

## RESEARCH ARTICLE

## CAMAP: Artificial neural networks unveil the role of codon arrangement in modulating MHC-I peptides presentation

Tariq Daouda<sup>1,2</sup>✉, Maude Dumont-Lagacé<sup>1,3</sup>, Albert Feghaly<sup>1</sup>, Yahya Benslimane<sup>1,3</sup>, Rébecca Panes<sup>1,4</sup>, Mathieu Courcelles<sup>1</sup>, Mohamed Benhammedi<sup>1,3</sup>, Lea Harrington<sup>1,3</sup>, Pierre Thibault<sup>1,5</sup>, François Major<sup>1,6</sup>, Yoshua Bengio<sup>6</sup>, Étienne Gagnon<sup>1,4</sup>, Sébastien Lemieux<sup>1,2,6</sup>, Claude Perreault<sup>1,3</sup>

**1** Institute for Research in Immunology and Cancer, Université de Montréal, Montréal, Canada, **2** Department of Biochemistry, Université de Montréal, Montréal, Canada, **3** Department of Medicine, Université de Montréal, Montréal, Canada, **4** Department of Microbiology, Infectiology and Immunology, Université de Montréal, Montréal, Canada, **5** Department of Chemistry, Université de Montréal, Montréal, Canada, **6** Department of Computer Science and Operations Research, Université de Montréal, Montréal, Canada

✉ These authors contributed equally to this work.

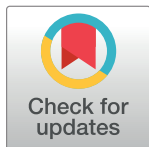
✉a Current address: Broad Institute of MIT and Harvard, Cambridge, Massachusetts, United States of America

✉b Current address: Center for Cancer Research, Massachusetts General Hospital, Charlestown, Massachusetts, United States of America

✉c Current address: Department of Medicine, Harvard Medical School, Boston, Massachusetts, United States of America

✉d Current address: Center for Immunology and Inflammatory Diseases, Massachusetts General Hospital, Charlestown, Massachusetts, United States of America

\* [tdaouda@broadinstitute.org](mailto:tdaouda@broadinstitute.org)



## OPEN ACCESS

**Citation:** Daouda T, Dumont-Lagacé M, Feghaly A, Benslimane Y, Panes R, Courcelles M, et al. (2021) CAMAP: Artificial neural networks unveil the role of codon arrangement in modulating MHC-I peptides presentation. PLoS Comput Biol 17(10): e1009482. <https://doi.org/10.1371/journal.pcbi.1009482>

**Editor:** Marco Punta, San Raffaele Hospital: IRCCS Ospedale San Raffaele, ITALY

**Received:** February 13, 2021

**Accepted:** September 27, 2021

**Published:** October 22, 2021

**Copyright:** © 2021 Daouda et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** The data underlying the results presented in the study are available from: - Human B-LCL: RNA-Seq data can be accessed on the NCBI Bioproject database (<http://www.ncbi.nlm.nih.gov/bioproject/>; accession PRJNA286122). - Human PBMC: RNA-sequencing data for human PBMC were extracted from healthy donors in Zucca et al (2019) and can be accessed under the GEO accession number GSE106443 and GSE115259, while MAPs were extracted from Murphy et al (2017). - Human B721.221: The B721.221 dataset was retrieved from Abelin et al

## Abstract

MHC-I associated peptides (MAPs) play a central role in the elimination of virus-infected and neoplastic cells by CD8 T cells. However, accurately predicting the MAP repertoire remains difficult, because only a fraction of the transcriptome generates MAPs. In this study, we investigated whether codon arrangement (usage and placement) regulates MAP biogenesis. We developed an artificial neural network called Codon Arrangement MAP Predictor (CAMAP), predicting MAP presentation solely from mRNA sequences flanking the MAP-coding codons (MCCs), while excluding the MCC *per se*. CAMAP predictions were significantly more accurate when using original codon sequences than shuffled codon sequences which reflect amino acid usage. Furthermore, predictions were independent of mRNA expression and MAP binding affinity to MHC-I molecules and applied to several cell types and species. Combining MAP ligand scores, transcript expression level and CAMAP scores was particularly useful to increase MAP prediction accuracy. Using an *in vitro* assay, we showed that varying the synonymous codons in the regions flanking the MCCs (without changing the amino acid sequence) resulted in significant modulation of MAP presentation at the cell surface. Taken together, our results demonstrate the role of codon arrangement in the regulation of MAP presentation and support integration of both translational and post-translational events in predictive algorithms to ameliorate modeling of the immunopeptidome.

(2017); RNA sequencing data can be accessed under the GEO accession number GSE93315. - Murine CT26: RNA-Seq data can be accessed under the GEO accession number GSE111092. Mass spectrometry data can be found on the ProteomeXchange Consortium via the PRIDE partner repository (human B-LCL: PXD004023 and murine CT26: PXD009065 and 10.6019/PXD009065). - Murine EL4: MAP dataset was extracted from Murphy et al (2017) and EL4 RNA sequencing dataset was extracted from Sidoli et al (2019) and can be accessed under the GEO accession number GSE125384. All figures were generated using R's package "ggplot2". Source code for pyGeno (<https://github.com/tariqdaouda/pyGeno>, doi: 10.12688/f1000research.8251.2) and Mariana (<https://github.com/tariqdaouda/Mariana>) are freely available online. Furthermore, a pyTorch implementation of CAMAP source code, as well as training and evaluation datasets are available on github (<https://github.com/tariqdaouda/CAMAP>).

**Funding:** This work was supported by grants from the Canadian Cancer Society (attributed to CP, <https://www.cancer.ca/>, number 705604 and 705714), the Oncopole (attributed to CP, <https://oncopole.ca/>) and the Leukemia & Lymphoma Society of Canada (attributed to CP, <https://www.llscanada.org/>). CP's lab is supported in part by The Katelyn Bedard Bone Marrow Association (<https://www.givemarrow.net/>). MDL was supported by a studentship from the Canadian Institute of Health Research (<https://cihr-irsc.gc.ca/>). Y.B was supported by a fellowship from the Cole Foundation (<https://www.colesfoundation.org/>). LH's lab is supported by the Canadian Institutes of Health Research operating grants (<https://cihr-irsc.gc.ca/>, #148936 and 163051). EG lab is supported by the Canadian Institute for Health Research operating grant (<https://cihr-irsc.gc.ca/>, MOP-133726). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

## Author summary

MHC-I associated peptides (MAPs) are small fragments of intracellular proteins presented at the surface of cells and used by the immune system to detect and eliminate cancerous or virus-infected cells. While it is theoretically possible to predict which portions of the intracellular proteins will be naturally processed by the cells to ultimately reach the surface, current methodologies have prohibitively high false discovery rates. Here we introduce an artificial neural network called Codon Arrangement MAP Predictor (CAMAP) which integrates information from mRNA-to-protein translation to other factors regulating MAP biogenesis (e.g. MAP ligand score and transcript expression levels) to improve MAP prediction accuracy. While most MAP predictive approaches focus on MAP sequences *per se*, CAMAP's novelty is to analyze the MAP-flanking mRNA sequences, thereby providing completely independent information for MAP prediction. We show on several datasets that the integration of CAMAP scores with other known factors involved in MAP presentation (i.e. MAP ligand score and mRNA expression) significantly improves MAP prediction accuracy, and further validate CAMAP learned features using an *in-vitro* assay. These findings may have major implications for the design of vaccines against cancers and viruses, and in times of pandemics could accelerate the identification of relevant MAPs of viral origins.

This is a *PLOS Computational Biology Methods* paper.

## Introduction

In jawed vertebrates, virtually all nucleated cells present at their surface major histocompatibility complex class-I (MHC-I) associated peptides (MAPs), collectively referred to as the immunopeptidome [1,2]. MAPs play a central role in shaping the adaptive immune system, as they orchestrate the development, survival and activation of CD8 T cells [3]. Moreover, recognition of abnormal MAPs is essential to the elimination of virus-infected and neoplastic cells [4]. Therefore, systems-level understanding of MAP biogenesis and molecular composition remains a central issue in immunobiology [5,6].

The generation of the immunopeptidome can be conceptualized in two main events: (a) the generation of MAP candidates (i.e. peptides of appropriate length for MHC-I presentation) through protein degradation, and (b) a subsequent filtering step through the binding of MAP candidates to the available MHC-I molecules. Rules that regulate the second event have been well characterized using artificial neural networks (ANN) and weighted matrix approaches [7,8]. However, accurately predicting which peptides will ultimately reach MHC-I molecules following a multistep processing in the cytosol and endoplasmic reticulum remains an open question [6]. Most efforts at modeling MAP generation have focused on post-translational events and their regulation by the amino acid sequence of MAPs and of directly adjacent residues (typically 10-mers at the N- and C-termini). While the consideration of preferential sites of proteasome cleavage has proven useful to enrich for MAP candidates [9], it remains insufficient for MAP prediction, due to prohibitive false discovery rates [10–12].

A large body of evidence suggests that a substantial portion of MAPs are produced co-translationally [13–15], deriving from defective ribosomal products (DRiPs), that is, polypeptides

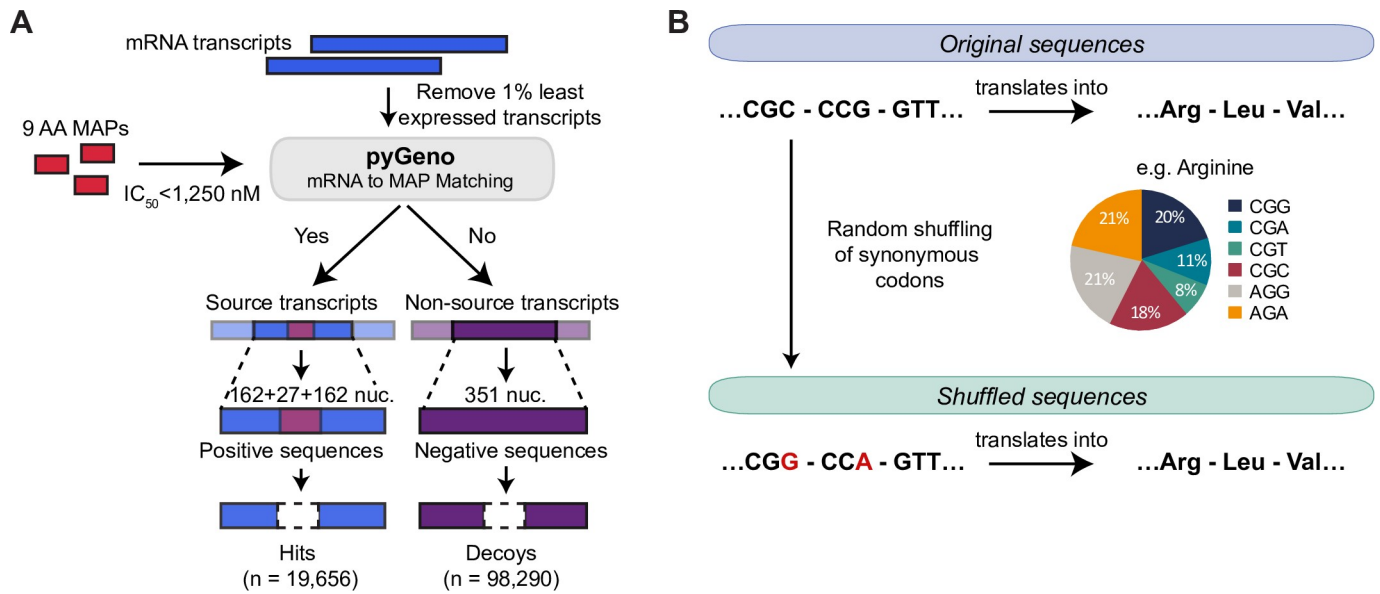
that fail to achieve a stable conformation during translation and are consequently rapidly degraded. This concept was initially supported by two observations: (i) viral MAPs can be detected within minutes after viral infection, much earlier than their associated proteins half-life [16], and (ii) MAP presentation correlates more closely with translation rate than with overall protein abundance [17,18]. In addition, while all proteins contain peptides that are predicted to bind MHC-I molecules, mass spectrometry analyses have revealed that the immunopeptidome is not a random excerpt of the transcriptome or the proteome [1,19]. Indeed, proteogenomic analyses of 25,270 MAPs isolated from B lymphocytes of 18 individuals showed that 41% of expressed protein-coding genes generated no MAPs [19]. These authors also provided compelling evidence that the presentation of MAPs cannot be explained solely by their affinity to MHC-I alleles and their transcript expression levels, while ruling out low mass spectrometry sensitivity as an explanation for the non-presentation of the strong binders. Because (i) MAPs appear to preferentially derive from DRiPs and (ii) codon usage influences both precision and efficiency of protein synthesis [20,21], we hypothesized that codon usage in the vicinity of MAP-coding codons (MCCs) might significantly contribute to the regulation of MAP biogenesis. We developed an artificial neural network called Codon Arrangement MAP Predictor (CAMAP), trained to identify MCCs flanking regions. We then used CAMAP to uncover key codon features that characterize mRNA sequences encoding for MAPs (i.e. source) when compared to sequences that do not (i.e. non-source).

