# Supporting Information

**Supplementary Tables**

Supplementary Table S1:List of experiments in the non-HLA build and the protease used.

Supplementary Table S2: Listing of peptides that are mapped to Ribo-seq ncORFs from the Human non-HLA PeptideAtlas 2023-06 build

Supplementary Table S3: Listing of Ribo-seq ncORFs annotated in the Human non-HLA PeptideAtlas 2023-06 build

Supplementary Table S4: Listing of peptides that are mapped to Ribo-seq ncORFs from the Human HLA PeptideAtlas 2023-11 build

Supplementary Table S5: Listing of detected  Ribo-seq ncORFs from the Human HLA PeptideAtlas 2023-11 build

Supplementary Table S6: List of HLA build MS runs, including the HLA type of each MS run.

Supplementary Table S7: List of 677 HLA-I peptides, including their sequence, best allele, the 22 features that the model used for training, and the output probabilities from the model.

Supplementary Table S8: List of 7,264 ncORFs along with the features that were used to train machine learning models and output probabilities of the model.

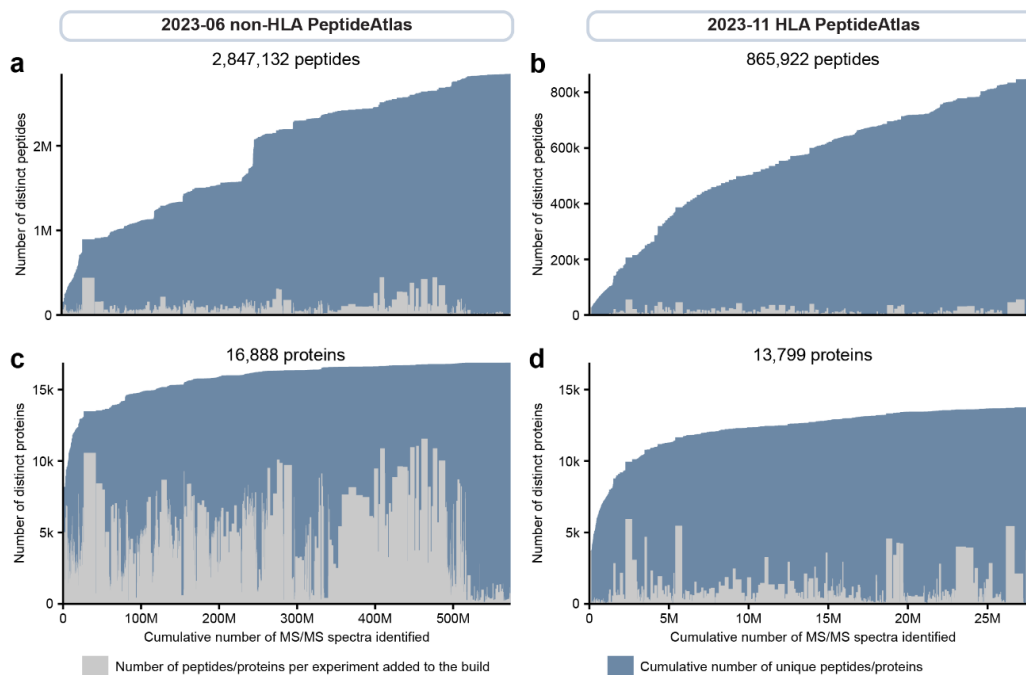Supplementary Table S9: FDR metrics for the non-HLA build analysis

