</sup> in CD4<sup>+</sup> cells in the range of 0.5-10 %) <sup>15,41</sup>. We can therefore conclude that the reported strategy represents a highly efficient method to define immunogenic regions and postulate peptide pools. Future studies should aim at scoring the tradeoffs of considering different peptide lengths and/or peptides to aim at covering a broader spectrum of the total response.

We prioritized minimizing the impact of additional MHCII restrictions for validating the immunodominant profile of these candidates selected. Thus, we opted for a very specific MHCII background (DRB1\*03:01/DRB1\*07:01) with minimal allotype variation. The set of 23 peptides (over 59) consisting on two sub-pools covering both DRB1\* restrictions displayed an excellent performance for re-calling immune responses by each pool, their combination, or the dual peptide combinations considered. We could assign a clear restriction by either or both allotypes for more than 2/3 of the tested peptides by measuring their affinities to the MHCII. However, our selection of candidates included 1 peptide with no measurable binding affinity for these allotypes under the assayed conditions (S<sub>933-947</sub>), and 5 with relatively low binding affinity (S<sub>414\_430</sub>, S<sub>623\_639</sub>, S<sub>1150-1166</sub>, N<sub>339-361</sub> and M<sub>91-114</sub>). All of the low-affinity binders except for the M<sub>91-114</sub>, have been identified by the reconstituted antigen processing system. The stringent identification and filtering criteria applied in the experimental setup led us to consider poor-binders as intermediate or final products of antigen processing hot-spots (e.g. preferentially selected upon binding and proteolysis), or mis-assigned peptides in the final candidate list refinement. Exemplarily, the region of the Spike protein selected in the DRB1\*07:01 experiments covers S<sub>937-951</sub> (representing up to 0.17 of the TIC of this antigen) instead of S<sub>933-947</sub> selected by us. While this may have allowed us covering other restrictions, we may have neglected a relevant binding register in this case. Interestingly, these candidates may yield relevant individual-specific responses. Under these premises one could speculate that either personalized antigen processing profiles or TcR repertoires could be responsible for this effect. Together, these results validate the performance of the proposed approach for the identification of immunodominant epitopes. Importantly, despite the previous extensive work using overlapping peptide pools, and specifically addressing this MHCII background<sup>9</sup>, our approach is still able to reveal new reactivities, e.g. N<sub>81\_99</sub>, S<sub>985\_1006</sub> and O3a<sub>23\_41</sub>.

Seminal studies have demonstrated the impact of antigen processing constraints on the selection of immunodominant epitopes focusing on the DRB1\*01:01 allotype in relation to Influenza's HA<sup>5</sup> and recently the entire HIV<sup>6</sup> proteome. Thus, antigen processing constraints could significantly inform epitope discovery<sup>37,42-44</sup>. However, the main limitations of these studies are: i) their restricted MHCII diversity (e.g. how general are the observations found for one allotype), and ii) the experimental and logistic limitations for epitope mapping/validation (e.g. assessment of T cell responses from relevant samples). Previous studies into the impact of antigen-specific constraints on immunodominance pointed out structural features of regions that are preferentially selected for display, showing a relation to folding and solvent accessible surface area (SASA)<sup>42,43</sup>. Our work capitalizes on the extensive data resulting from the impact of SARS-CoV-2 pandemics. This work allows us to design a dedicated and controllable experimental workflow to provide a comprehensive overview of the impact of antigen processing constraints on a broad panel of DRB1\* allotypes. Indeed, we could assign specific epitope selection patterns depending on



antigenic features. The present work elaborates on a limited context-dependent analysis for a relatively limited number of antigens. However, future works elaborating on expanding knowledge on CD4<sup>+</sup> T cell epitopes<sup>27</sup> combined with novel tools as AlphaFold enabling efficient structure predictions<sup>45</sup> will probably grant us with more accurate information.