## Results

### Dataset description

We analyzed a previously published dataset consisting of MAPs presented on B lymphoblastoid cell line (B-LCL) by a total of 33 MHC-I alleles from 18 subjects [19,22]. Because we were searching for features that influence MAP generation and not the binding of MAP to MHC-I molecules, we elected to analyze the MCC flanking sequences only and excluded the MCCs *per se* from our positive (hits) and negative (decoys) sequences (Fig 1A). To facilitate data analysis and interpretation, we restricted our hit dataset to MAPs with a length of 9 amino acids, for a total of 19,656 9-mer MAPs (which represents 78% of MAPs in this dataset). We next created a decoy dataset from transcripts that generated no MAPs, by randomly selecting 98,290 9-mers from these transcripts. Finally, we used pyGeno [23] to extract MCCs flanking regions corresponding to both hit and decoy MAPs, which constituted our final dataset for CAMAP. Of note, each sequence in the final dataset is unique and derives from the canonical reading frame.

In addition, in order to investigate the relative importance of codon vs. amino acid usage in MAP biogenesis, we generated a dataset of shuffled sequences (for both positive and negative datasets) in which original codon sequences were randomly replaced by synonymous codons according to their usage frequency in the human transcriptome (Fig 1B). This transformation was performed to ensure that both neural networks received the same number of parameters as input, preventing the introduction of a favorable bias for the codon network. The random shuffling causes any codon-specific feature to be shared among synonyms, thereby causing the shuffled codon distribution to reflect the amino acid usage (see [Materials and Methods](#) for more details). Indeed, codon distributions in the shuffled datasets more closely reflected those of their corresponding amino acid than in the original dataset (S1 Fig), with 92% of codons in the shuffled dataset showing a strong correlation ( $R^2 > 0.95$ ) with the amino acid distribution, compared to only 69% in the original dataset ( $p < 2 \times 10^{-16}$ , S2 Fig). It also eliminates differences in codon usage between the hit and decoy datasets. Importantly, this shuffling does not affect the resulting amino acid sequence thereby preserving all potential amino acid-related



**Fig 1. Construction of the dataset.** (A) Transcripts expressed in B-LCL from 18 subjects were considered as source or non-source transcripts depending on their match with at least one MAP. Because we were searching for features that might influence MAP generation and not the binding of MAP to MHC-I, we focused our attention on mRNA sequences adjacent to the nine MCCs (i.e. up to 162 nucleotides on each side of MCCs). (B) Creation of the shuffled dataset. Codons were randomly replaced by a synonymous codon according to their respective frequencies (i.e. codon usage) in the transcriptome. The random shuffling causes any codon-specific feature to be shared among synonyms, thereby causing the shuffled codon distribution to reflect the amino acid usage. It also eliminates the differences in synonymous codon usage between the hit and decoy datasets. Importantly, both the original sequence and its shuffled version translate into the same amino acids.

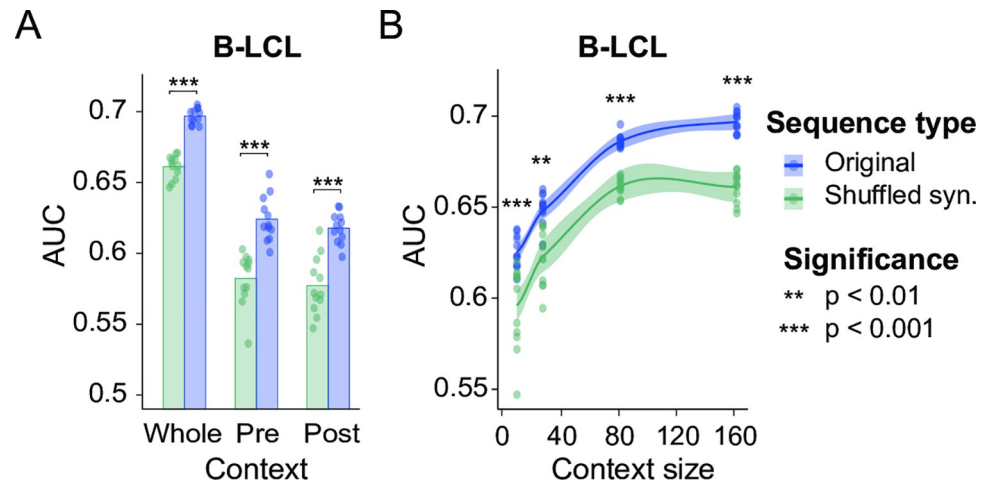
<https://doi.org/10.1371/journal.pcbi.1009482.g001>

motifs. Distributions of each codon in the original VS shuffled datasets and comparison to its corresponding amino acid can be found in S3–S20 Figs.

### CAMAP links codon usage to MAP presentation

To assess the importance of codon usage in MAP biogenesis, we reasoned that if codons bear important information that is operative at the translational rather than the post-translational level, then: (i) CAMAP trained to identify MCCs flanking regions should consistently perform better when trained on original codon sequences than on shuffled codon sequences (reflecting amino acid sequences), and (ii) synonymous codons should have different effects on the prediction. To test these hypotheses, CAMAP received as inputs MCCs flanking regions from hit and decoy sequences from either the original or shuffled datasets. It was then trained to predict the probability that individual input sequences were MCCs flanking regions (i.e. hit) rather than sequences from the negative dataset (i.e. decoy, S21A Fig).

We compared CAMAP performance when predicting MAP presentation from original codon sequences, versus shuffled sequences representing amino acid arrangement. To evaluate the robustness of our approach, 12 different CAMAPs were trained in parallel, with different train-validation-test splits of the dataset. Our results show that predictions were consistently better when CAMAP received the original codons rather than the shuffled sequences (Fig 2A). CAMAPs receiving information from both pre-MCCs and post-MCCs sequences (i.e. whole MCC flanking context) also performed better than when receiving only pre- or post-MCCs context (Figs 2A and S21B and S21C), suggesting that pre- and post-MCCs context are not redundant. Indeed, we found a weak correlation between the prediction scores of CAMAPs trained only with pre- or post-MCCs sequences (S22 Fig). In addition, CAMAPs receiving longer sequences performed better than those receiving shorter sequences (Fig 2B). Because



**Fig 2. CAMAP predictions on MAP-flanking sequences.** (A) Area under the curve (AUC) score for CAMAPs trained with whole MCCs context, versus CAMAPs trained with only pre- or post-MCCs context. All CAMAPs presented here were trained with a context size of 162 nucleotides. (B) AUC for CAMAPs trained with codon context sizes of 9, 27, 81 and 162 nucleotides (context here refer to mRNA sequences flanking the MCCs).

<https://doi.org/10.1371/journal.pcbi.1009482.g002>

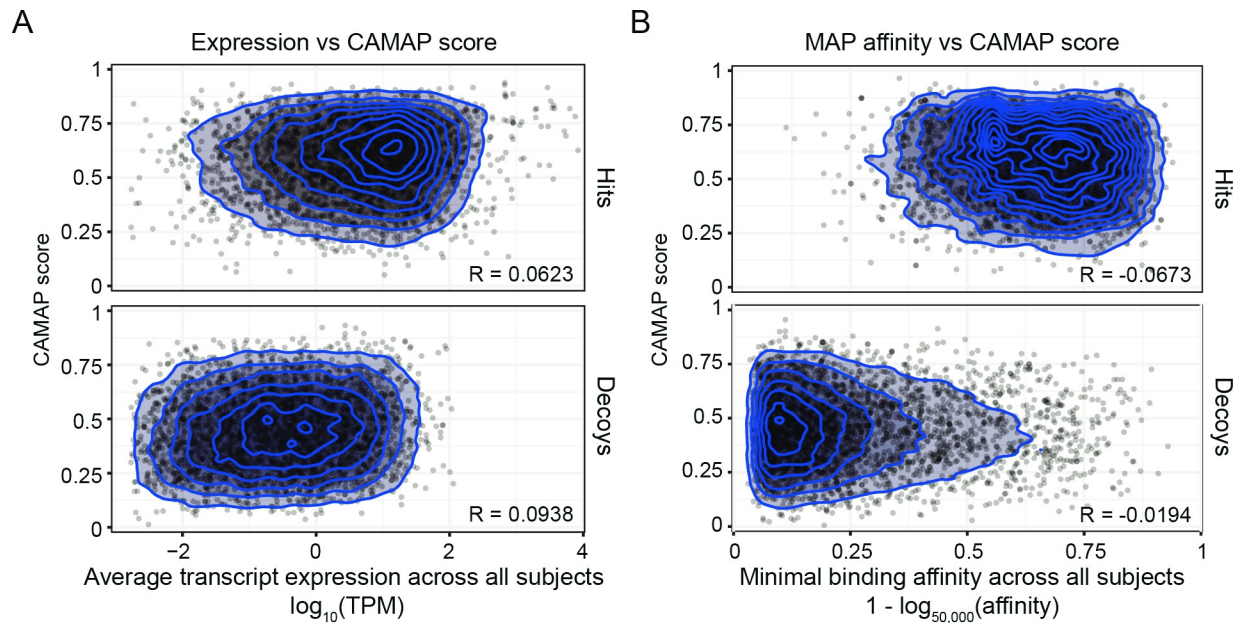
sequences located far upstream and downstream of the MCCs (i.e. in ranges exceeding the direct influence of proteases) are informative regarding MAP presentation, it supports the existence of factors unrelated to protein degradation modulating MAP presentation.

### CAMAP predictions are independent of MAP binding affinity or transcript expression level

Both MAP binding affinity to the MHC-I molecule and the level of gene expression are predictive of MAP presentation [19]. Because codon usage has been shown to be different in highly expressed genes, we wanted to verify whether the codon-specific rules captured by CAMAP were associated with potential biases in our positive dataset, which is enriched in highly expressed genes. We first show that there is no correlation between gene expression levels and CAMAP scores in both the positive and negative datasets ( $R < 0.1$ , Fig 3A). This was true for both average expression levels across our samples (Fig 3A), and for samples individually (see S23 Fig). Secondly, we trained CAMAP networks using a decoy dataset that mirrored the positive dataset gene expression level (S24A Fig) and showed similar results: CAMAP trained on original codon sequences performed better than CAMAP trained on shuffled sequences (S24B Fig). These results show that the codon-specific rules captured by CAMAP trained on original sequences are independent of gene expression levels.

While CAMAP does not receive any information about the MCC *per se*, we stipulated that the presence of MHC-I binding motifs in the MCCs in the positive dataset might be associated with biases in the MAP-flanking regions, which could also influence CAMAP training. Therefore, to evaluate the presence of this potential bias, we first evaluated the correlation between CAMAP scores and MAPs binding affinity. Again, our result showed no correlation between CAMAP scores and MAP binding affinity, both when considering the minimal binding affinity of each MAP to the MHC-I alleles contained in our dataset (Fig 3B) or when considering each allele individually (S25–S26 Figs). Secondly, we trained CAMAP networks using a decoy dataset that mirrored the positive dataset MAP binding affinities (S27A Fig). Again, CAMAPs trained on original codon sequence performed better than CAMAPs trained on shuffled sequences (S27B Fig). These results show that codon-specific rules captured by CAMAP





**Fig 3. Correlation between CAMAP prediction score and (A) transcript expression level and (B) MAP binding affinity.** CAMAP used here was trained on original codon sequences using a context size of 162 nucleotides (both pre- and post-MCCs context).

<https://doi.org/10.1371/journal.pcbi.1009482.g003>

trained on original sequences are independent of MAP binding affinities and of potential biases in codon usage of MAP-flanking sequences associated with the presence of an MHC-I binding motif in the MCCs.