Supplementary Table S10: FDR metrics for the HLA build analysis
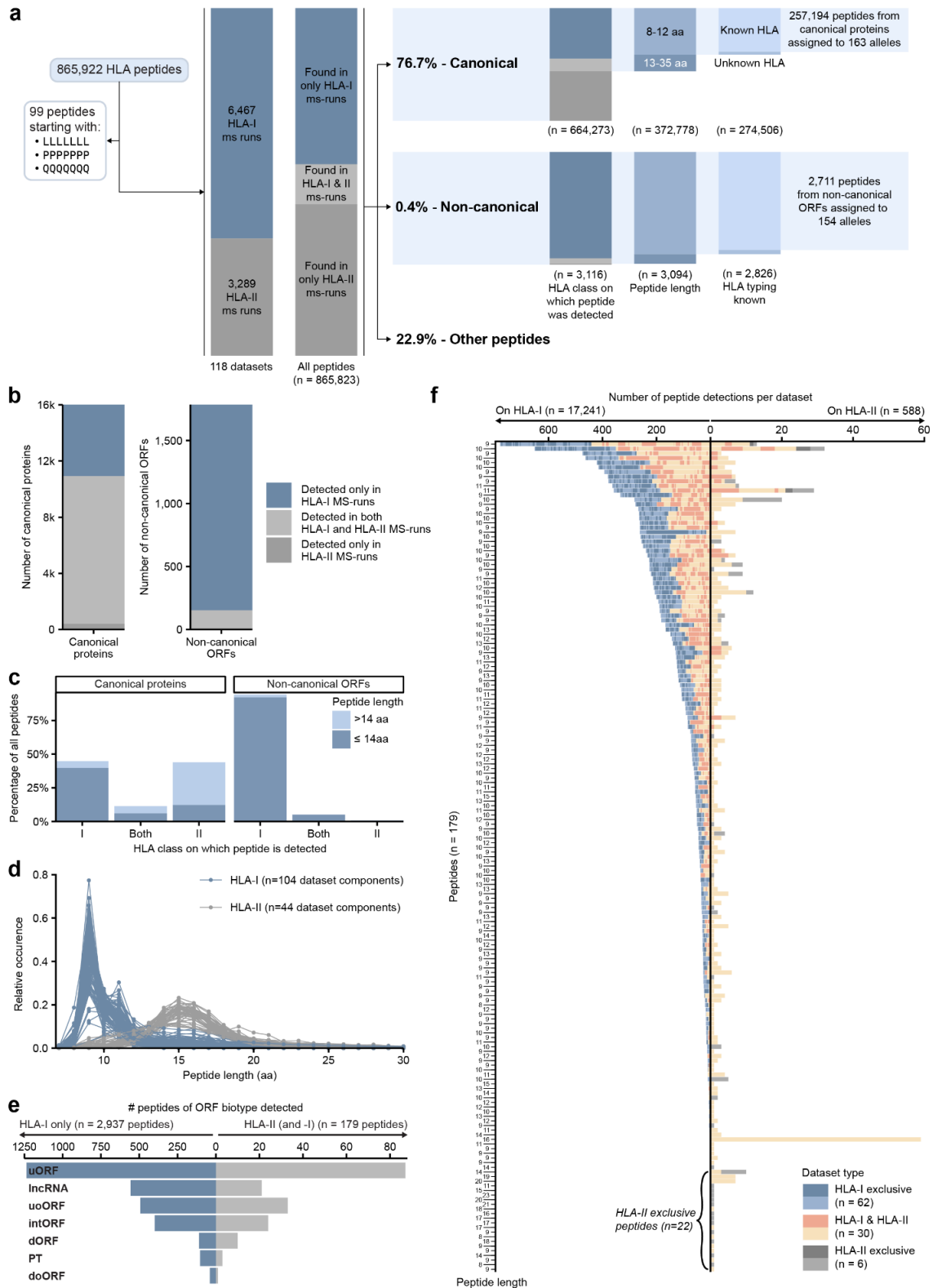