Our method demonstrates the potential of integrating experimental and *in silico* approaches to access and understand immunogenicity. We conclude that the impact of allotype-specific parameters such as DM-susceptibility<sup>46</sup>, antigen-related features (e.g. SASA or protein foldedness)<sup>44,45</sup> and experimentally definable processing constraints (e.g. proteolytic degradation) should be funneled into the next generation pipelines of effective MHCII-immunodominance prediction. Integrating these diverse factors into prediction platforms holds the potential to improve our ability to tailor vaccines and therapies in an HLA haplotype manner. Thus, this perspective might become increasingly relevance for the future design of more personalized T cell responses during the course of a new pandemic or in cancer immunotherapy.

## Materials and Methods

### *In silico* candidate epitope prediction

The SARS-CoV-2 genome sequence was obtained from the NCBI database <https://www.ncbi.nlm.nih.gov/nuccore/1798174254>. We extracted the sequences of the proteins orf1ab, S, orf3a, E, M, orf6, orf7a, orf7Bb, orf8, N, and orf10 based on the reference genome. We used a sliding window size of fifteen amino acids and a step of one amino acid for the following analysis (9591 peptides). Potential SARS-CoV-2 epitopes were identified using a novel selection workflow based on the integration of prediction algorithms for peptide-MHC class II binding and immunogenicity. The peptide-MHC class II binding/presentation prediction was performed using three different algorithms, netmhciipan<sup>33</sup>, maria<sup>34</sup>, and neonmhc2<sup>35</sup>, with percentile score cut-off of 10 %. Immunogenicity scores used as flagging criteria were determined using the IEDB tool<sup>36</sup> applying an immunogenicity score cutoff value of 50 %, yielding 1643 immunogenic peptides. Potential SARS-CoV-2-derived epitopes were identified as the top-ranked overlapping candidates for each allotype of the eleven MHC class II allotypes. Next, allele-specific lists of peptides with a minimum length of 15 residues were defined for each SARS-CoV-2 protein.

### Peptides and viral antigens

Peptides were purchased from GL Biochem (Shanghai) Ltd (10 mg purity > 95 %). Lyophilized peptides were diluted at a final concentration of 10mM in DMSO and subsequently diluted in PBS. When necessary, the pH was adjusted to 7.4.

Nucleocapsid and Spike Glycoproteins expressed in HEK cells were purchased from Sino biological (Cat Numbers. 40588-V08B for Nucleocapsid and 40589-V08B1 for Spike) as C-terminal His-tagged. Lyophilized proteins were reconstituted according to the manufacturer specifications.

### Protein methods

MHC proteins (HLA-DRs and HLA-DM) were expressed as previously described<sup>47</sup>. Briefly, HLA-DR cDNAs are cloned into the pFastBacDual vector and include Leu-Zippers in their C-termini as well as a sequence encoding for the CLIP peptide followed by Thrombin cleavage site and a G4S linker in the N-termini of the DRB1 polypeptides. Furthermore, the DRA polypeptide encodes as well a Biotin Acceptor Sequence in the C-termini of the corresponding Leu-Zipper. Expression was achieved by infection of Sf9 cells at an MOI of 5 and harvesting the cells after 4 days. Protein purification was achieved by immunoaffinity chromatography using a L243-FF-Sepharose resin casted in house. In all cases HLA-DR proteins were cleaved with Thrombin and gel-filtrated using a Sephadex S200. For the reconstituted *in vitro* system, these proteins were also cleaved with V8



protease to remove the Leu-Zippers. HLA-DM cDNAs are cloned into pFastBacDual and include a Flag-Tag in the C-termini of HLA-DMA chain. Purification in this case was achieved using an immunoaffinity M2-Sepharose resin, protein was eluted using Glycine pH3.5. After dialysis, the protein was concentrated (Vivaspin MWCO 10kDa) and gel filtrated (Sephadex S200).

### Peptide binding experiments

Peptide binding affinities to selected HLA-DR molecules were determined by competition experiments using fluorescently labelled reporter peptides. Reporter peptide binding signal was measured by FP. HLA-DR molecules expressed in insect cells were thrombin cleaved to facilitate peptide exchange. Competition experiments were set by adding 100 nM HLA-DR, 100nM reporter peptide (CLIP-FITC for DRB1\*07:01 or MBP-FITC for DRB1\*03:01) and titrated concentrations of the corresponding peptide in 50mM Citrate Phosphate buffer containing 150mM NaCl at pH 5.3. Each reaction was measured after 12h incubation at 37° C, and the corresponding IC50 values for each peptide were retrieved by fitting a sigmoidal function to the obtained data points.