We next evaluated the possibility of biases associated with many MAPs originating from conserved regions (e.g., found in multiple domains of the same domain family such as zinc fingers or kinases). We first evaluated MAPs that could originate from different transcripts within the transcriptome (i.e. transcripts with sufficient expression levels detected by RNA sequencing) as they are likely to represent conserved regions in the genome. While 79.9% of MAPs originated from unique contexts (S28A Fig), only 2.1% of MAPs had more than 3 possible origins, which represented 11.7% of the hit dataset (S28B Fig). These MAPs with several possible origins preferentially derived from zinc finger proteins, which are known to share homologous regions (S29 Fig). We therefore trained CAMAPs with datasets excluding entries encoding for MAPs that had  $>3$  or  $>10$  possible origins and compared their performance with that of CAMAPs trained without excluding these MAPs. Our results show that whatever the dataset used, CAMAP trained with original sequences always significantly outperformed CAMAP trained with shuffled sequences (S30 Fig). Taken together, these results suggest that the codon-specific rules captured by CAMAP are independent of potential homologies in the hit dataset, as they do not appear to influence CAMAP performance.

Non-source transcripts have been previously associated with higher GC content [19]. To evaluate the possibility of biases associated with GC-content, we evaluate the correlation between GC-content within MAP-flanking sequences (using a context size of 162 nucleotides) and CAMAP score. We found a negative correlation between an example's CAMAP score and its GC-content ( $R = -0.4$ ,  $p < 1 \times 10^{-16}$ , S31A Fig). This could be related to the fact that lower GC-content is associated with lower mRNA stability [24], which could lead to the formation of DRiPs. We then trained CAMAP networks using a decoy dataset that mirrored the positive dataset GC content (S32A Fig). Again, CAMAPs trained on original codon sequence

performed better than CAMAPs trained on shuffled sequences, even when corrected for GC content (S32B Fig), suggesting additional factors are at play in the regulation of MAP biogenesis.

Previous studies have shown that codon bias plays a significant role in translation efficiency (TE) by co-adaptation to the tRNA pool [25,26]. This effect has been quantified in the tRNA Adaptation Index (tAI), which is a metric for translation efficiency [27]. While hit and decoy datasets showed a similar distribution of mean tAI of the MAP-flanking sequences, we evaluated the possibility of a bias in translation efficiency in our datasets. Results show a weak correlation between mean tAI and CAMAP scores ( $R = 0.1767$  and  $0.2600$  for hits and decoys, respectively), suggesting that the signal captured by CAMAP could be in part attributable to codon usage impact on translation efficiency (S31B Fig).

We evaluated whether CAMAP would still be predictive for negative examples (decoys) derived from source transcripts. While most decoys derived from source transcripts showed a CAMAP score above 0.5 and are thus classified as source, they still had a lower CAMAP score compared to hits from source transcripts ( $p < 1 \times 10^{-16}$ , S33 Fig).

Finally, we evaluated how different shuffling methods affected CAMAP performance, to better understand the role of codon arrangement (i.e. usage and position) in CAMAP's predictions. The shuffling method used above (Figs 1 and 2) generates shuffled sequences by replacing each codon by one of its synonymous codons (including itself) according to the codon usage within the transcriptome. This shuffling method, which we designated as transcriptome shuffling (S34A Fig), eliminates any synonymous codon usage and GC content differences between the hit and decoy datasets. To evaluate the importance of codon positions in the MAP-flanking sequences, we generated a second shuffled dataset in which synonymous codons present within a given example (i.e. MAP-flanking sequences pre- and post-MCC context) are swapped with one another (S34B Fig). This shuffling method, which we call codon swap shuffling, preserves both the global codon usage, GC-content and amino acid sequences for each example within the hit and decoy datasets. Finally, we generated a third shuffled dataset in which only the third nucleotide of codons are swapped within each example, while preserving amino acid sequences and GC content (third nucleotide shuffling, S34C Fig). Here, the global codon usage is completely different from normal codon usage in humans, as the frequency of each synonymous codon is not taken into account during shuffling. The performance of CAMAP networks pre-trained on original (non-shuffled) datasets was then evaluated on each shuffled dataset.

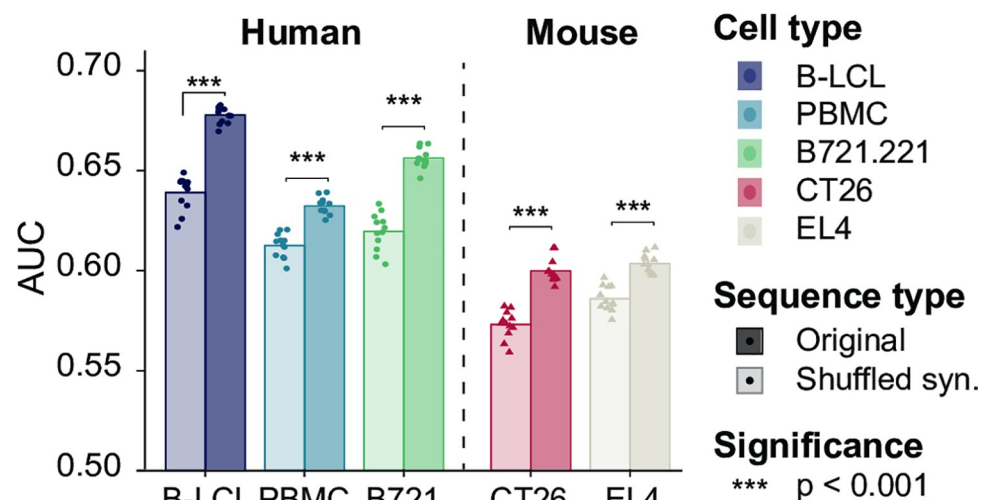
As shown in S35 Fig, the codon swap shuffling only slightly decreased CAMAP predictions compared to the original dataset (mean AUC from 0.689 for original to 0.687 for codon swap shuffling, respectively,  $p = 2.266 \times 10^{-4}$ ). This could be due to the limited number of synonymous codons available for swapping within a single example. This also suggests that the overall codon composition of the MAP-flanking sequences has more impact on CAMAP predictions than the codons' position within the MAP-flanking sequences. The third nucleotide shuffling also decreased CAMAP performances (mean AUC of 0.676) but remained above that of the transcriptome shuffled dataset (mean AUC of 0.647). These results suggest that the GC content has a significant impact on CAMAP performances, as the codon swap and third nucleotide shuffling both preserves the GC content within each example, while the transcriptome shuffling does not. This is also supported by the negative correlation between GC content and CAMAP scores, as shown in S31A Fig.

Taken together, these results support a role for codon arrangements, and especially for codon usage in the regulation of MAP biogenesis which cannot be reduced to codon's impact on gene expression levels, GC content or translation efficiency.

## Validation of CAMAP predictions on human and mouse datasets

We next validated our CAMAP trained on 9-mer MAPs derived from B-LCL using 5 datasets derived from different human and mouse cell types. All the validation datasets were described through proteogenomic analyses similarly to our B-LCL training datasets. However, all the validation datasets included MAPs of 8–11 mers, in contrast with the training dataset that contained only 9-mer MAPs. The validation datasets consisted of (i) our B-LCL dataset, this time including all peptide lengths [19,22], (ii) a dataset of human peripheral blood mononucleated cells (PBMCs) [28], (iii) a dataset of B-lymphoblastoid cells expressing unique HLA alleles (B721.221 [11]), (iv) murine colon carcinoma cell line (CT26) and (v) a murine lymphoma cell line (EL4, [28,29]). For all datasets, we created hit and decoy datasets of original and shuffled sequences using the same approach described above but including MAPs of 8–11 amino acids. Notably, CAMAPs trained on human sequences encoding 9-mers MAPs from one human cell type (i.e. B-LCL) could also predict presentation of 8–11 mers MAPs in other human cell types (Fig 4), as well as from mouse cell lines, albeit with lower performances (Fig 4). Here again, CAMAPs trained on original sequences consistently outperformed CAMAPs trained on shuffled sequences (Fig 4). These results show that the codon-specific rules derived by CAMAPs to predict MAP presentation are valid across different cell types, and can even be applied to another species, albeit with slightly lower performances. These results support a role for codons in the modulation of MAP presentation.

The lower performances of CAMAP trained with shuffled sequences (representing amino acid distribution) suggests that amino acids in MAP-flanking sequences are less informative than codons regarding MAP presentation. We formally quantified this difference in information using the Kullback-Leibler (KL) divergence (see [Materials and Methods](#) for more details). Most codons (47/61, 77%) showed greater KL divergence in the original dataset than the shuffled dataset, indicating that codon distributions contained more information with regards to MAP presentation than amino acid distributions (S36 Fig). These results suggest that codons in MAP-flanking regions play a role that is non-redundant with amino acids in MAP biogenesis.



**Fig 4. Validation of CAMAP predictions on 5 datasets derived from human and murine cell lines.** CAMAP prediction score for different datasets derived from humans (i.e. B-LCL, PBMCs and B721.221) or mouse (i.e. CT26 and EL4) cells. Of note, all CAMAPs were trained on B-LCL-derived sequences encoding for 9-mer MAPs only with a context size of 162 nucleotides. Results are reported for 8 to 11-mer MAPs derived from the 5 datasets. In all panels, 12 CAMAPs trained with original or shuffled synonymous sequences were compared (significance assessed using Student T test).

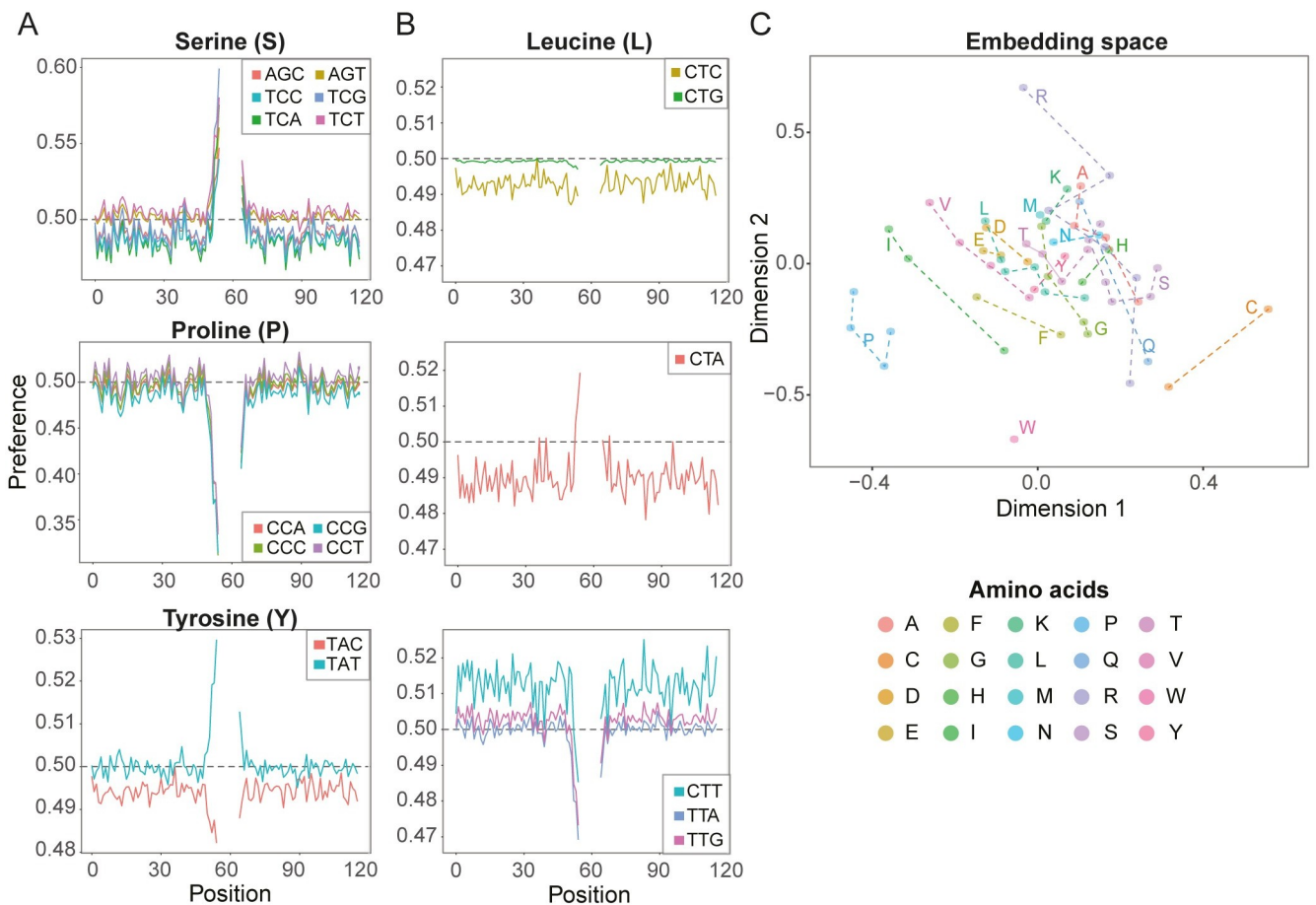
<https://doi.org/10.1371/journal.pcbi.1009482.g004>



### Codons closest to the MCC are most influential on CAMAP predictions

We wondered whether some regions were more influential on MAP presentation than others. To address this question, we retrieved the model preferences for each codon at each position. The preferences correspond to the prediction score of our best model (trained with original codon sequences for a context size of 162 nucleotides) when a single codon at a single position is provided as input (all other positions being set at [0,0] coordinates in the embedding space, see [Materials and Methods](#) for more details). The model's preferences are therefore a measure of each individual codon's propensity to increase or decrease the model's output probability as a function of its position relative to the MCCs. A value of 0.5 denotes a neutral preference, while negative and positive preferences correspond to values below and above 0.5, respectively. Preferences were obtained by feeding CAMAP sequences in which all codon values were masked, except for a single position that received a non-null codon label.