**Supplementary Documents**

Supplementary Document S1 - Discussion of ncORF detections in non-HLA data

Supplementary Document S2 - Discussion of machine learning results predicting detectability of ncORFs.
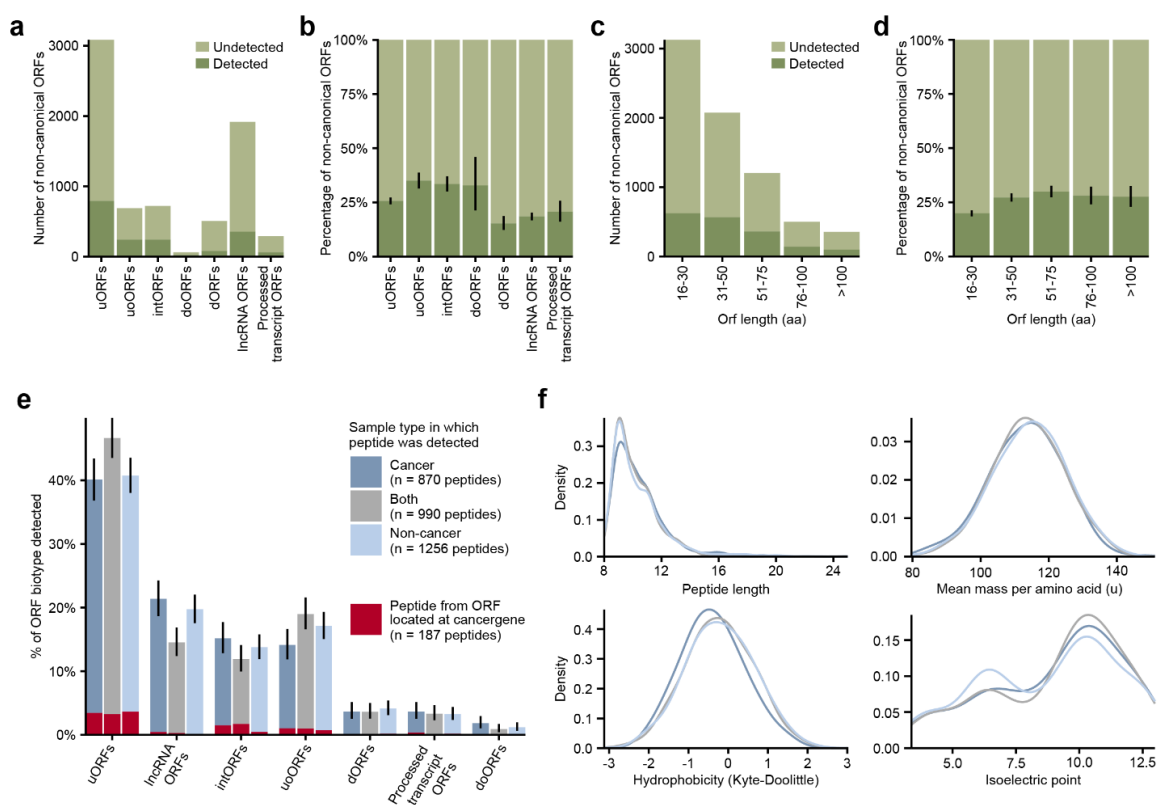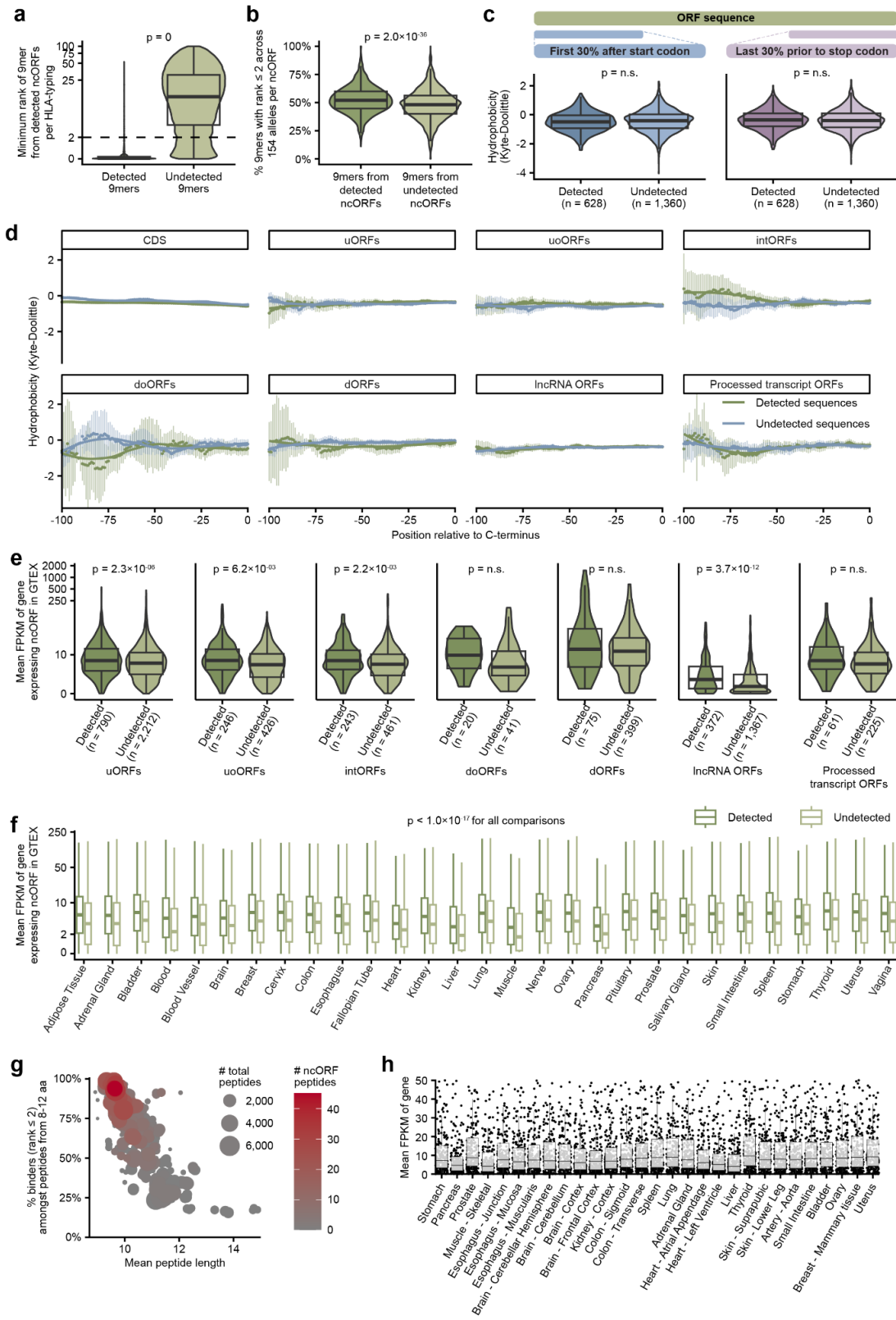
## Supplementary Figures