### *In vitro* reconstituted antigen processing system

The previously described cell-free reconstituted *in vitro* system<sup>25</sup> was modified according to the specific needs of the experiments<sup>26</sup>. HLA molecules together with the candidate antigens and the HLA-DM were incubated for 2 h at 37° C in citrate phosphate 50 mM pH 5.2 in the presence of 150 mM NaCl. Cathepsins were added to reaction mixtures after incubation with L-Cysteine (6 mM) and EDTA (5 mM). The final reaction mixture was incubated at 37° C for 2 to 5 hours. Afterwards the pH was adjusted to 7.5, and Iodoacetamide was added (25 µM). Immunoprecipitation (IP) of the pMHCII complexes was performed using L243 covalently linked to Fast Flow sepharose. Peptides were eluted from purified MHCII adding TFA 0.1 % to the samples. Peptides are separated from the MHCII molecules by using Vivaspin filters (10 kDa MWCO). Cathepsin B (Enzo), H (Enzo) at a molar ratio of 1:250, and S (Sigma) at a molar ratio of 1:500 were used in these experiments.

### Proteolytic degradation of antigens

Spike and Nucleocapsid proteins were incubated in the presence of cathepsins in the molar ratios indicated above. Reactions were performed at 37 °C citrate phosphate 50 mM pH 5.2 in the presence of 150 mM NaCl and stopped at t = 0, and t = 3 h, by adding Iodoacetamide, immediate transfer of the samples to ice. The pH was then adjusted to 7.5 by adding Tris-HCl 1 M pH8.0. Samples were splitted, and used for SDS-PAGE analysis, Western blotting and peptide identification by MS. For MS analysis, samples were dried in a SpedVac and treated as described below.

### LC-MS measurements

All samples were initially cleaned up by reverse phase C18 enrichment. The eluates were dried in a SpeedVac, and peptides were reconstituted in 20 µl of H<sub>2</sub>O containing acetonitrile (4 %), and TFA (0.05 %). 6 µl of these mixtures were analyzed using a reverse-phase capillary system (Ultimate 3000 nanoLC) connected to an orbitrap connected to a Q Exactive HF mass spectrometer (Thermo Fisher Scientific). Samples were injected and concentrated on a trap column (PepMap100 C18, 3 µm, 100 Å, 75 µm i.d. × 2 cm, Thermo Fisher Scientific) equilibrated with 0.05 % TFA in water. After switching the trap column inline, LC separation was performed on a capillary column (Acclaim PepMap100 C18, 2 µm, 100 Å, 75 µm i.d. × 25 cm, Thermo Fisher Scientific) at an eluent flow rate of 300 nL/min. Mobile phase A contained 0.1 % formic acid in water, and mobile phase B contained 0.1 % formic acid in 80 % acetonitrile/ 20 % water. The column was pre-equilibrated with 5 % mobile phase B followed by a linear increase of 5–44 % mobile phase B in 70 min. Mass spectra were acquired in a data-dependent mode utilizing a

single MS survey scan ( $m/z$  350–1,650) with a resolution of 60,000, and MS/MS scans of the 15 most intense precursor ions with a resolution of 15,000. The dynamic exclusion time was set to 20 seconds and automatic gain control was set to  $3 \times 10^6$  and  $1 \times 10^5$  for MS and MS/MS scans, respectively.