Interestingly, while codons closest to the MCCs were the most influential on CAMAP scores, some synonymous codons showed opposite effects, demonstrating that codon usage does not entirely recapitulate amino acid usage (Figs 5A, 5B and S37). The use of embeddings to encode codons has the advantage of arranging them into a semantic space, wherein codons with similar influences are positioned close to each other. Interestingly, most synonymous



**Fig 5. CAMAP interpretation of codon impact on MAP biogenesis.** Preferences for a network trained on a context of 162 nucleotides (54 codons) for (A) serine, proline and tyrosine codons, and (B) leucine codons. (C) Learned codon embeddings. Some synonymous codons, such as those encoding for Isoleucine (I), Cysteine (C) or Arginine (R) are located far from one another, while others tend to cluster together (e.g. Proline [P] and Glutamic acid [E]).

<https://doi.org/10.1371/journal.pcbi.1009482.g005>

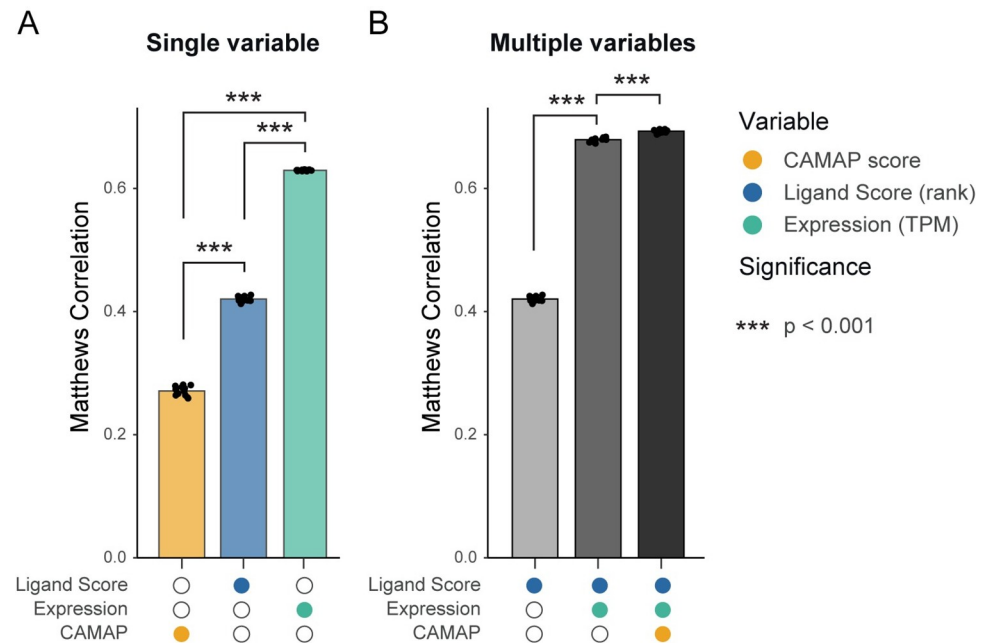
codons did not form clusters, with a notable exception being proline codons (Fig 5C). This finding indicates that for some codons, their effect on CAMAP prediction score may be closer to that of a non-synonymous codon than to that of one of its synonyms.

As expected, a significant part of the signal captured by CAMAP still originates from amino acids. As shown in Fig 2, the accuracy (AUC) of CAMAP trained with shuffled synonymous codons is closer to the AUC of CAMAP trained with original codons than it is to 0.5 (random score). In addition, we find that CAMAP has a higher preference for all codons encoding arginine (R) and lysine (K), indicating a strong tryptic-like specificity at the N-terminus (S37 Fig). CAMAP also showed a preference for alanine (A) codons both upstream and downstream of MAP coding codons (S37 Fig). This is in accordance with the enrichment of specific amino acids in the position flanking of MAP cleavage sites described by Abelin et al (2017). We also observe a lower preference score for proline (P), probably due to the conformational rigidity of this amino acid, and of acidic residue (glutamic acid [E] and aspartic acid [D]), also in accordance with Abelin et al (2017) (S37 Fig).

### Combination of binding affinity, expression levels and CAMAP scores increases MAP prediction accuracy

We next compared MAP prediction capacities of CAMAP scores to that of MAP predicted ligand score (ranks as predicted by NetMHCpan4.0) and mRNA transcript expression levels. We used ligand scores as predicted by NetMHCpan4.0, which was shown to possess the best predictive capacities for naturally processed peptides compared to other predictive algorithms [30]. Because MAP binding to the MHC molecule is essential for its presentation at the cell surface, we elected to only compare hits and decoys encoding potential binders, i.e. with a ligand score of  $\leq 1\%$  for at least one allele in the B-LCL dataset. Using a linear regression model, we compared the predictive capacity of each single parameter using Matthews correlation coefficient, which measures the quality of binary classifications [31]. Of note, only the predictions on the test set were used to evaluate the Matthew correlation coefficient in our different models, and the cutoff value for a prediction to be considered positive was set at 0.5.

Because only potential binders were analyzed here, the mRNA expression levels had the highest predictive capacity, then followed by ligand scores (second) and CAMAP scores (third, Fig 6A). As expected due to the multiplicative relationship between MAP ligand score and expression levels in predicting naturally processed MAPs [11], combining both variables greatly increased prediction performances (Fig 6B). Furthermore, we show that including CAMAP scores to the regression model further increased predictive performances (Fig 6B). We next computed how many predicted peptides would need to be tested to capture 1, 5, 10 or 50% of hits in the B-LCL dataset. Results presented in Table 1 show that using only NetMHCpan4.0 ligand scores (ranks) leads to a very high false positive rate (FPR) at 76.0% when targeting the top 1%. Adding the expression levels greatly increased prediction accuracy and decreased the FPR to 39.1% for the top 1% hits. When adding CAMAP scores as a third variable, the number of peptides needed to capture 1% of hits greatly decreased, resulting in a very low FPR at 8.9%. Similar trends were observed when targeting 5 or 10% of hits, although with higher FPR (see Table 1). Only when targeting 50% of hits does the performance decrease when adding the CAMAP score. Similarly, adding CAMAP scores to expression levels and ligand scores also ameliorated prediction accuracies for the two other human datasets introduced above (B721.221 and PBMCs, see S2 Table). We hypothesized that the variation in performance observed for the 50% case was due to the limitations of logistic regression, rather than to the data *per se*. We thus additionally assessed the value of each prediction method using a multi-layer perceptron (MLP) classifier. Results show that for all of them, adding



**Fig 6. CAMAP prediction score contributes to the prediction of MAPs.** (A) Matthews correlation coefficient for MAP prediction using a single variable. (B) Matthews correlation coefficient for MAP prediction using multivariable regression models. The B-LCL dataset (all MAP lengths) was filtered for MAP with a minimal ligand score (rank) of 1% (NetMHCpan4.0).

<https://doi.org/10.1371/journal.pcbi.1009482.g006>

CAMAP score to the input improves prediction accuracy (Table 1). Similar results were obtained for the two other human datasets (S3 Table). Interestingly, using a non-linear classifier such as an MLP ameliorated predictions for both the model receiving ligand score and transcript expression, and the model receiving ligand score, transcript expression and CAMAP scores. This suggests that the relationship between binding score, transcript expression and CAMAP score is a complex one. Therefore, implying the most efficient way to combine the three is through the use of a non-linear classifier. These results show that combining CAMAP scores with the MAP’s ligand score (ranks) and its corresponding transcript expression level significantly improves prediction of MAP and facilitates identification of relevant epitopes through more accurate predictions.

**Table 1. Number of peptides needed to capture 1%, 5%, 10% or 50% of epitopes detected by mass spectrometry.** The lower the number of peptides needed to capture the respective number of epitopes, the better the performance of the prediction model. This is also illustrated by the percentage of false identification (false positive rate, FPR) reported here. Peptides were rank ordered according to regression scores, for a total of 512,423 unique peptides and 8,807 hits. Of note, only the maximal transcript expression was used for peptides with multiple potential origins.

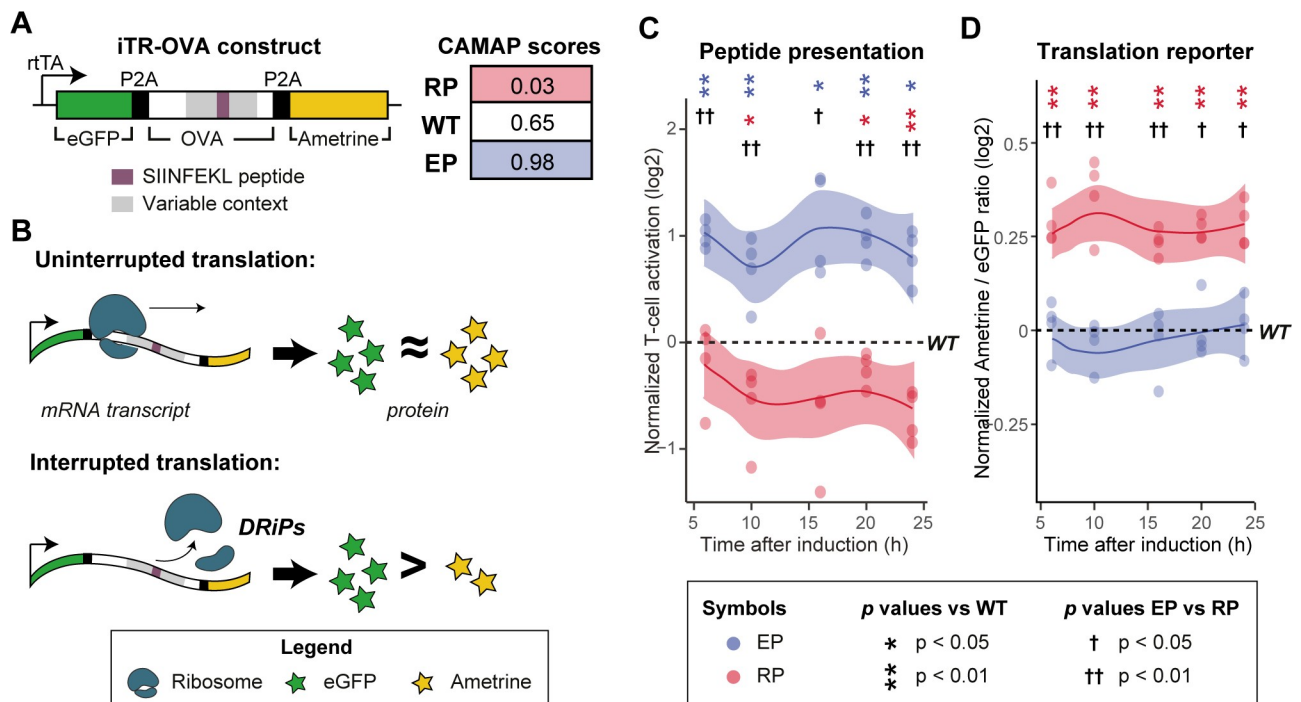
Linear regression	1% hits (n = 88)		5% hits (n = 440)		10% hits (n = 881)		50% hits (n = 4406)	
	n	FPR	n	FPR	n	FPR	n	FPR
NetMHCpan4.0	368 ± 20	76.0%	2,339 ± 70	81.1%	5,288 ± 164	83.38%	59,812 ± 878	92.6%
NetMHCpan4.0 + expression	145 ± 8	39.1%	630 ± 17	30.0%	1,264 ± 25	30.2%	11,734 ± 64	62.4%
NetMHCpan4.0 + expression + CAMAP	97 ± 5	8.9%	569 ± 17	22.4%	1,204 ± 24	26.8%	12,151 ± 171	63.7%
Multi-layer perceptron	1% hits (n = 88)		5% hits (n = 440)		10% hits (n = 881)		50% hits (n = 4406)	
	n	FPR	n	FPR	n	FPR	n	FPR
NetMHCpan4.0	360 ± 21	75.4%	2,289 ± 72	80.8%	5,239 ± 160	83.3%	59,659 ± 742	92.6%
NetMHCpan4.0 + expression	97 ± 4	8.8%	524 ± 9	16.0%	1,114 ± 12	20.9%	10,969 ± 113	59.9%
NetMHCpan4.0 + expression + CAMAP	92 ± 5	4.3%	506 ± 30	12.7%	1,067 ± 38	17.3%	10,397 ± 163	57.6%

