# Alternate RNA decoding results in stable and abundant proteins in mammals

Shira Tsour,<sup>1,2</sup> Rainer Machne,<sup>1</sup> Andrew Leduc,<sup>1</sup> Simon Widmer,<sup>1</sup> Jeremy Guez,<sup>3</sup> Konrad Karczewski,<sup>3</sup> & Nikolai Slavov<sup>1,4, $\boxtimes$ </sup>

<sup>1</sup>Departments of Bioengineering, Biology, Chemistry and Chemical Biology, Single Cell Proteomics Center, Northeastern University, Boston, MA 02115, USA; <sup>2</sup>Alnylam Pharmaceuticals, Cambridge, MA, USA; <sup>3</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA; Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA 02114, USA; <sup>4</sup>Parallel Squared Technology Institute, Watertown, MA, USA

⊠ Correspondence: nslavov@northeastern.edu

∈ Data & code: decode.slavovlab.net

Amino acid substitutions may substantially alter protein stability and function, but the contribution of substitutions arising from alternate translation (deviations from the genetic code) is unknown. To explore it, we analyzed deep proteomic and transcriptomic data from over 1,000 human samples, including 6 cancer types and 26 healthy human tissues. This global analysis identified 60,024 high confidence substitutions corresponding to 8,801 unique sites in proteins derived from 1,990 genes. Some substitutions are shared across samples, while others exhibit strong tissue-type and cancer specificity. Surprisingly, products of alternate translation are more abundant than their canonical counterparts for hundreds of proteins, suggesting sense codon recoding. Recoded proteins include transcription factors, proteases, signaling proteins, and proteins associated with neurodegeneration. Mechanisms contributing to substitution abundance include protein stability, codon frequency, codon-anticodon mismatches, and RNA modifications. We characterize sequence motifs around alternatively translated amino acids and how substitution ratios vary across protein domains, tissue types and cancers. The substitution ratios are positively associated with intrinsically disordered regions and genetic polymorphisms in gnomAD, though the polymorphisms cannot account for the substitutions. Both the sequence and the tissue-specificity of alternatively translated proteins are conserved between human and mouse. These results demonstrate the contribution of alternate translation to diversifying mammalian proteomes, and its association with protein stability, tissue-specific proteomes, and diseases.

#### Abbreviations

SAAP – substituted amino acid peptide (product of alternate translation); BP – base peptide translated according to the genetic code; RAAS – ratio of amino acid substituted peptide to the corresponding non-substituted peptide; CCRCC – clear cell renal cell carcinoma; BRCA – breast invasive carcinoma; UCEC – uterine corpus endometrial carcinoma; PDAC – pancreatic ductal adenocarcinoma; LUAD – lung adenocarcinoma; LSCC – lung squamous cell carcinoma

# Introduction

Genetic mutations, RNA editing, and "alternate translation" (mRNA translation deviating from the genetic code) may result in amino acid substitutions. Some substitutions arising from mutations profoundly change protein activity. For example, the V600E phosphomimetic substitution in BRAF destabilizes hydrophobic interactions and constitutively increases BRAF activity up to 500fold, inducing tumorigenesis<sup>1</sup>. Similarly, the H1047R substitution in the catalytic subunit  $p110\alpha$ of PI3K causes extensive cellular remodeling<sup>2</sup>. Substitutions introduced by RNA editing also alter protein functions, such as the A-to-I editing of the glutamate receptor mRNA<sup>3</sup>.

In the absence of genetic mutations and RNA editing detectable in RNA sequences, amino acid substitutions may arise from translation alternate to the genetic code. Such alternate decoding can occur when translating all open reading frames (short or long, annotated or not annotated as protein-coding) but remain less characterized. For decades, amino acid substitutions have been detected based on gel shifts of radioactively labeled proteins<sup>4</sup>. Mass spectrometry (MS) greatly increased the power of detecting substitutions<sup>5,6</sup>, and deep learning predictions of peptide elution times and MS fragmentation spectra are facilitating the validation of non-canonical amino acid sequences<sup>7–9</sup>. Translational substitutions are mostly studied in the context of amino acid starvation and antibiotic treatments as errors<sup>10–12</sup> but have also been detected in healthy organisms, mostly unicellular microorganisms, where the prevalence and the stability of alternatively translated proteins are presumed low<sup>12–14</sup>. They are often considered as translation errors resulting in low level molecular noise, though studies have suggested that non-cognate amino acid acylation may confer beneficial functions<sup>15,16</sup>. Another possible mechanism may involve mRNA modifications, such as pseudouridylation, which may recode stop codons to promote readthrough<sup>17,18</sup>, though its impact on the sequence of endogenous proteins is uncharacterized.

If the rate of alternate translation is low, the abundance and significance of its protein products may be low as well. However, the abundance of proteoforms harboring amino acid substitutions is not determined solely by the rate of their synthesis. It also depends on the rate of their degradation.

While most substitutions are likely to destabilize the substituted proteoform, some might stabilize them; if stabilized by substitution, a proteoform may accumulate to high abundance. Whether such protein stabilization contributes to elevated levels of proteins synthesized by alternate RNA decoding remains unknown. To explore this question, we systematically identified and quantified amino acid substitutions across healthy and cancer human and mouse tissues. We identified, validated and characterized thousands of abundant amino acid substitutions, including mechanisms contributing to their origin and their impact on protein stability.

# Results

# Systematic identification and validation of amino acid substitutions

To globally identify peptide modifications, we analyzed paired MS and RNA-seq data from 1,068 human samples described in Extended Data Fig. 1a,b. They included 6 cancer types (renal, uterine, pancreatic, breast, prostate, lung squamous cell carcinoma and lung adenocarcinoma) and matched normal adjacent tissue from the Clinical Proteomic Tumor Analysis Consortium (CPTAC)<sup>20–25</sup> and 26 tissue types from healthy individuals<sup>26</sup>, Fig. 1a. We included TMT-labeled and label-free MS data from different laboratories to explore peptide modifications across tissue types and disease states independent from biases specific to particular datasets. To minimize the influence of technical variation leading to missing data (peptides identified only in some samples), our analysis focuses on abundance ratios between quantified peptides, which are less affected by variations in data acquisition parameters.

To identify modified peptides, we first predicted a protein sequence database for each patient sample by using the genetic code to *in-silico* translate the corresponding transcriptome. Transcriptomes were assembled from paired-end Illumina reads (minimum of 120 million per cancer sample, 18 million per healthy tissue) using a workflow adapted from Galaxy<sup>27</sup>, Fig. 1a. The mean sequence coverage for all transcripts exceeds 98% (median 100%), corresponding to about 71,000 transcripts

per sample with 100% sequence coverage, Extended Data Fig. 1c,d. Using the patient-specific databases, we search the MS data with the dependent peptide algorithm of MaxQuant<sup>6</sup>, which implements ModifiComb developed by Savitski *et al.*<sup>5</sup>. This algorithm tests whether MS2 spectra not matched to a database sequence correspond to modifications of peptides identified in the sample. Such correspondence is reflected in systematic mass shifts in the precursor peptide ions and the associated peptide fragments<sup>5</sup>, Fig. 1b, and has been used previously to detect substitutions in bacteria and yeast<sup>13</sup>.

Our dependent peptide search identified almost 9 million modified peptides, 40% of which have mass shifts corresponding to 419 distinct known post-translational modifications (PTMs), Fig. 1c,d, Extended Data Fig. 1e; 124,000 of these peptides exhibited mass shifts consistent with amino acid substitutions and no other known modifications, Fig. 1c. The occurrence of known PTMs with biological and technical origins, such as M oxidation, deamidation, and carbamylation, is similar across cancer types, tissue types and datasets, Fig. 1d. These PTMs are reported in Supplemental Data Table 1 and can support much further analysis, which is beyond the scope of this article. Deeper datasets, such as the healthy tissues, have more identified peptides of all types, Fig. 1c.

# Sequence confirmation of amino acid substitutions

We focused on substituted amino acid peptides (which we call SAAP) supported by good RNA sequence coverage and peptide fragmentation spectra, Supplemental Fig. 1 and 2. We systematically applied rigorous filters to minimize false positives, including: (1) removed all peptides with a mass shift consistent with another known PTM; (2) validated SAAP discovered by dependent peptide search with a standard database-search, Fig. 1b,e; (3) required that at least 2 peptide fragments reflect a mass shift corresponding to a localized substitution, Fig. 1f, Extended Data Fig. 1f, Supplemental Fig. 3; (4) validated SAAP by comparing their retention time and fragmentation patterns to deep learning predictions, Fig. 1g, from the re-scoring pipeline Oktoberfest<sup>8</sup>, (5) excluded sequences that could be generated by a 6 reading frame translation of any transcript, Fig. 1e, (6) filtered SAAP at 1% FDR computed only across SAAP, Extended Data Fig. 1g, and (7) removed

SAAP corresponding to substitutions at trypsin cleavage sites (K, R); across all datasets, only 24 SAAP have K and R substitutions, and about 90% of them represent peptides with missed cleavages or with substitutions from  $K \to R$  or vice versa (Extended Data Fig. 1h; Supplemental Data Table 2), supporting the validity of the results. Our confidence in identified SAAP is bolstered by the validation with standard database-search, Fig. 1e, low mass errors for SAAP, Extended Data Fig. 1i, increased confidence after re-scoring, Fig. 1f, Extended Data Fig. 1g, identification of the same substitution site in a human tissue digested by multiple proteases, Fig. 1h, and accurate retention time prediction, as determined by DeepRTplus<sup>28</sup>, Extended Data Fig. 1j. The database search validated about 40% of the SAAP discovered by DP search, Fig. 1e; this result is expected since the database search does not use the spectra of the main peptides and thus has lower sensitivity. These SAAP were further filtered to remove peptides mapping to immunoglobulins and trypsin, since we could not be confident that these are true instances of alternate translation. For about half of all SAAP, the measured mass shift indicating a substitution is localized with high confidence to a single position (Supplemental Fig. 1, 2 and 3a), and thus cannot arise from multiple combinatorial modifications on different residues adding up to the mass shift; for the remaining SAAP, unknown combinatorial modifications cannot be rigorously disproved (see Methods for details). However, specific hypotheses for explaining observed mass shifts (such as  $Q \rightarrow G$  versus cleavage of terminal A) may be discriminated based on their fragment ion intensities. The empirical intensities are closer to those predicted for  $Q \rightarrow G$ , thus suggesting it is the more likely hypothesis, Supplemental Fig. 3b-f. Our extensive filtering likely removes many true SAAP, but it helps establish a robust set of about 9,000 unique SAAP that can support exploration of their biological implications, Supplemental Data Table 2. Blasting these peptides against the 120,000 Ensembl proteins or Uniprot confirmed that these sequences have not be previously predicted from genetic code translations.