### Mass spectrometry data processing

MaxQuant (v2.0.3.0) with implemented Andromeda peptide search engine was used for analyzing the raw MS and MS/MS data. All searches were done on the basis of unspecific protease cleavage, main ion search tolerance of 10 ppm and MSMS tolerance search of 50 ppm and enabling the feature “match between runs”. The reconstituted *in vitro* antigen processing samples were searched against a database containing the sequences of all SARS-CoV-2 proteins (note that Spike and Nucleocapsid were substituted for the sequences of the recombinant ones), cathepsins, MHCII and all reviewed *Spodoptera frugiperda* proteins (Uniprot, access on March 2020) as an internal control. The database used for the cathepsins digestion experiments included only the protein antigen sequence used, and the corresponding sequences of the cathepsins. In both cases a FDR of 0.01 (1 %) was used as on a decoy search. All identifications with an FDR higher than 0.01, reverse identifications and contaminants (identified by MaxQuant) were removed for data analysis. Each set of experiments was analyzed together, treating technical and biological replicates as independent samples.

All MS raw files from the reconstituted *in vitro* antigen processing experiments were processed as previously reported in<sup>26</sup>. In brief, allotype-specific subset identifications from the evidence file were submitted to the plateau webserver using the same database used for the peptide identification. Each of the consensus peptides identified was then used to determine replication and retrieve a MS1 relative intensity. These relative intensities were averaged throughout the different replicates, and unless otherwise indicated only peptides found in at least 2 technical replicates out of 2 independent experiments were considered. For the proteolytic degradation experiments we took the spectral counts for each peptide identified from the corresponding peptide.txt files.

### Donor recruitment and HLA-typing

Donors were recruited in the Medizinische Fakultät/Universitätsklinikum of the Otto-von-Guericke Universität Magdeburg. Informed consent was signed and agreed upon according to the ethic protocol approved by the corresponding university. HLA-typing was performed by the DKMS-LSL facility.

### PBMC isolation and fractionation

PBMCs were isolated from blood samples of either healthy individuals or donors recovered from SARS-CoV-2 infection via density-gradient sedimentation. PBMCs were either frozen in cell-freezing medium supplemented with 10 % DMSO or further processed for isolation of monocytes and CD4<sup>+</sup> T cells. PBMCs were thawed and counting of living cells was performed using trypan blue in a cell counting chamber. When stated PBMCs were used directly in specific experiments, and for others specifically isolated cellular fractions were used. In these cases, cells were isolated by magnetic cell separation using CD14 microbeads for monocytes, and subsequently CD4 microbeads for T helper cells (all Miltenyi Biotec, Bergisch Gladbach, Germany). The homogeneity of the cell preparations was controlled by flow cytometry. For flow cytometric analysis a total of  $1 \times 10^5$  PBMCs per well were seeded on 96-well plates, cultured in X-VIVO 15 medium (Lonza, Basel, Switzerland) supplemented with 4 % human AB plasma (Innovative Research, Novi, MI, USA), and provided with the corresponding peptide pools (17.5 ng/ml per peptide).

## ELISpot

Dual secretion of IFN $\gamma$  and IL10 was determined using the enzymatic Human IFN $\gamma$ /IL-10 Double-Color ELISPOT Kit (Cellular Technology, Shaker Heights, OH, USA) and using pre-isolated CD14<sup>+</sup> monocytes and CD4<sup>+</sup> T cells. 5x10<sup>4</sup> APC were seeded together with 1x10<sup>5</sup> CD4<sup>+</sup> T cells in the presence of relevant peptides at 2.5 ng/ul diluted in X-VIVO 15 medium (Lonza, Basel, Switzerland) supplemented with 4 % human AB plasma (Innovative Research, Novi, MI, USA) in 96-well plates. The cells were then pre-incubated for 48 h, washed, transferred to an ELISpot plate and further incubated for 60 h. The secreted cytokines were determined according to the manufacturer's instructions; counting of spots place on an ImmunoSpot analyzer (Cellular Technology). Dual secretion of cytokines was determined by overlapping the corresponding signals (IFN $\gamma$ -red and IL-10-blue). The extent of the response is directly correlated to the surface area covered by each signal, in this case determined for each colony counted.