<https://doi.org/10.1371/journal.pcbi.1009482.t001>

### Codon usage can modulate MAP presentation

To evaluate whether changing the codon arrangement in a MAP-coding sequence might directly lead to modulation of MAP presentation, we generated three variants of the chicken ovalbumin (OVA) protein containing the model MAP SIINFEKL [32]. One construct encoded the wild type OVA (OVA-WT). For the other two constructs, we used CAMAP (trained on original human B-LCL sequences; Fig 2) to generate two OVA variants *in silico*, both encoding for the same OVA protein but using different synonymous codons: one predicted to enhance SIINFEKL presentation (OVA-EP), the other predicted to reduce it (OVA-RP). Accordingly, the respective CAMAP scores for OVA-RP, OVA-WT and OVA-EP were: 0.03, 0.65, and 0.96 (Fig 7A). All variants encoded the same amino acid sequence but used different synonymous codons. Notably, the sole difference between the three constructs were the 162 nucleotides flanking each side of the SIINFEKL-coding codons (i.e. the RNA sequences coding for OVA<sub>202-256</sub> and OVA<sub>265-319</sub>, S1 Table and S38 Fig).

Because codon usage affects translation efficiency, theoretically leading to DRiP formation through premature translation arrest [20,21], we expected the variable regions of our construct to affect both translation rates and SIINFEKL presentation in our variants. Therefore, each construct also coded for two other proteins, eGFP and Ametrine, placed upstream and downstream of the OVA coding sequence, respectively (Fig 7A). While the Ametrine fluorescence



**Fig 7. Codon usage in MAP-flanking mRNA sequences can influence antigen presentation and translation efficiency.** (A) Design of the inducible Translation Reporter (iTR-OVA) constructs and CAMAP scores for OVA-WT, OVA-EP and OVA-RP sequences. (B) Schematic representation of possible translation events. When mRNA codon usage leads to efficient (uninterrupted) translation, similar amounts of eGFP and Ametrine proteins would be synthesized. When codon usage in the MAP-flanking regions enhances the frequency of translation interruption, a lower Ametrine/eGFP ratio would be observed. (C) Kinetics of SIINFEKL MAP presentation following induction of iTR-OVA constructs expression by doxycycline, measured in a T-cell activation assay. To remove the influence of differential expression levels on antigenic presentation and of varying proportion of transduced cells between samples, T-cell activation levels were normalized to the average Ametrine fluorescence intensity and to the proportion of eGFP+ cells (i.e. cells expressing the construct). (D) Translation efficiency as measured by Ametrine/eGFP ratio following iTR-OVA construct induction. For C and D, results are normalized over the WT sample from the same experiment (n = 4). Statistical differences at each time point were determined using bilateral paired Student T tests. Significance for the comparison against WT are indicated with \*, while comparison of EP vs RP is indicated with †. N.B.: Each replicate is shown with a dot, while the line and shaded area represent the average and 95% confidence interval, respectively.

<https://doi.org/10.1371/journal.pcbi.1009482.g007>

intensity reflected the translation rate of the whole construct, the ratio of Ametrine/eGFP fluorescence intensity was informative regarding the translation efficiency of the whole construct. Indeed, efficient translation of the full-length construct should produce equivalent quantities of Ametrine and eGFP proteins, while inefficient/interrupted translation of the construct (i.e. leading to DRiP formation) should decrease the Ametrine/eGFP ratio (Fig 7B). The three protein coding sequences were separated with P2A self-cleaving peptides [33], therefore allowing the co-synthesis of three separate proteins, controlled by the doxycycline-inducible Tet-On promoter. Importantly, the three proteins were tightly co-expressed because of the presence of only one start codon at the 5' end of the GFP protein, as shown by the very high correlation between eGFP and Ametrine fluorescence for each construct ( $R > 0.97$ , see S39 Fig). As we assumed that CAMAP scores reflected the probability of DRiP generation leading to increased MAP presentation, we expected the OVA-RP construct to show both reduced SIINFEKL presentation and enhanced translation efficiency compared to the OVA-EP and OVA-WT constructs. However, as both the OVA-EP and OVA-WT have CAMAP scores above the neutral threshold of 0.5 and closer to one another (0.98 and 0.65, respectively) compared to the OVA-RP construct (0.03), we expected OVA-EP and OVA-WT to behave more similarly.

We then used a SIINFEKL-H2-K<sup>b</sup> specific T-cell activation assay [34] to measure SIINFEKL presentation at the cell surface following doxycycline induction. Results for the T-cell activation assay were normalized by both the Ametrine mean fluorescence intensity and the percentage of transduced (eGFP+) cells in each specific sample, so that any difference in T-cell activation observed between our constructs could only be ascribed to synonymous codon variants in the SIINFEKL-flanking OVA codons. Two main findings emerged from our analyses. First, in accordance with CAMAP predictions, variation in codon usage led to a 2.3-fold difference in SIINFEKL presentation between the OVA-EP and OVA-RP variants, with OVA-WT in between (Fig 7C). Second, translation efficiency (Ametrine/eGFP ratio) was higher with OVA-RP than with OVA-EP or OVA-WT, while OVA-EP showed similar translation efficiency compared to OVA-WT (Fig 7D). Hence, synonymous codon variations led to slightly divergent outcomes in OVA-EP and OVA-RP: they modulated the levels of SIINFEKL presentation in both constructs, but enhanced translation efficiency could only be detected for OVA-RP. These data show that codon arrangement can modulate MAP presentation strength without any changes in the amino acid sequence and support a role for translation efficiency and DRiP formation in the modulation of MAP presentation.

## Discussion

Our analyses of large datasets using artificial neural networks and other bioinformatics approaches provide compelling evidence that codon usage regulates MAP biogenesis via both short- and long-range effects. While most MAP predictive approaches focus on MAP sequences *per se*, CAMAP's novelty is that it only receives the MAP-flanking mRNA sequences as input, and no information on the MAP itself, thereby providing completely independent information for MAP prediction.

The better prediction accuracy of CAMAPs trained with original codons rather than with shuffled synonyms supports the role of codon usage in modulating MAP biogenesis (Fig 2). A large body of evidence suggests that a significant proportion of MAPs derive from DRiPs produced during mRNA-to-protein translation [13–15]. Because codons influence both precision and efficiency of protein synthesis [20,21], it is likely that codon features in the mRNA sequences flanking MCCs might contribute to the regulation of MAP biogenesis through the generation of DRiPs. The functional link between codon arrangement and MAP biogenesis was illustrated by our *in vitro* analyses of SIINFEKL biogenesis, in which we were able to



modulate SIINFEKL presentation solely by substituting synonymous codons in mRNA regions flanking SIINFEKL codons, without changing the protein sequence (Fig 7). While the experimental data derives from a single model thus limiting the interpretability of our results, this points nonetheless to an interesting mechanism that could be exploited to enhance antigenic presentation in peptide-based immunotherapy (i.e. dendritic cells modified to express a specific MAP). The full extent of the contribution of CAMAP's learned codon features to DRiPs generation will need to be assessed in subsequent work.

Previous works have also demonstrated the role of transcript expression levels in the regulation of MAP biogenesis, with a higher transcript expression being associated with a higher likelihood of MAP presentation [11,19]. While codon usage is different in highly expressed genes, our results showed that the codon-specific signal captured by CAMAP is independent of transcript expression levels, thereby providing complementary and non-overlapping information regarding MAP presentation. Interestingly, CAMAP preferences were more influential for codons located close to the MCCs (Fig 5), which is consistent with the influence of adjacent amino acids on proteasome cleavage. However, the better performance of CAMAP trained with longer context size also points toward a long-range impact of codon usage on MAP presentation.

Further analyses will be needed to assess the full extent of codon arrangement's impact on both classic MAPs (i.e. derived from canonical reading frames of coding sequences) and cryptic MAPs (i.e. derived from non-canonical reading frames and non-coding sequences) [35,36], as well as the potential contribution of codons in non-coding regions (e.g. 5'- or 3'-UTRs) on the regulation of MAP presentation. However, our results show that the integration of CAMAP scores to the two best predictive factors for naturally processed MAPs led to a significant increase in prediction accuracy. Indeed, our MLP classifier showed that combining CAMAP scores to transcript expression levels and MAP ligand scores (ranks as predicted by NetMHCpan4.0), led to a lower FPR compared to a model combining only transcript expression levels and MAP ligand scores. Although predictions were not as accurate for the two other human datasets, adding CAMAP scores always resulted in improved prediction accuracy. Our results therefore support the combined use of ligand scores, transcript expression levels and CAMAP scores in MAP predictive algorithms. These results have important practical implications for cancer immunotherapy and peptide-based vaccines, where discovery of suitable target antigens remains a formidable challenge to this day [37,38].

## Materials and methods

### Dataset generation

We analyzed a previously published dataset consisting of MAPs presented on B lymphocytes by a total of 33 MHC-I alleles from 18 subjects [19,22]. Since this dataset was assembled using older versions of MHC-I binding prediction algorithms (i.e. using a combination of NetMHC3.4 for common alleles and NetMHCcons1.1 for rare alleles), we verified that the majority of MAPs in this dataset would also be predicted as binders using more recent algorithms (i.e. a rank  $\leq 2.0\%$  using NetMHC4.0 or NetMHCpan4.0). We found an overlap of  $>92\%$  between these methods (see S40 Fig), thereby validating this dataset for further analysis. In addition, we reasoned that a transcript should be considered as a genuine positive or negative regarding MAP biogenesis only if it was expressed in the cells. We therefore excluded from the dataset all transcripts with very low expression ( $<1^{\text{st}}$  percentile in terms of FPKM).

To facilitate data analysis and interpretation, we only included transcripts coding for MAPs with a length of 9 amino acids, for a total of 19,656 9-mer MAPs (which represents 78% of MAPs in this dataset). We then used pyGeno [23] to extract the mRNA sequences of

transcripts coding for these 9-mer MAPs, which constituted our source-transcripts (Fig 1A). We next created a negative (non-source) dataset from transcripts that generated no MAPs. Importantly, transcripts that encoded for MAPs of any length (i.e. 8 to 11-mer) were excluded from the negative dataset. We then randomly selected 98,290 non-MAP 9-mers from this negative dataset, and extracted their coding sequences using pyGeno. Of note, both positive and negative datasets were derived from the canonical reading frame of non-redundant transcripts.

We analyzed only the MAP context and excluded the MCCs *per se* from our positive (hits) and negative (decoys) sequences (Fig 1A). We limited our analyses of flanking sequences to 162 nucleotides (54 codons) on each side of MCCs, because longer lengths would entail the exclusion of >25% of transcripts (S41 Fig).

### Creation of the shuffled synonymous codon dataset

To create the shuffled synonymous codon dataset, each sequence was re-encoded by replacing each codon with itself or with a random synonym according to the human transcriptome usage frequencies. These frequencies were calculated using the annotations provided by *Ensembl* for the human reference genome GRCh37.75. Thus, all codon-specific features differing between the positive and negative datasets were removed from the shuffled datasets. Because codons were replaced by their synonymous codons, the shuffled sequences directly reflected amino acid usage in the positive and negative datasets.