## Quantification of amino acid substitutions

To provide context for the prevalence and possible functions of SAAP, we investigated the abundance of SAAP relative to their corresponding encoded base peptides (BP). The abundance of a

peptide depends both on the synthesis rates and the degradation rates of its parent proteins. At steady state, the ratio between the peptides with and without substitutions, i.e. SAAP/BP, which we term RAAS, equals the ratio of corresponding synthesis and degradation rates of the alternatively translated and encoded proteoforms, Fig. 2a. Most SAAP are expected to be synthesized at a low rate and degraded at a high rate, resulting in low SAAP abundance and low RAAS. However, even substitutions incorporated at low rates may accumulate to high level if they stabilize the resulting proteoforms.

While median RAAS estimates are consistent with previous reports (reviewed by ref.<sup>12</sup>), we observed that about 10% of SAAP have substitution ratios exceeding 1, Fig. 2b. Surprisingly, our estimates suggest that the alternatively translated proteoforms are the most abundant protein products for 378 proteins spanning diverse biological functions. This finding is highly unexpected and motivated us to further evaluate it based on multiple lines of evidence. First, a detailed inspection of the underlying RNA sequence data and mass spectra provide strong support for the absence of alternate alleles and the confidence in the assigned amino acid sequence, Supplemental Fig. 1 and 2. Accordingly, we find that the identification confidence to be the same for SAAPs having high and low RAAS, Extended Data Fig. 2a-c. Second, we tested if differences in peptide ionization can explain high RAAS and found that SAAP have about the same ionization efficiency as their corresponding BP according to the model proposed by Liigand, et al.<sup>29</sup>. Differences in a single amino acid residue have modest effects on ionization<sup>29</sup>, and even the largest deviations are too small to significantly affect RAAS estimates, Extended Data Fig. 2d-f. Third, we confirmed strong agreement in RAAS estimates for the same substitution quantified from different peptides, Fig. 2c, Extended Data Fig. 2g; the distinct peptides are derived from MS analysis of independent tonsil sample digestions with different proteases. Fourth, we ensured that high RAAS are not an artifact of having highly abundant BPs with missed cleavages, Extended Data Fig. 2h. Thus, all four lines of evidence support the accuracy of our RAAS estimates.

As a fifth validation of RAAS estimates, we tested whether they are consistent with the abundance of peptides shared between canonical and substituted proteoforms, as estimated from their precur-



Figure 1 | Identification and validation of amino acid substitutions (a) A schematic of the data processing workflow. RNA-seq data is *in-silico* translated to patient-specific protein databases, which are used to search patient-matched MS proteomics data for peptides with modifications. The resulting peptides are filtered and further validated. (b) Search strategy for discovering modified peptides in MS spectra with ModifiComb. Spectra with systematic mass shifts of the precursor and fragment ions from identified (base) peptides correspond to modified peptides<sup>5</sup>. (c) Number of database peptides (blue), modified peptides (any mass shift, orange), peptides with PTMs (mass shift of a known modification, green) or candidate SAAP (mass shift of a substitution, red). Distributions are across CPTAC datasets or across tissue samples (labelfree data). (d) The number of PTM peptides identified per 1,000 unmodified peptides for the most common PTMs. (e) Percentages of candidate SAAP that are discarded as potential products of 6-frame translation (orange) or validated by database-based search and used in further analysis (blue). (f) Distribution of the number of detected fragment ions providing evidence for each substitution. Only substitutions supported by 2 or more fragment ions were analyzed further. (g) Re-scoring of validated SAAP. Confidence of peptide sequence identification is higher after re-scoring with Oktoberfest<sup>8</sup> (y-axis) than initially determined in the validation search with Andromeda<sup>19</sup> (x-axis). (h) Identification of AAS in tonsil proteomes from multiple enzymatic digests. Percentage of SAAP validated in 1 (blue), 2 (orange) or 3 or more (purple) digests, binned by SAAP abundance. The percentage of BP observed in multiple digests is also shown as a control.



Figure 2 | Abundance of proteins with amino acid substitutions (a) Model for the abundance of the canonical (P1) and the alternatively translated (P2) proteoforms. P1 and P2 have base (BP) and substituted (SAAP) peptides, respectively, and shared peptides are common to both. P1/P2 is estimated with the SAAP/BP ratio.(b) Distribution of substitution ratios computed for each patient sample in 7 datasets. (c) Reproducibility of RAAS for substitutions quantified in multiple digests. Data points are individual SAAP. (d) The abundance of shared peptides follows the abundance trends for BP and SAAP predicted from our model (panel a) and the corresponding ratios denoted on the x-axis. (e) RAAS fold changes (relative RAAS) were computed for each SAAP identified in a pair of patients, and their medians displayed. The matrix is sorted by column and row means, and the barplot shows median relative RAAS for each row, corresponding to a patient. (f) Mean of relative RAAS computed between a patient relative to all other patients, as in (c), weighted by number of shared SAAP. N indicates number of patients in each dataset. (g) Rank sorted proteins with RAAS > 0.1. Proteins from functional groups that are significantly enriched for high RAAS are highlighted. (h)– (k) (Continued on the next page)

(Continued) (h) Copy number estimates of SAAP by the histone ruler method<sup>30</sup>. (i) RAAS as a function of the minimum number of codon-anticodon mismatches needed for incorporating the detected amino acid. (j) The median RAAS of all substitutions mapping to a codon is negatively correlated to the relative frequency of the codon. Ordinary least squares fit shown in red, Pearson correlation and associated p-value annotated in plot. See Extended Data Fig. 4 for details. (k) Protein degradation rates ( $\alpha$ ) are computed as the slope of  $ln(1 + h/l) = \alpha t$ , where h/l is the heavy-to-light peptide ratio. For some peptides with abundant substitutions, SAAP degradation rates are up to ten times slower than the degradation rate of the corresponding BP. (I) The ratio of degradation rates for SAAP relative to BP is inversely proportional to their RAAS in primary hepatocytes. Substitutions identified only in the hepatocytes are in gray, those also identified in the CPTAC and healthy human tissues in black, and the Pearson correlation for the union is in red.

sor intensities. The abundance of the shared peptides should reflect the cumulative abundance of both encoded (P1) and alternatively translated (P2) proteoforms, Fig. 2a, which is approximated by the most abundant proteoform, especially for very high and low RAAS. This expectation is strongly supported by the data, Fig. 2d. Specifically, we observe that for peptides with low RAAS, the BP peptide abundance is similar to the shared peptide abundance, while the SAAP abundance is orders of magnitude lower, Fig. 2d (bottom panel). As RAAS increases, the abundance of the SAAP becomes less distinct from the abundance of the shared peptides, while the BP abundance increasingly deviates from the abundance of shared peptides, Fig. 2d (top panel). For RAAS>1, substituted and shared peptides have similar abundance. Together, these trends bolster the conclusion that P1 proteoforms are dominant at low RAAS while P2 proteoforms are dominant at high RAAS.

Our sixth test for the reliability of high RAAS evaluated their consistency with the correlation patterns of BP or SAAP abundance to shared peptide abundance across patients, Extended Data Fig. 2i. These six complementary analyses indicate that the quantification derived from many peptides of different types (BP, shared peptides, and SAAP) at both MS1 and MS2 level all dovetail together and support that SAAP with RAAS>0 represent peptides from the most abundant proteoforms. Therefore, thousands of proteins have abundant proteoforms with amino acid substitutions. For some proteins, these are the most abundant proteoforms.

Having established confidence in the estimates of substitution ratios, we explore them across the different proteins, patients and cancer types. These ratios span  $10^8$  range at the level of individual

substitutions (Fig. 2b and Extended Data Fig. 3a,b) which shrinks to 100-fold range across patients (the median RAAS per patient), Extended Data Fig. 3c, and to 4-fold across cancer types, Extended Data Fig. 3d, Supplemental Data Tables 3,4. While most observed SAAP are specific to a dataset, there is notably higher overlap between the two lung cancer datasets and a cluster of shared SAAP that are also commonly found across the majority of patients in a dataset, Extended Data Fig. 3e,f. RAAS estimates based on peptides found in all datasets, which are not affected by differences in peptide detection, show significant ( $p < 10^{-12}$ ) variation in RAAS across cancer types, Extended Data Fig. 3g,h. Further, ratios of RAAS for SAAP intersected between pairs of patients show significant ( $p < 10^{-20}$ ) differences across datasets, with medians of relative RAAS varying over 2-fold from the lowest to the highest, Fig. 2e,f. These estimates suggest that there is biologically-driven RAAS variability across proteins, cancer types, and patients.

The proteins with high substitution ratios span multiple functional groups (Fig. 2g), including signal transduction, protein degradation and transcriptional regulation. Estimation of protein copy numbers by the histone ruler method<sup>30</sup> suggests that highly abundant alternatively translated proteoforms are present at hundreds to tens of thousands of copies per cell, Fig. 2h.