## Antibodies and reagents used in cell culture experiments

The following antibodies were purchased from the stated vendors and used according to the manufacturer specifications: for western-blotting Rabbit anti-6HisTag (Abcam ab1187). In flow cytometry experiments we used: anti-CD4, anti-CD137, anti-CD319 (all Miltenyi Biotec), anti-CD3, anti-PD-1, anti-IL-2, anti-TNF $\alpha$ , anti-IFN $\gamma$  (all Biolegend). And in case of Elispots we considered: anti-IFN $\gamma$  (capture and FITC-labelled detection antibodies), anti-IL-10 (capture and biotinylated detection antibodies), anti-FITC HRP (all Cellular Technology).

## Flow cytometry

To accumulate cytokines, cells were treated for 4 h with 5 mg/ml Brefeldin A (Merck, Darmstadt, Germany) after 140 h incubation time. Further, cells were briefly reactivated by addition of 10 ng/ml PMA and 1  $\mu$ g/ml ionomycin (all Merck) for 1h prior to flow cytometric analysis. Cells were then harvested, Fc-receptors blocked (FcR Blocking reagent, Miltenyi) and stained for extracellular markers (CD3, CD4, CD137, PD-1, and CD319), subsequently fixed with 2 % paraformaldehyde (Morphisto, Offenbach am Main, Germany), permeabilized with 0.5 % saponine (Merck), and stained for intracellular cytokines (IFN-g, IL-2, and TNF $\alpha$ ). Samples were acquired on a FACSCanto II flow cytometer with FACS-Diva software (v10, BD Bioscience).

## Analysis of antigen processing rules and antigenic peptide features

Each protein residue was assigned a value for each of the parameters considered: resistance to proteolytic degradation (Res\_Prot), relative surface accessible solvent area (rel\_sasa) and binding motif matching (BM). For the parameter resistance to proteases we considered the sum of spectral counts for each residue (see details above). For rel\_SASA we used the relative values calculated using the pdb PISA tool<sup>48</sup> using as input the structural models available for the Spike and the Nucleocapsid protein as of September 2021 in the Zhang lab website (<https://zhanggroup.org/COVID-19/>). For the binding motif parameter (Bind\_Motif) we considered all weak (coded as 1) and high (coded as 2) affinity binding cores (9mers) defined by NetMHCIIpan4.0<sup>33</sup> for each antigen and specific allotype. Note that Bind\_Motif is allele-specific and thus, we considered as well a cumulative value to represent the potential overlap between binding motifs for all allotypes used (BM\_ALL). Likewise, for this analysis all available entries in the IEDB described as epitopes, ligands or T cell Assay related (TcAs) were equally coded as binary vectors (0: not present, 1: part of a hit) (Accessed Sept 2022). The sum of appearances of each residue as a Hit facilitates defining their relative frequency. Furthermore, where additional data was available, e.g. binding affinity measured experimentally, this information was used to subset those hits. Finally, each candidate epitope identified by the reconstituted antigen processing system was coded at the amino acid level considering the MS1 intensity (e.g. Spike\_DRB1\*01:01\_n for candidate n derived from the Spike protein identified for DRB1\*01:01).

We worked with either normalized (to the maximum of each vector), count-based of hits or binary-based vectors, as stated in each case.

For the context dependent analysis, we determined the average for each parameter according to the values of each amino acid. As control we used a randomized pool of peptides.

$$y = \frac{1}{n} \sum_{i=1}^n x_i$$

$$x_{aa_i} = \begin{cases} Res\_Prot \\ rel\_SASA \\ Bind\_Motif \\ BM\_ALL \end{cases} \quad y_{Pep} = \begin{cases} Res\_Prot \\ rel\_SASA \\ Bind\_Motif \\ BM\_ALL \end{cases}$$

### Statistical analysis and model description

All analysis were carried out using R, version 4.2.1 over the RStudio suite unless otherwise stated. Binary logistic regression models were generated based on the general expression:

$$\text{logit}(p) \sim (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k) \\ p(x_i | Hit | x_i \text{predictors})$$

Each predicted candidate from the in silico tools (ipEp), identified experimentally (epEp) or its combination (pEp) were considered as single or combined explanatory variables ( $\beta_1$  to  $\beta_n$ ) to define a binary response output at the residue level ( $y = 0$  for negative, or  $y = 1$  for positive or selected hits). Different types of IEDB curated entries were considered as response variables following the criteria indicated (Ligand, Tetramers and HF\_TcAs). A more detailed description is also provided in the supplemental information.