### CAMAP architecture, sequence encoding and training

The first (input) layer received either MCCs flanking regions from the hit dataset or sequences of the same length contained in the decoy dataset (Fig 1A). The second layer (S21A Fig) was a codon embedding layer similar to that introduced for a neural language model [39]. Embedding is a technique used in natural language processing to encode discrete words, and has been shown to greatly improve performances [40]. With this technique, the user defines a fixed number of dimensions in which words should be encoded. When the training starts, each word receives a random vector-valued position (its embedding coordinates) in that space. The network then iteratively adjusts the words' embedding vectors during the training phase and arranges them in a way that optimizes the classification task. Notably, embeddings have been shown to represent semantic spaces in which words of similar meanings are arranged close to each other [40]. In the present work, we treated codons as words: each codon received a set of random 2D coordinates that were subsequently optimized during training. The third (output) layer delivered the probability that the input sequence was a MCCs flanking region (rather than a sequence from the negative dataset). Therefore, CAMAP only has two levels of processing, the embedding layer and the output layer. Finally, the output layer has no bias parameter, a feature that enabled us to easily compute the preferences.

CAMAPs were trained on sequences resulting from the concatenation of pre- and post-MCCs regions. Before presenting sequences to our CAMAPs, we associated each codon to a unique number ranging from 1 to 64 (we reserved 0 to indicate a null value) and used this encoding to transform every sequence into a vector of integers representing codons. Neural networks were first developed using the Python package Mariana [41] (<https://www.github.com/tariqdaouda/Mariana>) and reimplemented identically for this work using pyTorch (results presented herein derive from the pyTorch version). An *Embedding* layer was used to associate each label superior to 0 to a set of 2D trainable parameters; the 0 label represents a *null* (masking) embedding fixed at coordinates (0,0). As an output layer, we used a *Softmax* layer with two outputs (positive / negative). Because negative sequences are more numerous

than positive ones, we used an oversampling strategy during training. At each epoch, CAMAPs were randomly presented with the same number of positive and negative sequences. All CAMAPs in this work share the same architecture (S21A Fig), number of parameters and hyper-parameter values: learning rate: 0.001; mini-batch size: 64; embedding dimensions: 2; linear output without offset on the embedding layer; *Softmax* non-linearity without offset on the output layer.

For each condition (e.g. context size), the positive and negative datasets were randomly divided into three non-redundant subsets: (i) the training subsets containing 60% of the positive and negative examples, (ii) the validation and (iii) the test subsets each containing 20% of the positive and negative examples. Examples were assigned through a sequence redundancy removal algorithm, thereby ensuring that no example was assigned to multiple subsets. We used an early stopping strategy on validation sets to prevent over-fitting and reported average performances computed on test sets. We trained 12 CAMAPs for each combination of conditions, each one using a different random split of train/validation/test sets. To mask sequences either before or after the MCCs, we masked either half with *null* value.

### Kullback-Leibler divergence

The Kullback-Leibler (KL) divergence computes how well a given distribution is approximated by another distribution. Its value can be either positive or 0, a null value indicating that the two distributions are identical. Accordingly, a higher KL divergence for codon distributions vs. amino acid distributions would indicate that codon variations are not entirely accounted for by amino acid variations. KL divergence is not a metric, as it is neither symmetric nor does it satisfy the triangle inequality. It is nevertheless an accurate and most common way of comparing two probability distributions.

We defined the probability of having codon  $c$  at position  $i$  as a function of the number of occurrences of  $c$  at position  $i$ , divided by the total number of occurrences of that same codon:

$$Q_{(c,y,s)}(i) = \frac{N_{c,y,s}(i)}{\sum_j N_{c,y,s}(j)}$$

Here  $Q$  is a probability,  $N$  is a number of occurrences,  $c$  is a codon,  $y$  is a class (positive or negative),  $s$  indicates if codons have been shuffled (true or false),  $i$  is a position in sequence. For the remainder of the text we will use the following abbreviations:

$$P_c(i) = Q_{c,y=positive,s=false}(i)$$

$$D_c(i) = Q_{c,y=negative,s=false}(i)$$

$$PS_c(i) = Q_{c,y=positive,s=true}(i)$$

$$DS_c(i) = Q_{c,y=negative,s=true}(i)$$

We then used the KL divergence to compute how well  $P_c$  distributions approximate  $D_c$  distributions and  $PS_c$  distributions approximate  $DS_c$  distributions.

The KL divergence was defined as:

$$D_{KL}(P||Q) = \sum_i P(i) \log\left(\frac{P(i)}{Q(i)}\right)$$

We performed this calculation for both the original and the shuffled dataset, which we then compared together. If codons and amino acid distributions were equivalent, KL divergence between hits and decoys would be the same for both original and shuffled sequences, and codons would cluster along the diagonal.

### Extracting CAMAP's preferences

Preferences are a way of interpreting CAMAP networks that naturally derives from CAMAP's mathematical formalism. The validity of the preference method for CAMAP has been demonstrated when we used it to define the EP and RP constructs. By following this method, we were able to precisely engineer codon sequences that maximized or minimized the predictions of a CAMAP network. CAMAP only has two processing layers: an embedding layer and an output layer with no latent processing in between. Moreover, the output layer has no bias. In these very specific conditions, the manifold is the embedding layer. The word embedding layer encodes the importance of codons, while the weights of the output layer encode the importance of the position.

The output is thus computed as:

$$\text{Output} = E \times W$$

Where  $E$  is the vector of embeddings for a sequence and  $W$  is the weight matrix of the output layer. When masking all codons but one, we obtain the importance ( $I$ ) of that codon ( $c$ ) at that position ( $x$ ):

$$I_{c,x} = E_{c,x} \times W_x$$

Where  $W_x$  are the weights of the output layer for position  $x$ . The preference is just this quantity ( $I_{c,x}$ ) normalized to a probability (using a Softmax function) for ease of interpretation.

### Predicting MAP presentation with logistic regressions and MLP

The prediction capacity of CAMAP, NetMHCpan-4.0 ligand score and transcription expression (TPM) was tested in different combinations of those parameters (Ligand Score + Expression, Ligand score + Expression + CAMAP score) using the *LogisticRegressionCV* or *MLPClassifier* functions from the python package *sklearn* (*sklearn.linear\_model* and *sklearn.neural\_network*, v0.24.1). In each case, the dataset containing hits and decoy sequences was split into train and test datasets with a ratio of 0.7 to 0.3, respectively. Values for CAMAP score, Ligand Score and TPM were each scaled to a range of 0–1 in the train set using *MinMaxScaler* from *sklearn.preprocessing* and the same scaling model was applied to the test set afterwards. Regression analysis was performed using *LogisticRegressionCV* with a 10x cross-validation using the *lbfgs* solver with 1000 iterations. MLP with one hidden layer of 1000 neurons was trained using *MLPClassifier* for 4000 iterations using the *adam* solver with early stopping. Classes were balanced prior to training for both models. Matthew correlation coefficients were calculated using *matthews\_corrcoef* from *sklearn.metrics*. When a peptide had multiple sources (multiple transcripts or genes), only the maximum transcript expression is used.

### In vitro assay-inducible translation reporter (iTR)-OVA construct design

An inducible translation reporter was generated by flanking the truncated chicken ovalbumin (OVA) cDNA (amino acids 144–386) with EGFP-P2A (in 5') and P2A-Ametrine (in 3') cDNA sequences. MCCs flanking contexts for the EP and RP construct were synthesized as gBlocks (purchased from Integrated DNA Technologies). The fragments were amplified by PCR and

joined by Gibson assembly under a doxycycline-inducible Tet-ON promoter in a pCW backbone. Synthetic variants of the OVA coding sequence were generated in silico by varying synonymous codon usage in the MAP context regions (i.e. 162 nucleotides pre- and post-MCCs). Importantly, the amino acid sequence was preserved between the different variants; only nucleotide sequences in the MAP context (162 nucleotides on either side) differed. The sequences with the highest (EP) and the lowest (RP) prediction scores were selected for further in vitro validation and swapped into the iTR-OVA plasmid by Gibson assembly [42]. OVA-EP and OVA-RP sequences can be found in [S1 Table](#).

Important features of our inducible translation reporter construct and T cell activation assay were: (i) No changes in amino acid sequence between the three variants: only co-translational events can differ between the three variants, post-translational events being equivalent for the three constructs; (ii) Only one start codon, at the beginning of the eGFP coding sequence: this is important for the translation reporter aspect of our construct (i.e. Ametrine/eGFP ratio), to ensure that translation can only start at the 5'-end of the whole construct, and not at the beginning of the OVA or Ametrine coding sequences; (iii) Separation of the three proteins using P2A peptide: allows the inducible synthesis of three separate proteins in a highly correlated manner (as demonstrated in [S39 Fig](#)); also, the degradation of one protein will be independent from the others. As we hypothesized that codon usage might lead to DRiP formation, we did not want the degradation of OVA-derived polypeptide to induce degradation of attached eGFP or Ametrine, which would affect our translation reporter assay (Ametrine/eGFP ratio); (iv) Because transcript expression level impacts MAP presentation, we normalized T-cell activation results by both the number of transduced cells present in the samples (% of eGFP+ cells) and the Ametrine mean fluorescence intensity of eGFP+ cells (representing whole construct expression level). Because of these four features, any difference between the three constructs could be ascribed solely to synonymous codon variants in the SIINFEKL-flanking OVA codons.

### Stable cell line generation

Wildtype and transduced Raw-K<sup>b</sup> cells [43] were cultured in DMEM supplemented with 10% Fetal Bovine Serum (FBS), penicillin (100 units/ml), and streptomycin (100mg/ml). B3Z cells [44] were maintained in RPMI medium supplemented with 5% FBS, penicillin (100 units/ml), and streptomycin (100mg/ml).

Lentiviral particles were produced from HEK293T cells by co-transfection of iTR-OVA WT, EP or RP along with pMD2-VSVG, pMDLg/pRRE and pRSV-REV plasmids. Viral supernatants were used for Raw-K<sup>b</sup> transduction. Raw-K<sup>b</sup> OVA-WT, Raw-K<sup>b</sup> OVA-EP were sorted on Ametrine and GFP double positive population after 24h of doxycycline treatment (1 mg/ml).

### T-cell activation assay

Raw-K<sup>b</sup> OVA-EP, OVA-RP and OVA-WT cells were plated at a density of 250,000 cells/well in 24 well-plates 24h prior to doxycycline treatment (1 mg/ml). After the corresponding treatment duration, cells were harvested and fixed using PFA 1% for 10 minutes at room temperature and washed using DMEM 10% FBS. Raw-K<sup>b</sup> were then co-cultured (37°C, 5% CO<sub>2</sub>) in triplicates with the CD8 T cell hybridoma cell line B3Z cells at a 3:2 ratio for 16h (7.5 x 10<sup>5</sup> B3Z and 5 x 10<sup>5</sup> Raw-K<sup>b</sup>) in 96 well-plates. Cells were lysed for 20 minutes at room temperature using 50 µl/well of lysis solution (25mM Tris-Base, 0.2 mM CDTA, 10% glycerol, 0.5% Triton X-100, 0.3mM DTT; pH 7.8). 170 µl/well CPRG buffer was added (0.15mM chlorophenol red-β-d-galactopyranoside (Roche), 50mM Na<sub>2</sub>HPO<sub>4</sub>•7H<sub>2</sub>O, 35mM NaH<sub>2</sub>PO<sub>4</sub>•H<sub>2</sub>O, 9mM KCl, 0.9mM MgSO<sub>4</sub>•7H<sub>2</sub>O). β-galactosidase activity was measured at 575 nm using SpectraMax



190 Microplate Reader (Molecular Devices). In parallel, cells were analyzed by flow cytometry using a BD FACS CantoII for eGFP and Ametrine fluorescence.

## Supporting information

**S1 Fig. Codon distribution in the shuffled datasets more closely resembles that of amino acids, compared to the original datasets.** (A) Pearson correlation ( $R^2$ ) factors and (b) Kullback-Leibler (KL) divergence between positional distribution of codons and their corresponding amino acid in the shuffled (y axis) VS original (x axis) datasets. For all codons, the shuffled dataset showed greater correlations (A) and smaller KL divergence to their respective amino acid distributions than the original datasets ( $p < 1 \times 10^{-8}$ , assessed using unilateral paired Student T test).