#### **RAAS** depends on codon frequency, tRNA pairing and RNA modifications

Confident in our quantification of SAAP, we explored the dependence of RAAS on the number of nucleotide mismatches corresponding to a mRNA-tRNA pairing that can translate the detected SAAP. This analysis indicated a strong dependence between RAAS and the minimum number of codon mismatches required for the corresponding substitution, Fig. 2i, an effect consistently observed across the datasets, Extended Data Fig. 3i. Just as 3 base pair mismatches are less likely to occur, so too substitutions that would require complete codon mismatch have lower RAAS. We also observed that for some substitutions, such as  $T \rightarrow V$  in cancer, there is significantly higher abundance of the tRNA ligase loading the incorporated amino acid relative to the tRNA ligase loading the encoded amino acid, Extended Data Fig. 3j.

Exploring the dependence of substitution ratios on codons, we found that they are inversely proportional to relative codon frequency (Fig. 2j and Extended Data Fig. 4) and to the codon

stability coefficient (Extended Data Fig. 3k), which is an empirical measure of codon usage<sup>31</sup>. These observations are consistent with a previously proposed model<sup>32,33</sup> that less frequent codons are associated with a smaller tRNA pool, which may result in increased rate of substitution by the ribosome for a tRNA that is more readily available, increasing alternate decoding at these sites.

Furthermore, we explored RAAS dependence on RNA modifications detected by direct RNA sequencing using nanopores. Using data by McCormick *et al*<sup>34</sup>, we found that uracil modifications overlap significantly ( $p < 10^{-10}$ ) with substitution sites (Supplemental Fig. 4a) and the modification fraction correlates significantly to the substitution ratios, Supplemental Fig. 4c,d. Together, these data suggest that multiple RNA related mechanisms influence the rate of alternate decoding, and thus contribute to high substitution ratios.

#### Increased protein stability contributes to substitution abundance

While variation in the rate of alternate translation clearly contributes to RAAS, protein stability may contribute as well, Fig. 2a). Fortunately, this contribution can be directly quantified using metabolic pulse with stable isotope labeled amino acids. Thus, we analyzed data from such experiments with primary human liver-derived cells<sup>35</sup>. Applying the analysis pipeline from Fig. 1a, we identified 8,278 SAAPs with a RAAS distribution similar to the one observed from the TMT and label-free data, Extended Data Fig. 31. These similarities extend to detecting SAAP with identical sequences and substitutions as those detected in the other datasets, Supplemental Data Table 5. Remarkably, SAAP detected in hepatocytes have 5-fold higher overlap with SAAPs from the label-free healthy liver samples compared to other tissues (Supplemental Fig. 5a), suggesting tissue specificity. We validated these findings with an independent database search using the MSFragger pipeline<sup>36</sup>, with deep-learning rescoring enabled<sup>37</sup>, Supplemental Fig. 5b. These results generalize our observations across another type of MS data acquisition that offers further constraints (from stable isotope incorporation) on sequence identification.

Our analysis of the metabolic pulse data allows direct evaluation of the dependence between protein degradation rates and RAAS. Many substituted peptides have lower degradation rates than their un-substituted counterparts, Fig. 2k. This global trend is reflected in a strong inverse correlation ( $p < 10^{-10}$ ) between RAAS and degradation rates in hepatocytes, Fig. 2l. Similar inverse correlations are observed with primary B cells and natural killer cells, Extended Data Fig. 3m,n. These results indicate that substituted proteoforms with high RAAS are stable, i.e., have low protein degradation rates. This is consistent with the expectation that most substitutions likely destabilize proteins and result in low abundance SAAP, below our detection limit. Yet, the SAAP that are detected, particularly those with high RAAS, correspond to substitutions that increase the lifetimes of their proteoforms in living cells, Fig. 2l.

#### Substitution ratios depend on amino acid and tissue types

Next, we examined how RAAS depends on the substitution type defined by the specific combination of encoded and incorporated amino acids. We find that the median substitution ratios for different substitution types vary from  $10^{-4}$  to over 1, Fig. 3a,b and Extended Data Fig. 5a. Thus substitution type can explain much of the observed variation in substitution ratios. Still, ratios vary within a substitution type as well; this variation is relatively small for some types, such as  $H \to Q$ , and larger for others, such as  $E \to D$ , Fig. 3a. The majority of substitution types are represented by SAAP that are over 100-fold less abundant than the corresponding encoded BP, Fig. 3a and Extended Data Fig. 5b. Yet, some substitution types, such as  $S \rightarrow G$ , have median RAAS>1, Fig. 3a,b. Most substitution types with RAAS>1 are observed relatively infrequently in the data, and their average properties therefore are more influenced by individual SAAPs, Fig. 3a and Extended Data Fig. 5a. Interestingly, the substitution ratios of substitution types cluster by the chemical properties of the encoded amino acid, Fig. 3b,c and Extended Data Fig. 5c. Specifically, substitutions of polar amino acids have consistently higher RAAS than substitutions of other amino acid types. This effect is strongest for polar  $\rightarrow$  special substitutions. In contrast, substitutions of charged or hydrophobic amino acids result in SAAPs with lower ratios, Fig. 3b,c and Extended Data Fig. 5c.

The pairwise correlations of substitution-type median substitution ratios across datasets indicate



**Figure 3** | **Substitution ratios depend on amino acid and tissue types** (a) Number of SAAP (upper panel) and median and 90th percentile RAAS (middle panel) displayed for each substitution type across all datasets. RAAS distributions across all samples are displayed for 6 substitution types (lower panel) spanning the range of RAAS medians, highlighting the variability in RAAS across individual samples, even for the same substitution type. Substitution types displayed were chosen due to large number of data points. (b) RAAS dotplot for all amino acid substitution types that had a significantly low or high RAAS distribution ( $p \le 10^{-10}$ ) in at least one dataset (columns). The y-axis is grouped by chemical properties of the encoded amino acid. (c) RAAS dotplot as in (b) but by chemical properties of the encoded and incorporated amino acid. (d) Median RAAS profiles across all substitution types are strongly correlated across all dataset pairs (Pearson correlation, weighted by number of SAAP). (e) RAAS dotplot as in (b) for encoded and incorporated amino acids. (f) Results from ANOVA confirming variance in RAAS profiles is driven primarily by substitution type and secondarily by tissue type.

significant similarities, with no obvious global distinction between the label-free (Healthy) data and the TMT-labeled CPTAC data, Fig. 3d. Variance in substitution ratios for individual substitution types across datasets is highly dependent on substitution type. We find that some substitution types, such as  $E \rightarrow N$  have significant agreement in RAAS across datasets, Extended Data Fig. 5d, while others, such as  $P \rightarrow S$  show clear dataset dependence, Extended Data Fig. 5e. Similarities in

substitution ratios of substitution types across datasets is majorly driven by similarities in RAAS for SAAP with the same encoded amino acid, rather than RAAS consistency for incorporated amino acids, Fig. 3e and Extended Data Fig. 5f,g. These findings support significant dependence of RAAS on the substitution type, and we sought to quantify it using ANOVA. The results (Fig. 3f) indicate that substitution type, especially the encoded amino acid, explains the most variance in RAAS. Tissue type is also significantly associated with substitution ratios. Fold change analysis of RAAS for substitution types in a specific tissue (using both cancer and healthy data) relative to all other tissues indicates that some substitutions, such as  $G \rightarrow S$  in pancreas, have significantly  $(p < 10^{-15})$  tissue-specific RAAS distributions, Extended Data Fig. 5h,i.

#### Linking substitutions to protein function and structure

Having established the associations of substitution ratios with substitution type and tissue specificity, we next explored potential structural and functional consequences of alternate translation events, starting with differences between cancer and adjacent non-cancer tissue. While we did not have the power to detect global trends in RAAS with tumor status (Extended Data Fig. 6a,b), there are clear outlier SAAP for which substitution ratios in some tumor types are significantly higher than RAAS in adjacent non-cancer from the same patient. One striking example is a substitution in lamin isoform A protein, for which RAAS in tumor samples was significantly higher than in surrounding control tissues in 3 cancer types, Fig. 4a. Another example is the  $N \rightarrow G$  substitution in the serine/threonine-protein phosphatase PP1-beta catalytic subunit. For lung cancer patients, this substitution has consistently higher RAAS in tumor samples than in the matched controls, Extended Data Fig. 6c. Furthermore, specific substitution types, such as  $H \rightarrow N$  and  $T \rightarrow Q$ exhibit significantly different RAAS in tumor samples compared to matched surrounding tissues, Fig. 4b. Together, these results demonstrate cancer specific differences in the abundance of protein products of alternate RNA decoding.

To identify biological functions that are most affected by alternate translation, we next identified Gene Ontology (GO) protein sets enriched with highly abundant alternatively translated proteoforms, Fig. 4c (top panel) and Extended Data Fig. 7. These significant GO groups suggest a



Figure 4 | Sequential, structural and functional context of amino acid substitutions (a) An example SAAP that is significantly more abundant in tumor samples than in normal adjacent tissue from the same patient in 3 cancer types. (b) Volcano plot with substitution types that have significantly different RAAS in tumor than in matched normal samples. (c) RAAS dotplots for all SAAP in proteins annotated by the indicated GO terms (top panel) or all SAAP carrying the indicated sequence motifs (bottom panel). The top panel shows all GO groups that had at least one significantly high RAAS distribution ( $p \le 10^{-15}$ ). Individual proteins with high RAAS from each GO group are shown in Extended Data Fig. 8. (d) Sequence difference logos of selected subsets of substitution sites. AAS denotes the site of the substitution and numbers refer to adjacent positions in the protein sequence. The y-axis shows the Jensen-Shannon divergence of the selected set of sequences (number n of sequences is indicated) compared to all other sequences in our data; \*\*\* indicates enrichment significance  $p < 10^{-10}$ . (e) Clusters of substitutions in the 1D sequence of a thioredoxin domain (left) and an actin ATPase domain (right). Pfam annotations are shown as black lines, and predicted secondary structures (S4pred),  $\alpha$ -helices and  $\beta$ -sheets, are shown with gray bars. (f) A set of three substitutions distant in the 1D sequences but clustered together in the 3D structure of the Ras-related protein, RAP1A. The substitutions have high RAAS, reflected in their color-coding (as in the legend in (c)).

broad impact of substitutions on the proteome, including proteins functioning in gene expression, cellular organization, and signaling. As suggested by observations that alternatively translated proteasome subunits are highly abundant, Fig. 2e, Fig. 4b, protein catabolic processes are significantly enriched with high RAAS proteoforms, Fig. 4c, Extended Data Fig. 7. Interestingly, signaling proteins are also enriched for high RAAS, indicating potential downstream consequences, Fig. 4c and Extended Data Fig. 7. For example, an alternatively translated isoform of EZR is found in all datasets, which may affect cellular adhesion and migration processes.