**Supplementary Figure S1.** The number of distinct peptides and proteins as datasets were added to the Human non-HLA (left) and HLA (right) PeptideAtlas. (**a**) Over 2.8 million distinct peptides have been observed in the 573 million PSMs in the non-HLA build. Each rectangle is one of the 1,172 experiments. Blue rectangles represent the cumulative number of distinct peptides in the build, while the gray rectangles depict the total number of distinct peptides within each experiment. (**b**) Over 0.86 million distinct peptides have been observed in the 28 million PSMs in the HLA build. Each rectangle is one of the 592 experiments. Blue rectangles represent the cumulative number of distinct peptides in the build, while the gray rectangles depict the total number of distinct peptides within each experiment. (**c**) The blue rectangles depict the cumulative 16,888 canonical proteins that have been cataloged in the 2023-06 Human non-HLA PeptideAtlas, whereas the gray rectangles show the total number of proteins present in each of the 1,172 experiments. (**d**) The blue rectangles depict the cumulative 13,799 canonical proteins that have been cataloged in the 2023-11 Human HLA PeptideAtlas, whereas the gray rectangles show the total number of proteins present in each of the 592 experiments. Although the total number of peptides continues to increase steadily, progress in the number of proteins is now very slow. Over the last 100 million PSMs, the cumulative counts are increasing by ~2,000 peptides per million PSMs and ~1 newly identified protein per million PSMs.
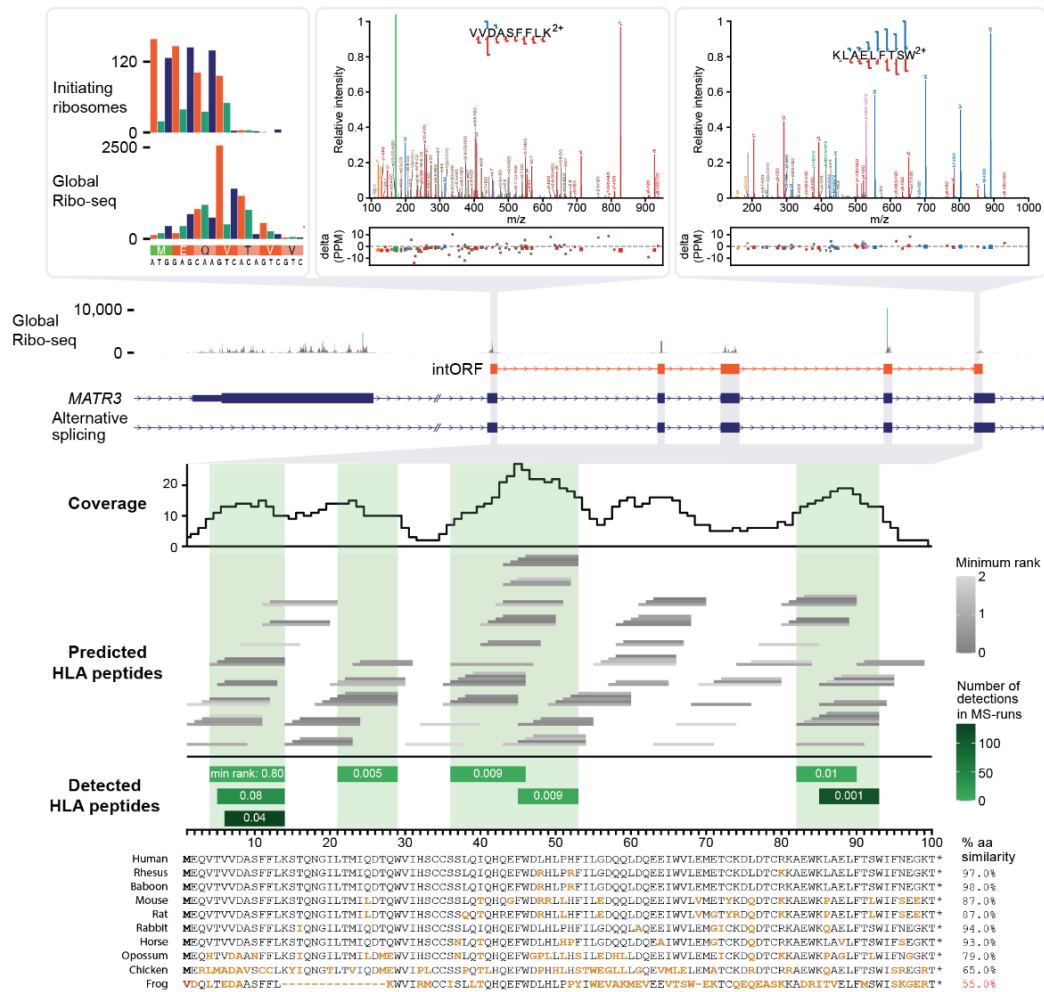
**Supplementary Figure S2.** Detection of ncORF peptides in HLA-I and HLA-II, and in cancer and non-cancer samples. (**a**) Schematic illustrating the total numbers of peptides (from both normal proteins and ncORFs) extracted from the total set of peptides. Depending on the analysis, peptides were further selected to those detected on HLA-I, had a length from 8-12 amino acids, and originated from an MS run with a known HLA-typing. The counts below each bar denote the number of distinct peptides. The distinction between "canonical", "non-canonical", and "other peptides" is defined in the methods. (**b**) Barplots showing for the detected canonical proteins (left) and detected ncORFs (right) whether their detected peptides were exclusively detected in HLA-I or HLA-II MS-runs, or in both. (**c**) Barplots showing for canonical proteins and ncORFs the percentage of peptides found exclusively in HLA-I and HLA-II MS-runs, or in both. Bars are colored by the peptide length being ≤14aa, or >14aa in length. (**d**) Line graph showing the peptide length distribution per dataset component split by HLA-class. (**e**) Bar plot showing the number of peptides detected per ORF exclusive to HLA-I samples (left) and those present in HLA-II samples, possibly in addition to HLA-I samples (right). Please note the x-axis scales differ by an order of magnitude between the left and right part of this panel. HLA-I and HLA-II peptide detection is not mutually exclusive as HLA-I peptides might be accidentally recovered from HLA-II pulldown experiments. (**f**) The frequency of peptide detection in HLA-I MS-runs (left) and HLA-II MS-runs (right) per peptide. Each alternating shade corresponds to a different dataset, with shades grouped by dataset type. The left axis denotes peptide lengths. Please note x-axis scales differ by an order of magnitude between the left and right part of this panel. Only peptides detected in at least one HLA-II sample are included. 22 of the 179 distinct peptides were exclusively detected in HLA-II samples. Fourteen of these peptides have a length 14 amino acids or greater, suggesting a potential presentation by HLA-II. This is still a minority in contrast to the total amount of 3,116 non-canonical ORF derived HLA peptides.