(TIF)

**S2 Fig. Distribution of Pearson's correlation factors calculated between codons and amino acids positional distributions in the original (green) and shuffled (coral) datasets.** 92% of codons in the shuffled dataset reflecting the amino acids distribution with a  $R^2 > 0.95$ , compared to only 69% in the original dataset ( $p < 5 \times 10^{-5}$ ).

(TIF)

**S3 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Alanine–A.**

(EPS)

**S4 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Cysteine–C.**

(EPS)

**S5 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Aspartic acid–D.**

(EPS)

**S6 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Glutamic acid–E.**

(EPS)

**S7 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Phenylalanine–F.**

(EPS)

**S8 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Glycine–G.**

(EPS)

**S9 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Histidine–H.**

(EPS)

**S10 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Isoleucine–I.**

(EPS)

**S11 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Lysine–K.**

(EPS)

**S12 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets Leucine-L.**

(EPS)

**S13 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Asparagine-N.**

(EPS)

**S14 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Proline-P.**

(EPS)

**S15 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Glutamine-Q.**

(EPS)

**S16 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Arginine-R.**

(EPS)

**S17 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Serine-S.**

(EPS)

**S18 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Threonine-T.**

(EPS)

**S19 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Valine-V.**

(EPS)

**S20 Fig. Distribution of amino acid and codon usage per position in the original VS shuffled datasets for Tyrosine-Y.**

(EPS)

**S21 Fig. CAMAP architecture and detailed predictions.** (A) Architecture of the ANN used in this work. (B) Results for the AUC on all train, validation and test subsets. Grey areas represent the 95% confidence intervals. (C) Distributions of output probabilities of CAMAPs used to calculate correlations in [S22 Fig](#).

(TIF)

**S22 Fig. Correlation between CAMAP prediction score trained only with pre-MCC or post-MCC sequences.** For each sequence in the test set we calculated the average prediction score given by CAMAPs in each condition, and calculated the Pearson correlation using the R software. Densities were calculated on all points and drawn using ggplot2. Only a random subset of the points is represented in the figures to limit their size.

(TIF)

**S23 Fig. Absence of correlation between CAMAP prediction score and transcript expression levels in 4 individual B-LCL samples (each derived from a different subject).**

(TIF)

**S24 Fig. Training of CAMAP on dataset selected to reflect positive dataset's distribution in expression levels.** (A) Distribution of transcript expression levels for normal datasets (related to Fig 2) and the dataset used here to retrain CAMAP. As shown in this figure, the decoy dataset was selected to mirror the distribution of transcript expression level in the hit dataset. (B) CAMAP performance (measured by the AUC) when trained using the decoy dataset that mirrors the transcript expression levels of the hit dataset. Significance was assessed using bilateral paired Student T test ( $p = 5.36 \times 10^{-7}$ ).  
(TIF)

**S25 Fig. Absence of correlation between CAMAP prediction score and binding affinities for individual alleles for decoys.**  
(TIF)

**S26 Fig. Absence of correlation between CAMAP prediction score and binding affinities for individual alleles for hits.**  
(TIF)

**S27 Fig. Training of CAMAP on dataset selected to reflect positive dataset's distribution in binding affinities.** (A) Distribution of binding affinities for normal dataset (related to Fig 2) and the corrected dataset used to retrain CAMAP. As shown in this figure, the decoy dataset was selected to mirror the distribution of binding affinities in the hit dataset. (B) CAMAP performance (measured by the AUC) when trained using the decoy dataset that mirrors the binding affinities of the hit dataset. Significance was assessed using bilateral paired Student T test ( $p = 1.21 \times 10^{-9}$ ).  
(TIF)

**S28 Fig. Evaluation of homology in hit dataset and its impact on CAMAP performance.** (A) Proportion of unique MAPs that can be ascribed to a single origin, 2–3, 4–10 or >10 possible origins. (B) Proportion of entries in the hit dataset that encode for MAPs with a single origin, 2–3, 4–10 or >10 possible origins  
(TIF)

**S29 Fig. Gene families overrepresented in hits with >3 possible origins.**  
(TIF)

**S30 Fig. CAMAP performance (AUC) when trained using either all hits (left), hits with 10 possible origins or less (center) or hits with 3 possible origins or less (right).**  
(TIF)

**S31 Fig. Correlation between CAMAP prediction score and (A) transcripts' GC content (%) and (B) tRNA gene copy number.** Pearson's R correlation score is shown on the graphs.  
(TIF)

**S32 Fig. Training of CAMAP on dataset selected to reflect positive dataset's distribution in GC content.** (A) Distribution of transcripts' GC content for normal dataset (related to Fig 2) and the corrected dataset used to retrain CAMAP. As shown in this figure, the decoy dataset was selected to mirror the distribution of GC content in the hit dataset. (B) CAMAP performance (measured by the AUC) when trained using the decoy dataset that mirrors the GC content of the hit dataset. Significance was assessed using bilateral paired Student T test ( $p = 1.07 \times 10^{-4}$ ).  
(TIF)

**S33 Fig. CAMAP scores for hits, for decoys derived from source transcripts and for decoys derived from non-source transcripts.** CAMAP was trained using hits from source transcripts

and decoys from non-source transcripts only. Significance was assessed using bilateral unpaired T test.

(TIF)

**S34 Fig. Shuffling methods.** (A) Transcriptome shuffling: generates shuffled sequences by replacing each codon by one of its synonymous codons (including itself) according to the codon usage within the transcriptome, regardless of the codon used within each example. (B) Codon swap shuffling: synonymous codons present within a given example are swapped with one another. (C) Third nucleotide shuffling or N3: the third nucleotide of codons are swapped within each example to preserve amino acid sequences and GC content. Here, the global codon usage is completely different than normal codon usage in humans, as the frequency of codons is not taken into account during shuffling.

(TIF)

**S35 Fig. CAMAP performances for different shuffled datasets.** The performance of CAMAP networks ( $n = 12$ ) pre-trained on original (non-shuffled) datasets was evaluated on three shuffled datasets according to the methods depicted in [S34 Fig](#). Significance was assessed using a paired bilateral T test.

(TIF)

**S36 Fig. Kullback-Leibler divergence between hit and decoy datasets in original codon (y-axis) or shuffled synonymous codon sequences (x-axis).** Shuffled sequences represent amino acid usage, as codon-specific information are removed with synonymous codon shuffling.

(TIF)

**S37 Fig. Preferences per position for all codons for CAMAP trained with original sequences.** See [Materials and Methods](#) for more details.

(TIF)

**S38 Fig. OVA-construct alignment, showing point mutations (red lines) in the mRNA sequences flanking the SIINFEKL MCC.** (A) Comparison of the OVA-EP nucleotide sequence to the wildtype OVA sequence. The OVA-EP and OVA-WT sequences have 93.3% nucleotide identity for a total of 78 modified nucleotides. (B) Comparison of the OVA-RP nucleotide sequence to the wildtype OVA sequence. The OVA-EP and OVA-WT sequences have 92.6% nucleotide identity, for a total of 86 modified nucleotides. Mutations, shown in red, are located only in the 162 nucleotide regions flanking the SIINFEKL coding codons. Of note, the SIINFEKL coding codons (nucleotides 772–799) were not modified between the 3 constructs.

(TIF)

**S39 Fig. Correlations between eGFP and Ametrine fluorescence intensity at the single cell level.** Single cell eGFP and Ametrine fluorescence intensities measured at 10 hours post-induction are shown for the OVA-WT (A), OVA-EP (B) and OVA-RP (C) constructs. N.B.: only transduced cells are shown (eGFP+ cells).

(TIF)

**S40 Fig.** Validation of MHC-I associated peptides (MAP) dataset from Pearson H. *et al.* (2016) using the new versions of MAP binding affinity prediction algorithm NetMHC4.0 (A) and NetMHCpan4.0 (B).

(TIF)

**S41 Fig. Percentage of transcript ineligibility as a function of context size.** Transcript length corresponds to  $C \times 2 + 27$ , where  $C$  is the context size in nucleotides and 27 the length of the

MCCs. Related to [Fig 1A](#).  
(TIF)

**S1 Table. Nucleotide sequences of the EP and RP constructs.** SIINFEKL MCCs are shown in bold, while the variant regions (pre- and post-MCCs flanking sequences, context size of 162-nucleotides) are in blue and italics. Related to [Fig 7](#).

(DOCX)

**S2 Table. Number of peptides needed to capture 1%, 5, 10 and 50% of epitopes detected by mass spectrometry in B721.221 and PBMC cell lines using a logistic regression model.** The lower the number of peptides needed to capture the respective number of epitopes, the better the performance of the prediction model. This is also illustrated by the percentage of false identification (false positive rate, FPR) reported here. Peptides were rank ordered according to regression scores. Of note, only the maximal transcript expression was used for peptides with multiple potential origins.

(DOCX)

**S3 Table. Number of peptides needed to capture 1%, 5, 10 and 50% of epitopes detected by mass spectrometry in B721.221 and PBMC cell lines using a multi-layer perceptron (MLP).** The lower the number of peptides needed to capture the respective number of epitopes, the better the performance of the prediction model. This is also illustrated by the percentage of false identification (false positive rate, FPR) reported here. The non-linear classifier led to better predictions when using ligand score and transcript expression, with or without CAMAP scores compared to logistic regression model. Peptides were rank ordered according to regression scores. Of note, only the maximal transcript expression was used for peptides with multiple potential origins.

(DOCX)

## Acknowledgments

The B3Z CD8+ T cell hybridoma cell line was a kind gift from Nilabh Shastri. RAW-K<sub>b</sub> cells were kindly provided by Michel Desjardins.

## Author Contributions

**Conceptualization:** Tariq Daouda, Maude Dumont-Lagacé, Yahya Benslimane, Sébastien Lemieux, Claude Perreault.

**Data curation:** Tariq Daouda, Albert Feghaly, Claude Perreault.

**Formal analysis:** Tariq Daouda, Maude Dumont-Lagacé, Albert Feghaly, Yahya Benslimane, Rébecca Panes, Mathieu Courcelles, Mohamed Benhammadi, Sébastien Lemieux, Claude Perreault.

**Funding acquisition:** Lea Harrington, Étienne Gagnon, Sébastien Lemieux, Claude Perreault.

**Investigation:** Tariq Daouda, Maude Dumont-Lagacé, Albert Feghaly, Yahya Benslimane, Mathieu Courcelles, Mohamed Benhammadi, Sébastien Lemieux, Claude Perreault.

**Methodology:** Tariq Daouda, Maude Dumont-Lagacé, Yahya Benslimane, Rébecca Panes, Yoshua Bengio, Étienne Gagnon, Sébastien Lemieux, Claude Perreault.

**Project administration:** Sébastien Lemieux, Claude Perreault.



**Resources:** Lea Harrington, Pierre Thibault, François Major, Étienne Gagnon, Sébastien Lemieux, Claude Perreault.

**Software:** Tariq Daouda, Albert Feghaly, Sébastien Lemieux.

**Supervision:** Lea Harrington, Pierre Thibault, François Major, Étienne Gagnon, Sébastien Lemieux, Claude Perreault.

**Validation:** Tariq Daouda, Albert Feghaly, Rébecca Panes, Yoshua Bengio, Sébastien Lemieux, Claude Perreault.

**Visualization:** Tariq Daouda, Maude Dumont-Lagacé, Albert Feghaly, Sébastien Lemieux, Claude Perreault.

**Writing – original draft:** Tariq Daouda, Maude Dumont-Lagacé.

**Writing – review & editing:** Albert Feghaly, Yahya Benslimane, Rébecca Panes, Mathieu Courcelles, Mohamed Benhammedi, Lea Harrington, Pierre Thibault, François Major, Yoshua Bengio, Étienne Gagnon, Sébastien Lemieux, Claude Perreault.