To explore the context of substitutions, we investigated sequence patterns associated with alternate translation. Sequence logos and amino acid enrichment around substitution sites and along tryptic peptides revealed four enrichment patterns common to different substitution types, Extended Data Fig. 8a-c. We observed that many highly abundant SAAP have a  $Q \rightarrow G$  substitution. While the mass shift for this substitution is also consistent with A cleavage from the peptide Nterminus, the fragment ion intensities are better explained by  $Q \rightarrow G$ , Supplemental Fig. 3b-f. Through sequence motif analysis, we identified 768 SAAPs significantly enriched for a "KRAQ" motif, in which a K, R, A, and G are enriched immediately before a  $Q \rightarrow G$  or  $Q \rightarrow A$  substitution, Fig. 4d. This sequence motif has substitution ratios of about 0.1 across all datasets (Fig. 4c bottom panel), significantly ( $p < 10^{-10}$ ) above the median RAAS.

Furthermore, we found that substitutions surrounded by CC, MM, or WW tend to have high RAAS and all show enrichment for specific substitution types, Fig. 4d. For example, CCxCC corresponds to high probability of Q or N to be substituted by A or G. Additional motifs with high RAAS substitutions, include an enrichment of M in +/-1 and +/-2 positions relative to the substitution site, Fig. 4d, Fig. 4c (bottom panel). This motif is generally found in the context of substitutions of T and S. A spatial proximity between reversibly sulfoxidized M and S/T kinase targets has been previously noted<sup>38</sup>, but a biological function or effects on protein stability of these loci have not been described. Another similarly interesting motif identified is W in +/- 1/2 position relative to the substitution, Fig. 4d, Fig. 4c (bottom panel). W is both the rarest and metabolically most costly amino acid<sup>39</sup>, and tryptophan codons have been observed to cause ribosome stalling and frame-shifting translation errors ("W bumps") in melanoma<sup>40</sup>. Here, we find W in the context of substitutions of branched chain amino acids, which are themselves known to be dysregulated in starvation and diabetes<sup>41</sup>, and particularly in PDAC<sup>42,43</sup>. These results suggest the existence of amino acid sequence patterns that are significant targets of alternate translation.

These observations led us to explore more broadly the spatial clustering of substitutions, both

within structural domains and within regions having high density of substitutions in the primary amino acid sequence (1D) or in the 3D structure of proteins. Analysis of Pfam domains confirmed enrichment of abundant substitutions in several structural domains, including those related to significant functional protein sets, such as the proteasome subunit domain, Ras family domain (signaling), and RNA recognition motif, Extended Data Fig. 9a. Many substitutions occur close to each other, with over 1,200 base peptides having multiple substituted peptides, Extended Data Fig. 9b. An example of such clustered substitutions is shown in Fig. 4e with the significantly enriched thioredoxin domain, Extended Data Fig. 9a. The ratios of these substitutions with both high and low RAAS occur in the actin domain of ACTG2 (Fig. 4e), as well as in MZB1 and the glycolytic domain of ALDOB, Extended Data Fig. 9c,d. We were also able to identify structural domains with clusters of alternate translation in 3D. Notably, the significant Ras family domain (Extended Data Fig. 9a) has three substitution sites with high RAAS that cluster together in the 3D structure of the RAP1A protein, Fig. 4f. Other interesting domains with clusters of alternate translation include the ribosomal and proteasomal protein complexes, Extended Data Fig. 9e,f.

#### **Predictors of substitution ratios**

To place alternate translation in broader biological context, we examined protein attributes correlated with median RAAS. We find that RAAS is positively correlated to protein length, Fig. 5a, and negatively correlated to protein half-life, Fig. 5b, features which are known to negatively correlate to one another by thermodynamic principles<sup>44</sup>. This suggests that shorter, more stable proteins, are less likely to have alternatively translated proteoforms that are more stable than their encoded form, with the caveat that these features and RAAS are also correlated to protein abundance, Extended Data Fig. 10a. Additionally, we observe that substitution sites with high RAAS are associated with protein regions that are more highly disordered and more lowly conserved, Fig. 5c and Extended Data Fig. 10b,c. This intuitively suggests that more highly conserved and structured protein regions are less likely to have alternatively translated proteoforms that are more stable than their canonical counterparts. These features, along with substitution site position and binding site information were used to successfully predict RAAS distributions using a random forest model, Fig. 5d.



Figure 5 | Predictors of substitution ratios (a) Protein RAAS is significantly positively correlated to protein length. Data points are proteins with substitutions. Protein RAAS is computed as the median RAAS across all SAAP identified in the protein. (b) Protein RAAS is significantly negatively correlated to protein half-lives from ref.<sup>35</sup>. Data points are proteins with substitutions. (c) RAAS dotplots (as in Fig. 4c) for substitution sites grouped by bins of a predicted disorder score (rows, based in an IUpred3 prediction) and sequence conservation (columns, based on MMseqs2 alignments). (d) An XGboost model using protein features described in (a-c) is able to accurately predict median RAAS per protein. Each feature used increases the model's predictive power. An average split of the data into test and train had a predictive Pearson correlation of 0.44. Additional details can be found in Methods. (e) Allele frequency in the gnomAD database for all possible missense variants in alternatively translated codons. (f) Observed / expected ratios for all missense variants in alternatively translated codons, determined from analysis of gnomAD database (see Methods). Missense variants are less constrained with increasing RAAS ( $p < 10^{-4}$ ). Data point colors correspond to RAAS quartiles as in (e). (g) 89 SAAP sequences from human tissues were found to be conserved in mouse lung, kidney and pancreas tissues with significantly correlated RAAS values. (h) Substitution type is the most significant driver of variability in RAAS values across the SAAP plotted in (e), followed by species and tissue type (ANOVA).

While these trends are broadly informative their interpretation has a caveat: Low ratio substitutions on lowly abundant proteins are less likely to be detected, leading to overestimation of the *median* RAAS for lowly abundant proteins, which may induce indirect correlations of features associated with protein abundance.

The strong association of RAAS values with protein sequence conservation led us to explore the association with allele frequency variation in the human population using the gnomAD database<sup>45</sup> (see Methods). We estimated the population frequency of missense variants in alternatively translated codons. The results for all possible allele variants within alternatively translated codons (Fig. 5e) and for the subset of nucleotides that may lead to the incorporated amino acid (Extended Data Fig. 10d) indicate that most alternate translation sites have no observed allele variants in gnomAD. This result provides further evidence that the substitutions identified in our data do not arise from genetic polymorphisms translated according to the genetic code. Furthermore, we found that codons corresponding to high RAAS substitutions have more allele variation than expected, suggesting that they are less constrained than codons corresponding to low RAAS substitutions, Fig. 5f. This trend is even stronger when considering only the subset of alleles that may encode the incorporated amino acid, Extended Data Fig. 10e. This result reinforces the association with protein conservation, namely that alternate translation occurs with higher frequency at sites that are more tolerant to sequence variation.

#### **Conservation across species**

Lastly, we assessed conservation of alternate translation across species. We ran our pipeline from Fig. 1a on label free proteomics data from 3 mouse tissues<sup>14</sup> and identified 1,708 SAAPs, 89 of which are identical with human SAAPs analyzed above, Extended Data Fig. 10f. This intersection demonstrates conservation, though it likely underestimates the degree of conservation since (1) we did not account for protein sequence homology differences and (2) technical factors affect peptide detectability. To more quantitatively assess the conservation of alternate translation, we further analyzed RAAS values. Their significantly high correlation ( $p < 10^{-10}$ ) across human and mouse strongly support species conservation, Fig. 5g. The total variance in substitution ratios of these shared SAAP are most significantly driven by substitution type, followed by species and tissue

type, Fig. 5h.

# Discussion

Our results demonstrate that alternate translation produces thousands of stable and abundant proteins in human and mouse tissues. These proteoforms differ in sequence, similar to proteoforms resulting from missense mutations. However, there are important differences. First, mutations introduce about 44 missense substitutions per cancer<sup>46</sup> while alternate translation introduces thousands of sequence changes. Second, mutations affect all protein copies templated by the mutated allele while alternative decoding affects only a fraction, though the fraction can be large and amount to  $\mathcal{O}(10^4)$  copies per cell, Fig. 2h. Furthermore, some substitutions are significantly more abundant in tumors (Fig. 4a), suggesting a disease association. Another potential link to disease are the highly abundant substitutions in numerous proteins associated with neurodegenerative diseases, such as TDP43, FUS and VCP. These substitutions are identified across all datasets and are reported in Supplemental Data Table 6.

Multiple mechanisms likely contribute to the amino acid substitutions that we quantified. Many substitutions, especially those with low ratios, likely reflect limited fidelity of aminoacyl-tRNAs recognizing their cognate codons. This mechanism is supported by higher ratios for aminoacyl-tRNAs pairing based on fewer nucleotide mismatches (Fig. 2i and Extended Data Fig. 3i) and for rare codons, Fig. 2j, Extended Data Fig. 3k, and Extended Data Fig. 4. High ratio substitutions likely involve sense codon recoding that may involve tRNA and mRNA modifications. Indeed, mRNA pseudouridylation recodes stop codons<sup>17,18</sup> and cytidine acetylation affects translational fidelity<sup>47</sup>. We found that U modifications overlap significantly ( $p < 10^{-10}$ ) with the amino acid substitutions reported here, Supplemental Fig. 4. This is consistent with its potential to recode sense codons previously observed with exogenous synthetic constructs<sup>48</sup>. Such recoding might also be regulated by UTR sequences altering codon pairing analogous to the way the Selenocysteine Insertion Sequence alters the reading of the UGA stop codon into incorporating selenocysteine<sup>49</sup>, though our data do not provide direct evidence for such UTR sequences. Similarly, modifications

of transfer RNA, ribosomes<sup>50–52</sup>, RNA binding proteins and structures delaying elongation may contribute to sense codon recoding<sup>53</sup>. A major mechanism determining the abundance of the substituted proteins is their stability, which we directly estimated based on protein degradation rates, Fig. 2.