**Supplementary Figure S3.** (**a-d**) Comparisons of the detected and undetected non-canonical ORFs. (**a**) The total number of non-canonical ORF per ORF type and the number of ORFs for which a peptide was observed. (**b**) As in (**a**), but now shown in percentages. (**c**) The total number of non-canonical ORF grouped by length and the number of ORFs for which a peptide was observed. (**d**) As in (**c**), but now shown in percentages. Black lines on bar graphs indicate 95% confidence intervals. (**e**) Bar plot showing the proportion of Ribo-seq ORF-derived HLA peptides detected per biotype, categorized by whether the peptide was exclusively identified in immunopeptidomics analyses of cancer tissues or cell lines, non-cancer samples, or both. No significant changes in ORF biotype recovery are observed between these sample types. Peptides originating from ncORFs located on a known cancer gene are colored red. Black lines on bar graphs indicate 95% confidence intervals. (**f**) Density plots comparing the Ribo-seq ORF-derived HLA peptides differentiated on sample type (as depicted in (**e**)): cancer, non-cancer, or both. The plots compare peptides by their length, mass, hydrophobicity (Kyte-Doolittle scale), and isoelectric point. No significant changes between the distributions of these density plots can be observed.

**Supplementary Figure S4.** Potential determinants of ncORF detection. **(a)** Violin plot comparing for all MS-runs grouped by HLA-typing the minimum binding prediction rank for detected and undetected 9mers. **(b)** Violin plot comparing all 9mers from detected and undetected non-canonical ORFs. The y-axis shows per ORF the percentage of 9mers with a NetMHCpan rank ≤ 2 across all 154 alleles associated with ncORF peptides. **(c)** Violin plots similar to **(4a)** comparing the hydrophobicity by the Kyte-Doolittle scale between detected and undetected ncORFs for the first 30% of the ncORF sequence after the start codon, or the last 30% of the ncORF sequence. Statistical tests were performed with the two-sided Wilcoxon test, reported p-values were adjusted for multiple testing with Bonferroni correction. **(d)** Comparison of the hydrophobicity similar to **(4b)** between detected and undetected ncORFs/CDS per ncORF biotype. Each dot represents the average hydrophobicity of the amino acids at that position and the 14 amino acids before that position per ncORF biotype or CDS grouped by whether these were detected or not in the immunopeptidomics data. The lines were fitted using Local Polynomial Regression Fitting. Vertical bars represent 95% confidence intervals. Note that because ncORFs are mostly smaller than 100 aa, confidence intervals get larger with increasing C-terminus offset. **(e)** Comparison of the expression levels of detected and undetected ncORFs similar to **(4c)**, but split per biotype. On the y-axis, the mean FPKM in GTEX of genes expressing an ncORF is shown on a pseudo-log scale. 326 ncORFs for which the gene id was not present in GTEX are not shown. Significance was determined using two-sided Wilcoxon tests, reported p-values were adjusted for multiple testing with Bonferroni correction. **(f)** Comparison of the expression levels of detected and undetected ncORFs similar to **(e)**, but split per tissue. Outliers are not shown in the graph. Significance was determined using two-sided Wilcoxon tests, and p-values were adjusted for multiple testing with Bonferroni correction. All comparisons were found to be significant. **(g)** Dot plot similar to **(3e)**, for MS-runs originating from the HLA-ligand-atlas. The plot visualizes the correlation between mean peptide length and the percentage of predicted binders amongst peptides with a length between 8 and 12 amino acids (NetMHCpan rank ≤ 2) per MS run. Dot size corresponds to the total number of peptides per MS-run. Dot color corresponds with the percentage of non-canonical ORF-derived peptides per MS-run. Statistical tests were performed with the two-sided Wilcoxon test, reported p-values were adjusted for multiple testing with Bonferroni correction. **(h)** Comparison of the GTEx expression of 224/277 genes from which ncORFs in the HLA ligand atlas originate (53 genes with ncORFs in the HLA ligand atlas were not present in GTEx). GTEx tissues comparable to those from the HLA ligand atlas were selected, and sorted in the same way as in **(4e)**. Each represents the mean FPKM of a gene across these tissue samples in GTEx. Only genes with a mean FPKM lower than 50 are plotted for clarity, but all 224 genes were included for the boxplots.

**Supplementary Figure S5.** Overview of data available for c5norep142, an intORF in the *MATR3* gene. Ribo-seq data shows the initiation of translation at the methionine translation initiation codon (green), as determined by enrichment of ribosomes at the TIS. Two peptide spectral matches for HLA-I peptides VVDASFFLK and KLAELFTSW are shown having nearly complete sequence coverage (USIs are mzspec:PXD037270:Liv32_1176935F:scan:33690:VVDASFFLK/2 and mzspec:PXD011628:PBMC009_msms37:scan:16281:KLAELFTSW/2, respectively). The lowest panel shows the position of all 8 peptides that were observed in the immunopeptidomics data. The color shading indicates the number of MS runs in which each peptide was observed. The middle panel shows all peptides that are predicted with NetMHCpan to be observable in the MS runs (i.e. they are predicted to bind with NetMHCpan score <2 to at least one allele in one of the samples in which peptides were observed). The top part shows the number of predicted binding peptides in which each amino acid was located. Green shadings indicate which part of the ORF sequence was

45

observed. Except for the region near the offset of 62, detected peptides occurred in the regions with the highest numbers of predicted binders.