Goodness of fit and the performance of the models were estimated by calculating the Akaike Information Criterion (AIC) and considering the Area Under the Curve (AUC) from Receiver Operator Curves (ROC). Likewise, we considered the positive predictive value of each model as indicated below:

$$PPV = \frac{nTP}{nTP + nFP}$$

### References

1. Sercarz, E.E., Lehmann, P.V., Ametani, A., Benichou, G., Miller, A., Moudgil, K. Dominance and Crypticity. *Annu Rev Immunol.* **11**, 729-66 (1993).
2. Wieczorek, M. *et al.* Major Histocompatibility Complex (MHC) Class I and MHC Class II Proteins: Conformational Plasticity in Antigen Presentation. *Front Immunol.* **8**, 292 (2017).
3. Álvaro-Benito, M. & Freund, C. Revisiting nonclassical HLA II functions in antigen presentation: Peptide editing and its modulation. *HLA.* **96**, 415–429 (2020).
4. Miller, M.A., Ganesan, A.P.V., Luckashenak, N., Mendonca, M. & Eisenlohr, L.C. Endogenous antigen processing drives the primary CD4+ T cell response to influenza. *Nat Med.* **21**, 1216–1222 (2015).
5. Cassotta, A. *et al.* Deciphering and predicting CD4+ T cell immunodominance of influenza virus hemagglutinin. *J Exp Med.* **217**(10), e20200206 (2020).
6. Sengupta, S. *et al.* A cell-free antigen processing system informs HIV-1 epitope selection and vaccine design. *J Exp Med.* **227**(7), e20221654 (2023).