Our results reveal only a subset of alternatively translated proteins because (1) we tested a limited sequence space (a single amino acid deviation from genetic code prediction), (2) we analyzed only a subset of all peptide ions, and (3) limitations exist in interpreting the mass spectra of analyzed ions<sup>54</sup>. Indeed, even the deepest MS datasets cannot achieve full protein sequence coverage<sup>55</sup>. In the data-dependent-acquisition datasets analyzed here, the peptide ions were selected for fragmentation in order of their abundance, and therefore less abundant peptides are less likely to be fragmented and identified. As a result, lowly abundant substituted peptides are less likely to be identified and their ratios are more likely to be missing from our data; thus, lower ratios are under-represented in our RAAS distributions. Even more limiting, only a minority of all fragmented peptide ions are assigned to confident sequence, which is a general challenge in MS proteomics<sup>8,9,56</sup>. This sequence assignment remains limiting when all detectable ions are sampled by data-independent-acquisition methods, as current methods can assign confident sequence to less than 25% of the peptide ions<sup>57</sup>.

The sequence assignment challenge is heightened by our priority on minimization of false positives, which likely increased false negatives. For example, we excluded 2,734 SAAP whose sequence might correspond to alternative reading frames and DNA sequences annotated as noncoding. Yet this hypothetical sequence correspondence is unlikely to represent a templating relationship. We also excluded thousands of SAAPs that did not pass our very stringent filters (Fig. 1) even though they are likely correctly identified. We chose these conservative and rigorous filters to increase confidence in the remarkably large number of abundant SAAPs. The mass spectra of some SAAPs with high ratios contain fragment ions corresponding to fragmentation of each peptide bond and are thus consistent with only one amino acid sequence, Supplemental Fig. 1 and 2. However, some mass spectra are incomplete and provide lower confidence. The substitutions of such sites are supported by fragmentation patterns (Supplemental Fig. 3) and their significant associations with mRNA modifications (Supplemental Fig. 4), protein disorder, conservation, and allele polymorphisms, Fig. 5. Future research is needed to balance the trade-offs between false positives and negatives, improve estimated of false discovery rates of SAAPs, and expand the set of confident SAAP whose biological significance can be interpreted.

Sites on lowly abundant proteins with low RAAS are less likely to be detected. This may results in non-ignorable missing data for some analysis, such as establishing the biological association between median RAAS per protein and protein abundances. To mitigate these confounding effects, we focused our conclusions on ratios between quantified substitution sites to establish ratio variation across patients and cancer types, Fig. 2e,f and Fig. 4a,b.

Protein products of alternate translation are challenging to detect since most interpretation of MS data is based on database-searching that tests hypotheses about the detection of protein products whose sequence matches nucleic acid sequences; hypotheses about the presence of protein products with sequences not predicated by the genetic code are not commonly tested. Furthermore, if such proteins are identified, they need to be further evaluated by a battery of additional tests, and downstream analysis as reported here. For example, previous reports of substitutions<sup>14,54</sup> did not test if they have genetic origin and did not focus on their validation, abundance and stability. We expect that advances in *de novo* protein sequencing and increased protein sequence coverage will identify many more substitutions. The abundance of many substitutions reported here is high enough (Fig. 2) to quantify them in single cells<sup>58,59</sup>, especially if their analysis is prioritized<sup>60</sup>.

Our results provide direct evidence for the abundance and stability of proteins with amino acid substitutions (Fig. 2) and associative evidence for their functional roles, Fig. 4 and Fig. 5. To avoid assuming functions or mechanisms, we chose the neutral phenomenological term "alternate translation". While some SAAPs may merely reflect limits of translational fidelity and proteostasis, others may have evolved biological functions as previously suggested<sup>15,16</sup>, consistent with regulated sense codon recoding discussed above. The high abundance, conservation across species,

and associations of some SAAPs with cancer and protein domains (Fig. 4) imply functional significance, and this possibility needs to be explored directly by future research.

The conservation of SAAPs between human and mouse (Fig. 4) contrasts with limited conversation between mammals and unicellular microorganisms, including yeast and bacteria. Previous analysis of substitutions in yeast and bacteria did not discover highly abundant substitutions<sup>13</sup>, and we confirmed this result.

# Methods

# Sample-specific protein databases

RNA-sequencing data was downloaded from NIH Genomic Data Commons (CPTAC data, portal.gdc.cancer.gov) or from dataset repository (label-free data, E-MTAB-2836). Reads were sorted by name and paired ends were written to two .fastq files (CPTAC only) with samtools v1.10<sup>61</sup>. Adapters were trimmed, reads were filtered for phred quality > Q28 and QC data was printed with fastp v0.23.4<sup>62</sup>. The remaining reads were aligned to GRCh38 reference genome with hisat2 v2.2.1<sup>63</sup> with –dta flag for downstream trancriptome assembly processing. The alignment was converted from .sam to .bam, sorted by position and indexed with samtools  $v1.10^{61}$ . Aligned reads were processed into a custom protein database with a 2-pronged approach: De-novo transcriptome assembly and single nucleotide variant (SNV)/insertion-deletion (indel) calling. First, the sorted and indexed hisat2 alignment was used to assemble a de novo transcriptome with stringtie v2.2.1<sup>64</sup>, which captures alternative splice variants. A transcript was called at a given locus if its abundance was at least 0.1% of the most abundant transcript at that locus, and if there was at least 1 independent read for that transcript. The assembly was compared to the reference genome using gffcompare v0.11.5<sup>65</sup> and the annotated result was filtered for coding regions (CDS) with gffread v0.12.7<sup>65</sup> and converted into .bed format using a Galaxy-sourced python script (gffcompare\_to\_bed.py)<sup>27</sup>. A second Galaxy-sourced script (translate\_bed.py) was used to translate the CDS regions of the annotated de novo assembled transcriptome into protein sequences (based

on the genetic code) and generate a fasta database of the translated sequences. Second, SNV and indel variants were called using freebayes v1.3.4<sup>66</sup> with the sorted hisat2 alignment as input, and default parameters. customProDB v1.41.0<sup>67</sup> was then run to translate the SNV and indel calls into protein fasta databases. Fasta databases from the two variant calling processes were merged into a single protein database for each sample using another Galaxy-sourced python script (fasta\_merge\_files\_and\_filter\_unique\_sequences.py)<sup>27</sup>. Sample-specific protein databases for all samples in a TMT experiment (CPTAC data) or all samples of a given tissue type (label-free data) were further merged. This pipeline was implemented on the Linux command line using the Discovery high performance computing cluster (MGHPCC, Holyoke, MA).

#### Blast

All proteins in the custom protein databases were blasted against the RefSeq human proteome using the blastp command from ncbi-blast+ command line tool at:

ftp.ncbi.nlm.nih.gov/blast/executables/blast+/2.16.0/). The top matches for each protein were used to identify the most likely source of ambiguous proteoform translations.

# Proteomics data processing

LC-MS proteomics data was downloaded from NIH Proteomic Data Commons (CPTAC data, proteomic.datacommons.cancer.gov/pdc/) or from dataset repository (label-free human data, PXD010154, and label-free mouse data, PXD030983), in the form of .raw files. For the discovery search, raw LC-MS data for each TMT set (or tissue for label-free data) was searched against the corresponding custom protein database using MaxQuant v1.6.17.0<sup>6</sup> with dependent peptide (DP) search option set to True and dependent peptide FDR set to 1%. For the validation search, the SAAP sequences were appended to the sample-specific FASTA file, the DP search option was set to False, and match between runs was enabled. PSM, peptide and protein FDR were all set to 1%. To run DP search with TMT-labeled data, we set LCMS run type to "Standard" and included TMT labels as a fixed modification at the N-terminus and lysine residues of peptides. To search label-free data, the TMT modification was omitted. Cysteine carbamidomethylation was included as a fixed modification, and methionine oxidation and N-terminal acetylation were added as variable modifications. Other parameters were set as defined by the original publications and were consistent across CP-TAC datasets. See project Github for a sample configuration file. Searches were implemented on the Linux command line using the Discovery high performance computing cluster (MGHPCC, Holyoke, MA).

#### **Candidate SAAP and PTM identification**

DP search results were mined for peptides exhibiting a mass shift that corresponds to either a known PTM (unimod.org) or to a potential AAS using custom python scripts including functions adapted from ref.<sup>13</sup>. Only a single modification per peptide was considered. The mass error threshold between the observed mass shift and the mass shift of a given modification was required to be below 5 ppm. For PTMs, we required that an amino acid known to be modified by that PTM was present in the peptide sequence. For AAS, we further required that that amino acid was assigned a positional probability of modification by MaxQuant. Peptides that could be explained by both an AAS and a PTM were not considered as candidate SAAP, e.g., we filtered out  $Q \rightarrow E$  and  $N \rightarrow D$  as such changes may arise as part of deamidation. Base peptides were allowed to have more than one associated modified peptide. The few identified peptides with substitutions of K and R are filtered out due to the potential impact of miscleavage on abundance estimation. Peptide sequences with potential AAS were compared against a six-frame translation of all regions of the human genome and any matches were discarded as potential products of non-canonical translation.