## References

1. Granados DP, Laumont CM, Thibault P, Perreault C. The nature of self for T cells—a systems-level perspective. *Curr Opin Immunol*. 2015; 34: 1–8. <https://doi.org/10.1016/j.coi.2014.10.012> PMID: 25466393
2. Caron E, Espona L, Kowalewski DJ, Schuster H, Ternette N, Alpizar A, et al. An open-source computational and data resource to analyze digital maps of immunopeptidomes. Chakraborty AK, editor. *eLife*. 2015; 4: e07661. <https://doi.org/10.7554/eLife.07661> PMID: 26154972
3. Davis MM, Krogsgaard M, Huse M, Huppa J, Lillemeier BF, Li Q. T Cells as a Self-Referential, Sensory Organ. *Annu Rev Immunol*. 2007; 25: 681–695. <https://doi.org/10.1146/annurev.immunol.24.021605.090600> PMID: 17291190
4. Schumacher TN, Schreiber RD. Neoantigens in cancer immunotherapy. *Science*. 2015; 348: 69–74. <https://doi.org/10.1126/science.aaa4971> PMID: 25838375
5. Caron E, Vincent K, Fortier M-H, Laverdure J-P, Bramoullé A, Hardy M-P, et al. The MHC I immunopeptidome conveys to the cell surface an integrative view of cellular regulation. *Mol Syst Biol*. 2011; 7: 533. <https://doi.org/10.1038/msb.2011.68> PMID: 21952136
6. Neefjes J, Jongsma MLM, Paul P, Bakke O. Towards a systems understanding of MHC class I and MHC class II antigen presentation. *Nat Rev Immunol*. 2011; 11: 823–836. <https://doi.org/10.1038/nri3084> PMID: 22076556
7. Bassani-Sternberg M, Gfeller D. Unsupervised HLA Peptidome Deconvolution Improves Ligand Prediction Accuracy and Predicts Cooperative Effects in Peptide–HLA Interactions. *J Immunol*. 2016; 197: 2492–2499. <https://doi.org/10.4049/jimmunol.1600808> PMID: 27511729
8. Nielsen M, Andreatta M. NetMHCpan-3.0; improved prediction of binding to MHC class I molecules integrating information from multiple receptor and peptide length datasets. *Genome Med*. 2016; 8: 33. <https://doi.org/10.1186/s13073-016-0288-x> PMID: 27029192
9. Tenzer S, Peters B, Bulik S, Schoor O, Lemmel C, Schatz MM, et al. Modeling the MHC class I pathway by combining predictions of proteasomal cleavage, TAP transport and MHC class I binding. *Cell Mol Life Sci CMLS*. 2005; 62: 1025–1037. <https://doi.org/10.1007/s00018-005-4528-2> PMID: 15868101
10. Nielsen M, Lundegaard C, Lund O, Keşmir C. The role of the proteasome in generating cytotoxic T-cell epitopes: insights obtained from improved predictions of proteasomal cleavage. *Immunogenetics*. 2005; 57: 33–41. <https://doi.org/10.1007/s00251-005-0781-7> PMID: 15744535
11. Abelin JG, Keskin DB, Sarkizova S, Hartigan CR, Zhang W, Sidney J, et al. Mass Spectrometry Profiling of HLA-Associated Peptidomes in Mono-allelic Cells Enables More Accurate Epitope Prediction. *Immunity*. 2017; 46: 315–326. <https://doi.org/10.1016/j.immuni.2017.02.007> PMID: 28228285
12. Capietto A-H, Jhunjunwala S, Delamarre L. Characterizing neoantigens for personalized cancer immunotherapy. *Curr Opin Immunol*. 2017; 46: 58–65. <https://doi.org/10.1016/j.coi.2017.04.007> PMID: 28478383
13. Antón LC, Yewdell JW. Translating DRiPs: MHC class I immunosurveillance of pathogens and tumors. *J Leukoc Biol*. 2014; 95: 551–562. <https://doi.org/10.1189/jlb.1113599> PMID: 24532645

14. Wei J, Kishton RJ, Angel M, Conn CS, Dalla-Venezia N, Marcel V, et al. Ribosomal Proteins Regulate MHC Class I Peptide Generation for Immunosurveillance. *Mol Cell*. 2019; 73: 1162–1173.e5. <https://doi.org/10.1016/j.molcel.2018.12.020> PMID: 30712990
15. Yewdell JW, Antón LC, Bennink JR. Defective ribosomal products (DRiPs): a major source of antigenic peptides for MHC class I molecules? *J Immunol*. 1996; 157: 1823–1826. PMID: 8757297
16. Croft NP, Smith SA, Wong YC, Tan CT, Dudek NL, Flesch IEA, et al. Kinetics of Antigen Expression and Epitope Presentation during Virus Infection. *PLOS Pathog*. 2013; 9: e1003129. <https://doi.org/10.1371/journal.ppat.1003129> PMID: 23382674
17. Milner E, Barnea E, Beer I, Admon A. The Turnover Kinetics of Major Histocompatibility Complex Peptides of Human Cancer Cells\*. *Mol Cell Proteomics*. 2006; 5: 357–365. <https://doi.org/10.1074/mcp.M500241-MCP200> PMID: 16272561
18. Hassan C, Kester MGD, de Ru AH, Hombrink P, Drijfhout JW, Nijveen H, et al. The Human Leukocyte Antigen-presented Ligandome of B Lymphocytes\*. *Mol Cell Proteomics*. 2013; 12: 1829–1843. <https://doi.org/10.1074/mcp.M112.024810> PMID: 23481700
19. Pearson H, Daouda T, Granados DP, Durette C, Bonneil E, Courcelles M, et al. MHC class I-associated peptides derive from selective regions of the human genome. *J Clin Invest*. 2016; 126: 4690–4701. <https://doi.org/10.1172/JCI88590> PMID: 27841757
20. Cannarozzi G, Schraudolph NN, Faty M, von Rohr P, Friberg MT, Roth AC, et al. A Role for Codon Order in Translation Dynamics. *Cell*. 2010; 141: 355–367. <https://doi.org/10.1016/j.cell.2010.02.036> PMID: 20403329
21. Plotkin JB, Kudla G. Synonymous but not the same: the causes and consequences of codon bias. *Nat Rev Genet*. 2011; 12: 32–42. <https://doi.org/10.1038/nrg2899> PMID: 21102527
22. Granados DP, Rodenbrock A, Laverdure J-P, Côté C, Caron-Lizotte O, Carli C, et al. Proteogenomic-based discovery of minor histocompatibility antigens with suitable features for immunotherapy of hematologic cancers. *Leukemia*. 2016; 30: 1344–1354. <https://doi.org/10.1038/leu.2016.22> PMID: 26857467
23. Daouda T, Perreault C, Lemieux S. pyGeno: A Python package for precision medicine and proteogenomics. *F1000Research*. 2016; 5. <https://doi.org/10.12688/f1000research.8251.2> PMID: 27785359
24. Courel M, Clément Y, Bossevain C, Foretek D, Vidal Cruchez O, Yi Z, et al. GC content shapes mRNA storage and decay in human cells. Weis K, Manley JL, editors. *eLife*. 2019; 8: e49708. <https://doi.org/10.7554/eLife.49708> PMID: 31855182
25. Kanaya S, Yamada Y, Kinouchi M, Kudo Y, Ikemura T. Codon Usage and tRNA Genes in Eukaryotes: Correlation of Codon Usage Diversity with Translation Efficiency and with CG-Dinucleotide Usage as Assessed by Multivariate Analysis. *J Mol Evol*. 2001; 53: 290–298. <https://doi.org/10.1007/s002390010219> PMID: 11675589
26. Hanson G, Coller J. Codon optimality, bias and usage in translation and mRNA decay. *Nat Rev Mol Cell Biol*. 2018; 19: 20–30. <https://doi.org/10.1038/nrm.2017.91> PMID: 29018283
27. Waldman YY, Tuller T, Shlomi T, Sharan R, Ruppin E. Translation efficiency in humans: tissue specificity, global optimization and differences between developmental stages. *Nucleic Acids Res*. 2010; 38: 2964–2974. <https://doi.org/10.1093/nar/gkq009> PMID: 20097653
28. Murphy JP, Konda P, Kowalewski DJ, Schuster H, Clements D, Kim Y, et al. MHC-I Ligand Discovery Using Targeted Database Searches of Mass Spectrometry Data: Implications for T-Cell Immunotherapies. *J Proteome Res*. 2017; 16: 1806–1816. <https://doi.org/10.1021/acs.jproteome.6b00971> PMID: 28244318
29. Laumont CM, Vincent K, Hesnard L, Audemard É, Bonneil É, Laverdure J-P, et al. Noncoding regions are the main source of targetable tumor-specific antigens. *Sci Transl Med*. 2018; 10. <https://doi.org/10.1126/scitranslmed.aau5516> PMID: 30518613
30. Paul S, Croft NP, Purcell AW, Tschärke DC, Sette A, Nielsen M, et al. Benchmarking predictions of MHC class I restricted T cell epitopes in a comprehensively studied model system. *PLOS Comput Biol*. 2020; 16: e1007757. <https://doi.org/10.1371/journal.pcbi.1007757> PMID: 32453790
31. Powers DMW. Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation. *ArXiv201016061 Cs Stat*. 2020 [cited 27 May 2021]. Available: <http://arxiv.org/abs/2010.16061>
32. Dersh D, Yewdell JW, Wei J. A SIINFEKL-Based System to Measure MHC Class I Antigen Presentation Efficiency and Kinetics. In: van Ender P, editor. *Antigen Processing: Methods and Protocols*. New York, NY: Springer; 2019. pp. 109–122. [https://doi.org/10.1007/978-1-4939-9450-2\\_9](https://doi.org/10.1007/978-1-4939-9450-2_9) PMID: 31147936
33. Kim JH, Lee S-R, Li L-H, Park H-J, Park J-H, Lee KY, et al. High Cleavage Efficiency of a 2A Peptide Derived from Porcine Teschovirus-1 in Human Cell Lines, Zebrafish and Mice. *PLOS ONE*. 2011; 6: e18556. <https://doi.org/10.1371/journal.pone.0018556> PMID: 21602908

34. Shastri N, Gonzalez F. Endogenous generation and presentation of the ovalbumin peptide/Kb complex to T cells. *J Immunol.* 1993; 150: 2724–2736. PMID: [8454852](#)
35. Laumont CM, Daouda T, Laverdure J-P, Bonneil É, Caron-Lizotte O, Hardy M-P, et al. Global proteogenomic analysis of human MHC class I-associated peptides derived from non-canonical reading frames. *Nat Commun.* 2016; 7: 10238. <https://doi.org/10.1038/ncomms10238> PMID: [26728094](#)
36. Zanker DJ, Oveissi S, Tschärke DC, Duan M, Wan S, Zhang X, et al. Influenza A Virus Infection Induces Viral and Cellular Defective Ribosomal Products Encoded by Alternative Reading Frames. *J Immunol.* 2019; 202: 3370–3380. <https://doi.org/10.4049/jimmunol.1900070> PMID: [31092636](#)
37. Ehx G, Perreault C. Discovery and characterization of actionable tumor antigens. *Genome Med.* 2019; 11: 29. <https://doi.org/10.1186/s13073-019-0642-x> PMID: [31039809](#)
38. Sette A, Fikes J. Epitope-based vaccines: an update on epitope identification, vaccine design and delivery. *Curr Opin Immunol.* 2003; 15: 461–470. [https://doi.org/10.1016/s0952-7915\(03\)00083-9](https://doi.org/10.1016/s0952-7915(03)00083-9) PMID: [12900280](#)
39. Bengio Y, Ducharme R, Vincent P, Janvin C. A neural probabilistic language model. *J Mach Learn Res.* 2003; 3: 1137–1155.
40. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature.* 2015; 521: 436–444. <https://doi.org/10.1038/nature14539> PMID: [26017442](#)
41. Daouda T. Mariana: The Cutest Deep learning Framework. 2015. Available: Available: <https://www.github.com/tariqdaouda/Mariana>
42. Gibson DG, Young L, Chuang R-Y, Venter JC, Hutchison CA, Smith HO. Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nat Methods.* 2009; 6: 343–345. <https://doi.org/10.1038/nmeth.1318> PMID: [19363495](#)
43. Bell C, English L, Boulais J, Chemali M, Caron-Lizotte O, Desjardins M, et al. Quantitative Proteomics Reveals the Induction of Mitophagy in Tumor Necrosis Factor- $\alpha$ -activated (TNF $\alpha$ ) Macrophages\*. *Mol Cell Proteomics.* 2013; 12: 2394–2407. <https://doi.org/10.1074/mcp.M112.025775> PMID: [23674617](#)
44. Karttunen J, Sanderson S, Shastri N. Detection of rare antigen-presenting cells by the lacZ T-cell activation assay suggests an expression cloning strategy for T-cell antigens. *Proc Natl Acad Sci.* 1992; 89: 6020–6024. <https://doi.org/10.1073/pnas.89.13.6020> PMID: [1378619](#)