7. Wragg, K.M. *et al.* Establishment and recall of SARS-CoV-2 spike epitope-specific CD4<sup>+</sup> T cell memory. *Nat Immunol.* **23**, 768–780 (2022).
8. Becerra-Artiles, A. *et al.* Broadly recognized, cross-reactive SARS-CoV-2 CD4 T cell epitopes are highly conserved across human coronaviruses and presented by common HLA alleles. *Cell Rep.* **39**(11), 110952 (2022).
9. Tarke, A. *et al.* Comprehensive analysis of T cell immunodominance and immunoprevalence of SARS-CoV-2 epitopes in COVID-19 cases. *Cell Rep Med.* **2**(2), 100204 (2021).
10. Loyal, L. *et al.* Cross-reactive CD4<sup>+</sup> T cells enhance SARS-CoV-2 immune responses upon infection and vaccination. *Science.* **374**(6564), eabh1823 (2021).
11. Niessl, J., Sekine, T. & Buggert, M.T cell immunity to SARS-CoV-2. *Semin Immunol.* **55**, 101505 (2021).
12. Mather, M. W., Jardine, L., Talks, B., Gardner, L. & Haniffa, M. Complexity of immune responses in COVID-19. *Semin Immunol.* **55**, 101545 (2021).
13. Sette, A. & Crotty, S. Adaptive immunity to SARS-CoV-2 and COVID-19. *Cell.* **184**, 861–880 (2021).
14. Braun, J. *et al.* SARS-CoV-2-reactive T cells in healthy donors and patients with COVID-19. *Nature.* **587**, 270–274 (2020).
15. Mateus, J. *et al.* Selective and cross-reactive SARS-CoV-2 T cell epitopes in unexposed humans. *Science.* **370**(6512), 89-94 (2020).
16. Sahin, U. *et al.* COVID-19 vaccine BNT162b1 elicits human antibody and TH1 T cell responses. *Nature.* **586**, 594–599 (2020).
17. Moderbacher, C.R. *et al.* Antigen-Specific Adaptive Immunity to SARS-CoV-2 in Acute COVID-19 and Associations with Age and Disease Severity. *Cell* **183**, 996-1012.e19 (2020).
18. Bacher, P. *et al.* Low-Avidity CD4<sup>+</sup> T Cell Responses to SARS-CoV-2 in Unexposed Individuals and Humans with Severe COVID-19. *Immunity.* **53**, 1258-1271.e5 (2020).
19. Poon, M.M.L. *et al.* SARS-CoV-2 infection generates tissue-localized immunological memory in humans. *Sci Immunol.* **6**, 9105 (2021).
20. Heide, J. *et al.* Broadly directed SARS-CoV-2-specific CD4<sup>+</sup> T cell response includes frequently detected peptide specificities within the membrane and nucleoprotein in patients with acute and resolved COVID-19. *PLoS Pathog.* **17**, e1009842 (2021).
21. Painter, M.M. *et al.* Rapid induction of antigen-specific CD4<sup>+</sup> T cells is associated with coordinated humoral and cellular immunity to SARS-CoV-2 mRNA vaccination. *Immunity.* **54**, 2133-2142.e3 (2021).
22. Goel, R.R. *et al.* Distinct antibody and memory B cell responses in SARSCoV-2 naïve and recovered individuals following mRNA vaccination. *Sci Immunol.* **6**, 1–19 (2021).
23. Low, J.S. *et al.* Clonal analysis of immunodominance and crossreactivity of the CD4 T cell response to SARS-CoV-2. *Science.* **372**, 1336–1341 (2021).
24. Wragg, K.M. *et al.* Establishment and recall of SARS-CoV-2 spike epitope-specific CD4 + T cell memory. *Nat Immunol.* **23**, 768–780 (2022).
25. Hartman, I.Z. *et al.* A reductionist cell-free major histocompatibility complex class II antigen processing system identifies immunodominant epitopes. *Nat Med.* **16**, 1333–1340 (2010).
26. Ebner, F., *et al.* CD4<sup>+</sup> Th immunogenicity of the *Ascaris* spp. secreted products. *NPJ Vaccines.* **5**, 25 (2020).
27. Vita, R. *et al.* The Immune Epitope Database (IEDB): 2018 update. *Nucleic Acids Res.* **47**, D339–D343 (2019).
28. Grifoni, A. *et al.* Characterization of Magnitude and Antigen Specificity of HLA-DP, DQ, and DRB3/4/5 Restricted DENV-Specific CD4<sup>+</sup> T Cell Responses. *Front Immunol.* **10**, 1568 (2019).
29. Gourraud, P.-A. *et al.* HLA Diversity in the 1000 Genomes Dataset. *PLoS One.* **9**, e97282 (2014).