## Mass shift localization

When the MS2 spectra of SAAPs contain fragments generated from cleavages of the precursor ions just before and just after the substituted amino acid, the mass shift can be confidently assigned to a single amino acid residue. The confidence of such localization is reflected in the positional probability estimated by MaxQuant, Supplemental Fig. 3a. When the positional probability is close to 1, the measured mass shift can be confidently assigned to single residue, a substituted amino acid or another modification of the residue. However, MS2 spectra often do not contain a

full set of fragments, especially from lowly abundant peptides, sequences that fragment poorly or small fragment ions (e.g., b1) below the low limit of MS2 spectra.

In such cases, we used fragmentation patterns to determine whether an observed mass spectrum with uncertain positional probability aligns better the proposed SAAP sequence or with an alternative hypothesis that can explain the observed mass shift relative to the BP. Specifically, we predicted the fragment ion spectra for the SAAP and the peptide representing the alternate hypothesis using Prosit<sup>68</sup> as implemented in Koina<sup>69</sup> and shown in Supplemental Fig. 3b-f. For TMT-labeled peptides the Prosit\_2020\_intensity\_TMT model<sup>70</sup> was used with HCD fragmentation, and for label-free peptides the Prosit\_2020\_intensity\_HCD<sup>68</sup> model was used. Predicted spectra for the SAAP and the alternate hypothesis peptides were compared to the empirical spectra with a cosine similarity score computed with the CosineGreedy function from the matchms.similarity python package. Koina was also used to predict spectra for high-confidence SAAPs and their corresponding BPs as shown in Supplemental Fig. 1 and 2.

## **Rescoring SAAP confidence**

To validate if the discovered SAAP represent genuine peptides that can not be explained by more confident canonical PSMs, closed validation search outputs from MaxQuant v2.4.3 for the entire CCRCC cohort were rescored using the Prosit model and Percolator within the Oktoberfest pipeline<sup>8</sup>. Unfiltered database search results performed at 100% FDR including decoy PSMs as well as raw files were used as input. Normalized collision energy (NCE) calibration of the model was performed in the range 19-50. Peptide intensity prediction was performed with the model "Prosit\_2020\_intensity\_HCD" and indexed retention times were predicted with "Prosit\_2019\_irt". Matching of predicted intensities was performed using previously described similarity measures such as spectral angle and Pearson correlation. False discovery rate was estimated using the SVM Percolator<sup>71</sup> with its standard settings.

# Quantifying SAAP, BP, and RAAS

For CPTAC datasets, peptide abundance was computed at 2 levels: precursor and reporter ion. Precursor ion abundance corresponds to the total intensity of a peptide across all samples in a given mutiplexed TMT-labeled experiment (experiment-level). This value is reporter by MaxQuant in the "Intensity" column of the evidence.txt output file. Reporter ion abundance corresponds to the intensity of a peptide in an individual patient sample (sample-level). This value is computed by distributing the precursor ion intensity by the fractional contribution of each reporter ion in the experiment to the total reporter ion intensity for that peptide. This method was used to calculate abundance of all peptides reported here, including SAAP, BP, and other peptides. For the labelfree dataset, only precursor-level peptide abundances were computed. If the peptide is quantified in multiple charge states and across different off-line fractions, the intensities of all instances were summed up to estimate the overall abundance. Ratios of amino acid substitution (RAAS) were computed by dividing the SAAP abundance by the BP abundance and are generally reported on a  $\log_{10}$  scale. If the substitution impacts the cleavage probability, the distribution of the peptide across fully cleaved and miscleaved forms may be altered and impact RAAS estimation. To limit the impact of such artifacts, we removed from our analysis all base peptides that were also identified as part of longer miscleaved peptides, Extended Data Fig. 1h(i). Precursor ion, experimentlevel SAAP quantification data can be found in Extended Data Table 3. Reporter ion, sample-level SAAP quantification data can be found in Extended Data Table 4.

#### **Ionization efficiency estimation**

Ionization efficiencies of peptides were computed according to the model published by Liigand, *et al.* [29]. The authors established that for small peptides, such as those produced by trypsin digestion, the ionization efficiencies can be approximated as the sum of the ionization efficiencies of the amino acids comprising the peptide. They also provided empirical measures of the ionization efficiency of each amino acid. Therefore, the ionization efficiency of a SAAP or BP is estimated by the sum of the ionization efficiencies of the amino acids that make up the peptide.

#### **Evaluating the internal consistency of RAAS estimates**

To test the reliability of high RAAS, we evaluated their consistency with the correlation patterns of BP or SAAP abundance to shared peptide abundance across patients, Extended Data Fig. 3i. This test capitalized on the internal consistency of different peptides from the encoded (P1) and alternatively translated (P2) proteoforms, thus allowing triangulation across multiple measurements, Fig. 2. Yet, this power is diminished by the uncertainty of proteoform composition and thus the shared peptide associated with each pair of BP and SAAP. This correlation analysis used the intensities of sample-specific reporter ions that are unlikely to share acquisition artifacts with the precursor intensities used by our consistency test described in Fig. 2d. The data indicate that as RAAS increases, suggesting reduced relative proportion of the P1 proteoforms, the correlations between BP and shared peptides, CORR(BP, shared peptides), significantly decrease, Extended Data Fig. 3i (top panel). Meanwhile, the correlations between SAAP and shared peptides, CORR(SAAP, shared peptides), increase slightly for peptides with RAAS>0, Extended Data Fig. 3i (middle panel) and the difference between CORR(SAAP, shared peptides) and CORR(BP, shared peptides) increases with increasing substitution ratios, Extended Data Fig. 3i (bottom panel). Many SAAP with high substitution ratios correlate to shared peptides more strongly than their counterpart BP, an effect most significant for peptides with RAAS>1 (Fig. 3i). This trend indicates that our RAAS estimates are consistent with triangulated measurements across shared peptides, despite the caveat of uncertain proteoform structures.

## Validation in multiple digests

To establish confidence in the data and evaluate potential impacts of the protease used for protein digestion, we used the tonsil sample from the label-free dataset<sup>26</sup>, which the authors subjected to digestion by several proteases. Each digest was processed through the SAAP identification, validation and quantification pipeline described above. Substituted peptides found in each digest were mapped back to their protein to match substitution sites across digests and compute the percentage of substitutions identified in multiple digests. To determine the percentage of base peptides iden-

tified in multiple digests, we chose a random amino acid from the tryptic base peptide to search whether it is detected as part of a peptide identified from the alternative proteases.

# **SAAP degradation rates**

Raw LC-MS proteomics data from SILAC-labeled human cell lines ref.<sup>35</sup> was downloaded from the publication data repository (B cells: PXD008511, hepatocytes: PXD008512, monocytes: PXD008513, NK cells: PXD008515) and processed through the SAAP identification pipeline described above. Since there was no corresponding transcriptomic data, the DP search in MaxQuant was conducted with the human protein sequence fasta, downloaded from Uniprot. Lys8 and Arg10 were included as additional variable modifications to account for the SILAC-introduced heavy labels. After identifying candidate SAAP, these sequences were appended to the Uniprot database, along with the validated SAAP from the CPTAC and label-free<sup>26</sup> datasets. To increase confidence in our results by using an independent search engine, the validation search was done using MS-Fragger<sup>36</sup> with the SILAC3 workflow and MSBooster<sup>37</sup> enabled. All identified heavy and light peptides, including validated SAAP, BP and other peptides, were quantified by summing the precursor ion intensity over all charge states and modifications (other than heavy isotopes) and taking the median across the replicates. RAAS was calculated as described above. Degradation rates ( $\alpha$ ) were computed using linear regression as the slope of the  $\ln(1 + h/l) = \alpha t$ , where h/l is the ratio between the heavy (h) and light (l) peptides abundance across time.

# Protein set enrichment analysis

Protein RAAS values were computed as the median RAAS of all SAAP mapping to the protein. Gene ontology (GO, geneontology.org) biological processes enriched with high RAAS were identified by comparing the distribution of RAAS for proteins in a given biological process to the distribution of RAAS for all substituted proteins using a Kolmogorov-Smirnov test to obtain p-values, which were then adjusted with a Benjamini-Hochberg FDR correction. GO processes with identical sets of substituted proteins were manually combined. Proteins representing significant biological processes are highlighted in Fig. 2g.

#### Estimation of substituted protein copy number

Substituted protein copy numbers were estimated using the histone ruler approach introduced in ref.<sup>30</sup>. Briefly, the abundance of histones was computed as the median intensity of peptides corresponding to core histone proteins (H2A, H2B, H3, H4). This value was assumed to represent 30e6 protein copies, the approximate number of histone protein copies in a cell. SAAP protein copy numbers were estimated by multiplying the SAAP abundance, computed from precursor ion intensities as described above, by 30e6/histone protein abundance.

# **RAAS dotplots**

The global distribution of  $log_{10}(RAAS)$  is approximately log-normal. Therefore, the distributions for various subsets of the observed substitutions (substitutions corresponding to a specific codon, substitution type or amino acid property, all substitutions in a given protein or a GOslim annotated group of proteins, and substitutions classified by disorder and conservation score bins), were compared to the distribution of all other values from the same distribution (globally for Extended Data Fig. 4b,d, Fig. 5c and Extended Data Fig. 5a, or within cancer types for Fig. 3b,c,e, Fig. 4c, Extended Data Fig. 4c and Extended Data Fig. 5c,f) by two-sided Student's t-tests (R's t.test function with the Welch approximation, called from our custom function raasProfile). The p-value, sidedness (test statistic t < 0 or t > 0) and  $\log_{10}$  of the median RAAS value were recorded, and used to generate the RAAS dotplots (function dotprofile), where the dot size scales with  $-\log_{10}(p)$  up to a cut-off, and dot colors scale with the median RAAS value, as indicated by individual Figure legends. The number of RAAS values that were considered for each row and column group are indicated on the top and right axes. RAAS dotplots in the main manuscript often show only a subset of the analyzed categories filtered for having at least one significantly high or low RAAS distribution, as indicated in the Figure captions, and using the sortOverlaps function from the segmenTools R package (https://github.com/raim/segmenTools/ release RAAS\_preprint).

#### **Correlation of RAAS profiles**

Median RAAS across all SAAP with a given substitution type was computed for each substitution type in each dataset, a subset of which is displayed in Fig. 3b. Pearson correlations weighted by the number of SAAP associated with each substitution type were computed using vectors of these values for every pair of datasets. The resulting correlation heatmap of RAAS profiles for substitution types across datasets is shown in Fig. 2d. Unweighted Pearson correlations across datasets are similarly computed for RAAS profiles of encoded and incorporated amino acid types, as displayed in Extended Data Fig. 5f,g.

#### **ANOVA** analyses

To increase our confidence that variability in RAAS is biologically driven as opposed to driven by technical differences between datasets, we fit an ordinary least squares multiple regression model with precursor-level RAAS values for all SAAP in all datasets as the dependent variable and encoded amino acid type, incorporated amino acid type, substitution type, tissue type, and dataset type (TMT or label-free) as the independent variables. The contribution of each of the model variables to variance in RAAS was determined by an analysis of variance (ANOVA) test. The Fvalue, representing the fraction of explained variance, and p-value for each feature is displayed in Fig. 3f. A similar analysis was executed to determine the contribution of species (human/mouse) differences to RAAS variance relative to the contributions of substitution and tissue types, Fig. 5f.

## Substitution enrichment in tissue types

To determine enrichment of substitution types in specific tissues, we focused on tissues with both CPTAC and label-free healthy data, namely kidney (CCRCC), pancreas (PDAC), lung (LUAD, LSCC) and endometrium (UCEC). For each substitution type, we computed the log<sub>2</sub>(fold change) and p-value (t-test) between RAAS of all SAAP with that substitution type identified in a given tissue and the RAAS of all SAAP with that substitution type measured in all other tissues. p-values were adjusted with a Benjamini-Hochberg FDR correction to get q-values. Substitution types were considered significantly enriched with high or low RAAS values in a given tissue if the log<sub>2</sub>(fold

change) relative to all other tissues was > 1 or < -1, respectively, with a q-value of  $\leq 0.01$ . Substitution types with significantly high RAAS in a given tissue are highlighted in Extended Data Fig. 5h,i.

#### Comparison of substitution type RAAS in tumor and normal tissue

 $Log_2$  fold changes of RAAS between tumor and matched normal tissue for a given substitution type was computed by taking the mean of the  $log_2$  ratios of the RAAS for each SAAP with the substitution type in a patient tumor sample relative to the RAAS for the same SAAP in the patientmatched normal adjacent tissue sample. p-values were computed with a paired t-test. p-values were FDR-adjusted into q-values using the Benjamini-Hochberg method.

## Pfam domain analysis

The canonical isoforms of proteins with substitutions were mapped to Pfam structural protein domains using InterProScan v5.67-99.0, which was downloaded from EMBL-EBI (InterProScan) and run on the Discovery high performance computing cluster (MGHPCC, Holyoke, MA). The significance of finding amino acid substitutions in a given Pfam domain was computed by counting the number of substitutions identified in a domain and comparing to a bootstrapped distribution of substitution counts in the domain obtained by shuffling the position of the substitutions across the entire protein sequence. The p-value was computed as the probability of finding more substitutions in the domain than expected by chance (random shuffling of substitutions). Domains with significantly high or low RAAS were identified by comparing the distribution of RAAS for substitutions in a domain to the distribution of all RAAS values using a Kolmogorov-Smirnov test and Benjamini-Hochberg FDR correction. Pfam domains significantly enriched with substitutions of significantly high or low RAAS are displayed in Extended Data Fig. 6d.

# Substitution sequence context analysis

To analyze the sequence context of amino acid substitution sites, we generated a custom protein database of all proteins defined in the Ensembl human genome release GRCh38.110, supplemented

with all proteins with patient-specific mutations from the sample-specific protein database generation pipeline described above. All unique base peptides were searched against this database using blastp (NCBI's BLAST+, version 2.15.0+), and only blast hits with 100% sequence identity and no mismatches were considered. For these hits, we obtained the amino acid context (sequences surrounding the substitution sites) from the original Ensembl protein and the codon from the matched Ensembl transcript. If the substitution itself covered a patient-specific mutation the codon was not considered. GOslim annotations for all genes were downloaded via Ensembl BioMart queries on 2023-12-04. Relative codon frequencies were calculated from the full Ensembl transcripts of the set of all Ensembl proteins with identified AAS sites. The total count for each codon was divided by the total count of all codons for the same amino acid.

#### **Sequence Difference Logos**

Difference logos were generated with the R package DiffLogo<sup>72</sup> (version 2.26.0) using a customized DiffLogo function that allows for custom y-axis limits of the plots; and p-values were calculated with the packages' enrichDiffLogoObjectWithPvalues function, but are only shown for residues that were not part of the sequence selection (e.g. the substitution site when selecting sequences by encoded and incorporated amino acid). In total 7,086 unique sequences surrounding the substitution sites were considered, and for all difference logos, a selected subset of these sequences was compared to all other sequences.

# Sorted enrichment profiles

Amino acid enrichment around unique substitution sites (Extended Data Fig. 8a,b) and the location of substitution sites along base peptides (Extended Data Fig. 8(c)) were evaluated using functions of the R package segmenTools, as described in detail by Behle et al.<sup>73</sup>. Shortly: enrichments were calculated by cumulative hypergeometric distribution tests (function clusterCluster). Significantly enriched ( $p \le 10^{-5}$ ) amino acids or sites were sorted along the x-axis from top to bottom (function sortOverlaps). The field color scales with  $-log_2(p)$ , cut at  $p \le 10^{-10}$  and the white text indicates  $p \le 10^{-5}$  (function plotOverlaps).

#### **Protein structural features**

IUPred3<sup>74</sup> was used to calculate a "disorder" score for each protein in our database, with options -a to add the ANCHOR2 prediction of disordered binding sites, and -s medium to use the "medium" smoothing type. Pfam domains for each protein were predicted with HMMER3<sup>75</sup> (version 3.4 (Aug. 2023)), using a reporting cutoff -E 0.01 for a comprehensive output. Protein secondary structures were predicted by S4pred<sup>76</sup> (version 1.2.5, downloaded from https:// github.com/psipred/s4pred on 2024-01-26). Alternative scores for disordered (flDPnn<sup>77</sup> and disordered protein interaction sites (disordRDPbind<sup>78</sup>), as well as a MMseqs2<sup>79</sup>-based sequence conservation score were obtained *via* the DescribePROT database<sup>80</sup>. Download and calculation of protein structural features is documented at https://gitlab.com/raim/genomeBrowser/ -/tree/master/data/mammary (release RAAS\_preprint). ChimeraX<sup>81</sup> was used to display selected PDB structures and color substituted sites by their median RAAS value.

# **Predicting RAAS from sequence features**

To understand how much of the site to site variance on RAAS could be explained by known protein, local sequence, and amino acid site specific structural properties, we created an classifier trained on a variety of features to predict RAAS. XGboost<sup>82</sup>, R implementation version 1.7.7.1, was chosen for the classifier because it allowed for easily incorporating missing data features that were not completely known across the entire proteome. The data was randomly split into a train fraction, 80 %, and a test fraction, 20 % to evaluate model. The model was retrained 1000 times on random subsets of test and train data. The scatter plot was a representative sample reflecting the median correlation between train and test. The model was then retrained 100 times in the same manner leaving out one factor each time to estimate the added predictive power for each feature, defined as the decrease in correlation from the all feature model.

# Allele frequency of missense variants in alternatively translated codons

Allele frequencies in gnomAD v4.1.0  $(730,947 \text{ exomes})^{45}$  were calculated for two sets of missense substitutions: (1) all possible single-nucleotide missense variants in codons affected by alternative

translation (Fig. 5e), and (2) single-nucleotide mutations in the codon that result in the alternate amino acid when translated according to the genetic code (Extended Data Fig. 10d). Codons were grouped into quartiles based on their RAAS values. The site frequency spectrum (SFS) for each group was binned into eight allele frequency categories: 0 (monomorphic), 0 to  $10^{-6}$ ,  $10^{-6}$  to  $10^{-5}$ ,  $10^{-5}$  to  $10^{-4}$ ,  $10^{-4}$  to  $10^{-3}$ ,  $10^{-3}$  to  $10^{-2}$ ,  $10^{-2}$  to  $10^{-1}$ , and  $10^{-1}$  to 0.5.

## Missense variant constraint in alternatively translated codons

Constraint was calculated for three sets of missense substitutions in the gnomAD v4.1.0 dataset (730,947 exomes)<sup>45</sup>: (1) all possible single-nucleotide missense variants in codons affected by alternative translation (Fig. 5f), (2) single-nucleotide mutations that lead to the observed SAAP sequences when translated according to the genetic code (Extended Data Fig. 10e), and (3) singlenucleotide mutations that do not lead to the alternate amino acid in SAAP when translated according to the genetic code. For each of these sets, across all RAAS quartiles, constraint was measured by dividing the observed number of mutations by the expected number based on the Roulette mutational model<sup>83</sup>. This model predicts the expected number of mutations at each site by considering nucleotide context and known mutational processes in the human genome. Two control groups were used for comparison within each RAAS quartile. The first control group contained the most deleterious missense variants in genes from this RAAS quartile, as predicted by AlphaMissense<sup>84</sup> (top decile). These variants had an observed/expected ratio close to 0.5 indicating strong constraint (twice as constrained as synonymous mutations). The second control group contained the least deleterious missense variants in genes from this RAAS quartile (lowest decile), with an observed/expected ratio near 1 indicating minimal constraint (no more constrained than synonymous mutations). AlphaMissense is a machine learning-based tool that predicts the deleteriousness of missense variants by integrating functional annotations and evolutionary conservation data. The use of these control groups allows for benchmarking the constraint observed in the sets of interest against variants with known levels of deleteriousness.