30. Dhanda, S.K. *et al.* ImmunomeBrowser: a tool to aggregate and visualize complex and heterogeneous epitopes in reference proteins. *Bioinformatics*. **34**(22), 3931-3933 (2018).
31. Álvaro-Benito, M., Ebner, F., Bertazzon, M. & Morrison, E. CD4+ T Cell Epitope Identification from Complex Parasite Antigen Mixtures. *Methods Mol Biol*. **2673**, 89–109 (2023).
32. Álvaro-Benito, M., Morrison, E., Abualrous, E.T., Kuroпка, B. & Freund, C. Quantification of HLA-DM-Dependent Major Histocompatibility Complex of Class II Immunopeptidomes by the Peptide Landscape Antigenic Epitope Alignment Utility. *Front Immunol*. **9**, 872 (2018).
33. Reynisson, B., Barra, C., Kaabinejadian, S., Hildebrand, W.H., Peters, B. & Nielsen, M. Improved Prediction of MHC II Antigen Presentation through Integration and Motif Deconvolution of Mass Spectrometry MHC Eluted Ligand Data. *J Proteome Res*. **19**(6), 2304–2315 (2020).
34. Chen, B. *et al.* Predicting HLA class II antigen presentation through integrated deep learning. *Nature Biotechnol*. **37**, 1332–1343 (2019).
35. Abelin, J.G. *et al.* Defining HLA-II Ligand Processing and Binding Rules with Mass Spectrometry Enhances Cancer Epitope Prediction. *Immunity*. **51**, 766-779.e17 (2019).
36. Dhanda, S.K. *et al.* Predicting HLA CD4 Immunogenicity in Human Populations. *Front Immunol*. **9**, 1369 (2018).
37. Kim, A. *et al.* Divergent paths for the selection of immunodominant epitopes from distinct antigenic sources. *Nat Commun*. **5**, 1–16 (2014).
38. Lindestam Arlehamn, C.S. & Sette, A. Definition of CD4 immunosignatures associated with MTB. *Front Immunol*. **5**, 1–7 (2014).
39. Grifoni, A. *et al.* A Sequence Homology and Bioinformatic Approach Can Predict Candidate Targets for Immune Responses to SARS-CoV-2. *Cell Host Microbe*. **27**, 671-680.e2 (2020).
40. Holland, C.J., Cole, D.K. & Godkin, A. Re-directing CD4+ T cell responses with the flanking residues of MHC class II-bound peptides: The core is not enough. *Front Immunol*. **4**, 54167 (2013).
41. Grifoni, A. *et al.* Targets of T Cell Responses to SARS-CoV-2 Coronavirus in Humans with COVID-19 Disease and Unexposed Individuals. *Cell*. **181**, 1489-1501.e15 (2020).
42. Mettu, R.R., Charles, T. & Landry, S.J. CD4+ T-cell Epitope Prediction Using Antigen Processing Constraints. doi:10.1016/j.jim.2016.02.013.
43. Landry, S.J. *et al.* Structural Framework for Analysis of CD4+ T-Cell Epitope Dominance in Viral Fusion Proteins. *Biochemistry*. **62**(17), 2517-2529.
44. Kim, A.R., Boronina, T.N., Cole, R.N., Darrah, E. & Sadegh-Nasseri, S. Distorted Immunodominance by Linker Sequences or other Epitopes from a Second Protein Antigen during Antigen-Processing. *Sci Rep*. **7**, 1–11 (2017).
45. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature*. **596**, 583-589 (2021).
46. Abualrous, E.T. *et al.* MHC-II dynamics are maintained in HLA-DR allotypes to ensure catalyzed peptide exchange. *Nat Chem Biol*. **19**(10), 1196-1204 (2023).
47. Álvaro-Benito, M., Wieczorek, M., Sticht, J., Kipar, C. & Freund, C. HLA-DMA Polymorphisms Differentially Affect MHC Class II Peptide Loading. *J Immunol*. **194**(2), 803-16 (2014).
48. Krissinel, E. & Henrick, K. Inference of macromolecular assemblies from crystalline state. *J Mol Biol*. **372**, 774–797 (2007).

# Acknowledgements

This project was funded by a grant from the BMBF (01KI2072) to MCBW and CF, TRR186 (project A21N) to CF. MAB received funding from the DKMS-SLS (JHR-2021-01) and the Comunidad Autónoma de Madrid (2022-T1/BMD-23752). Funding was also received from BMBF LongCoCID (01EP2101C) (to MCBW) and the Covid19 program by the state of Saxony-Anhalt (FP1) to HL and MCBW. We also thank to other members of the working groups and to the Biosupramol Core facility, and Eliot Morrison for critical reading of the manuscript.

1020  
 1021 **Author contributions**  
 1022 MA-B: Conceptualization, experimental design, performed research, data analysis, wrote the  
 1023 manuscript.  
 1024 ETA: Conceptualization, experimental design, performed research, data analysis.  
 1025 HL: Experimental design, performed research, data analysis.  
 1026 SM: Performed research, analyzed data.  
 1027 JH: Performed research, analyzed data.  
 1028 JS: Conceptualization, data analysis, manuscript drafting.  
 1029 FK: Performed research.  
 1030 BK: Performed research.  
 1031 CC: Discussed results.  
 1032 MCB-W: Conceptualization, funding, manuscript drafting.  
 1033 CF: Conceptualization, funding, manuscript drafting.  
 1034  
 1035 **Competing interests**  
 1036 The authors declare no competing interests