#### **Overlap with RNA modification sites**

To explore whether uracil (U) modifications may contribute to alternate RNA decoding, we analyzed the nucleotide overlaps of codons corresponding to amino acid substitutions with sites of modified U residues in transcripts, as previously defined by nanopore sequencing<sup>34</sup>. We mapped all modified U sites from 6 cell lines to the coding sequences of the set of Ensembl MANE transcripts (Matched Annotation from NCBI and EBI), yielding 39,723 unique sites in 6,771 distinct transcripts. We also reduced the set of AAS from 7,069 to 6,967 sites that map to MANE transcripts. We then counted all U in the union of all 7,471 transcripts with AAS and/or  $\psi$  sites as the total set, and all 4,250 U in codons at AAS sites as the test set, and found 180 overlapping sites. To test whether this is a significant enrichment we used a cumulative hypergeometric distribution test (p[X > 179], in R: phyper (q=180-1, m=39723, n=2879635-39723, k=4250, lower.tail=FALSE). To further evaluate the calculated p-value, we plotted the p-value distribution over different counts in panel and indicate the real count and its associated p-value in red, Supplemental Fig. 4a.

#### Availability

Supporting information, data, and documentation is available at decode.slavovlab.net. The code is freely available at github.com/SlavovLab/decode.

#### Acknowledgments

We thank Dr. Mahlon Collins and Eunice Koo for help with analysis, Prof. Nuno Bandeira, Prof. Arjun Raj, Prof. Zoya Ignatova and Prof. Barry Karger for detailed feedback, and members of the Slavov laboratory for discussions and suggestions. The work was funded by an Allen Distinguished Investigator award through The Paul G. Allen Frontiers Group to N.S., an NIGMS award R01GM144967 to N.S., NCI awards UG3CA268117 and UH3CA268117 to N.S., and a MIRA award R35GM148218 from the NIGMS of to N.S.

# **Competing interests**

Nikolai Slavov is a founding director and CEO of Parallel Squared Technology Institute, which is a non-profit research institute.

# **Author contributions**

Study design, supervision, and raising funding: N.S.
Data analysis: S.T., R.M., A.L., S.W., and N.S.
gnomAD analysis: J.G., S.T., and K.K.
Initial draft: S.T., and N.S.
Writing: All authors approved the final manuscript.

# References

- 1. Cantwell-Dorris, E. R., O'Leary, J. J. & Sheils, O. M. BRAFV600E: implications for carcinogenesis and molecular therapy. *Molecular cancer therapeutics* **10**, 385–394 (2011).
- Hart, J. R. *et al.* The butterfly effect in cancer: a single base mutation can remodel the cell. en. *Proc. Natl. Acad. Sci. U. S. A.* **112,** 1131–1136. ISSN: 0027-8424, 1091-6490. http: //dx.doi.org/10.1073/pnas.1424012112 (Jan. 2015).
- 3. Wright, A. & Vissel, B. The essential role of AMPA receptor GluR2 subunit RNA editing in the normal and diseased brain. *Frontiers in molecular neuroscience* **5**, 34 (2012).
- Parker, J. & Friesen, J. D. "Two out of three" codon reading leading to mistranslation in vivo. *Molecular and General Genetics MGG* 177, 439–445 (1980).
- Savitski, M. M., Nielsen, M. L. & Zubarev, R. A. ModifiComb, a new proteomic tool for mapping substoichiometric post-translational modifications, finding novel types of modifications, and fingerprinting complex protein mixtures. en. *Mol. Cell. Proteomics* 5, 935–948. ISSN: 1535-9476. http://dx.doi.org/10.1074/mcp.T500034-MCP200 (May 2006).

- Cox, J. & Mann, M. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. en. *Nature Biotechnology* 26, 1367–1372. ISSN: 1546-1696. https://doi.org/10.1038/nbt.1511 (Nov. 2008).
- 7. Wilhelm, M. *et al.* Deep learning boosts sensitivity of mass spectrometry-based immunopeptidomics. *Nature communications* **12**, 3346 (2021).
- Picciani, M. *et al.* Oktoberfest: Open-source spectral library generation and rescoring pipeline based on Prosit. *Proteomics*, e2300112 (Sept. 2023).
- 9. Yang, K. L. *et al.* MSBooster: improving peptide identification rates using deep learningbased features. *Nature Communications* **14**, 4539 (2023).
- Pataskar, A. *et al.* Tryptophan depletion results in tryptophan-to-phenylalanine substitutants. *Nature* 603, 721–727 (2022).
- 11. Yang, C. *et al.* Arginine deprivation enriches lung cancer proteomes with cysteine by inducing arginine-to-cysteine substitutants. *Molecular Cell* **84**, 1904–1916 (2024).
- K, M. & M, I. Translational fidelity and mistranslation in the cellular response to stress. *Nature Microbiol.* 2 (2017).
- Mordret, E. *et al.* Systematic detection of amino acid substitutions in proteomes reveals mechanistic basis of ribosome errors and selection for translation fidelity. *Molecular Cell* 75, 427–441 (2019).
- 14. Giansanti, P. *et al.* Mass spectrometry-based draft of the mouse proteome. *Nature Methods* 19, 803–811. ISSN: 1548-7105. https://doi.org/10.1038/s41592-022-01526-y (2022).
- 15. Sun, L. *et al.* Evolutionary gain of alanine mischarging to noncognate tRNAs with a G4: U69 base pair. *Journal of the American Chemical Society* **138**, 12948–12955 (2016).
- 16. Netzer, N. *et al.* Innate immune and chemically triggered oxidative stress modifies translational fidelity. *Nature* **462**, 522–526 (2009).

- 17. Karijolich, J. & Yu, Y.-T. Converting nonsense codons into sense codons by targeted pseudouridylation. *Nature* **474**, 395–398 (2011).
- Dai, Q. *et al.* Quantitative sequencing using BID-seq uncovers abundant pseudouridines in mammalian mRNA at base resolution. *Nature biotechnology* 41, 344–354 (2023).
- Cox, J. *et al.* Andromeda: a peptide search engine integrated into the MaxQuant environment. *Journal of proteome research* 10, 1794–1805 (2011).
- Clark, D. J. *et al.* Integrated proteogenomic characterization of clear cell renal cell carcinoma. *Cell* 179, 964–983 (2019).
- 21. Krug, K. *et al.* Proteogenomic Landscape of Breast Cancer Tumorigenesis and Targeted Therapy. *Cell* 183, 1436–1456.e31. ISSN: 0092-8674. https://www.sciencedirect.com/science/article/pii/S0092867420314008 (2020).
- 22. Gillette, M. A. *et al.* Proteogenomic Characterization Reveals Therapeutic Vulnerabilities in Lung Adenocarcinoma. *Cell* 182, 200–225.e35. ISSN: 0092-8674. https://www. sciencedirect.com/science/article/pii/S0092867420307443 (2020).
- 23. Dou, Y. et al. Proteogenomic Characterization of Endometrial Carcinoma. Cell 180, 729– 748.e26. ISSN: 0092-8674. https://www.sciencedirect.com/science/article/ pii/S0092867420301070 (2020).
- Cao, L. *et al.* Proteogenomic characterization of pancreatic ductal adenocarcinoma. *Cell* 184, 5031–5052.e26. ISSN: 0092-8674 (2021).
- Satpathy, S. *et al.* A proteogenomic portrait of lung squamous cell carcinoma. *Cell* 184, 4348–4371.e40 (2021).
- 26. Wang, D. *et al.* A deep proteome and transcriptome abundance atlas of 29 healthy human tissues. *Molecular systems biology* **15**, e8503 (2019).
- Batut, B. *et al.* Community-Driven Data Analysis Training for Biology. *Cell Systems* 6, 752–758.e1. ISSN: 2405-4712 (2018).
- Ma, C. *et al.* Improved Peptide Retention Time Prediction in Liquid Chromatography through Deep Learning. *Analytical Chemistry* 90, 10881–10888 (2018).

- 29. Liigand, P., Kaupmees, K. & Kruve, A. Influence of the amino acid composition on the ionization efficiencies of small peptides. *Journal of Mass Spectrometry* **54**, 481–487 (2019).
- Wiśniewski, J. R., Hein, M. Y., Cox, J. & Mann, M. A "proteomic ruler" for protein copy number and concentration estimation without spike-in standards. *Molecular & cellular proteomics* 13, 1535–9484 (2014).
- 31. Wu, Q. *et al.* Translation affects mRNA stability in a codon-dependent manner in human cells. *Elife* **8** (Apr. 2019).
- 32. Drummond, D. A. & Wilke, C. O. Mistranslation-induced protein misfolding as a dominant constraint on coding-sequence evolution. *Cell* **134**, 341–352 (2008).
- Quax, T. E., Claassens, N. J., Söll, D. & van der Oost, J. Codon bias as a means to fine-tune gene expression. *Molecular cell* 59, 149–161 (2015).
- 34. McCormick, C. *et al.* mRNA psi profiling using nanopore DRS reveals cell type-specific pseudouridylation. *bioRxiv* (May 2024).
- 35. Mathieson, T. *et al.* Systematic analysis of protein turnover in primary cells. en. *Nat. Commun.* 9, 689. ISSN: 2041-1723. http://dx.doi.org/10.1038/s41467-018-03106-1 (Feb. 2018).
- 36. Kong Andy, T., Leprevost, F. V., Avtonomov, D. M., Mellacheruvu, D. & Nesvizhskii, A. I. MSFragger: ultrafast and comprehensive peptide identification in mass spectrometry-based proteomics. *Nature Methods* 14, 513–520. ISSN: 1548-7105. https://doi.org/10.1038/nmeth.4256 (2017).
- Yang, K. L. *et al.* MSBooster: improving peptide identification rates using deep learning-based features. *Nature Communications* 14, 4539. ISSN: 2041-1723. https://doi.org/10.1038/s41467-023-40129-9 (2023).
- 38. Veredas, F., Canton, F. & Aledo, J. Methionine residues around phosphorylation sites are preferentially oxidized in vivo under stress conditions. *Sci Rep* **7**, 40403 (Jan. 2017).
- 39. Barik, S. The Uniqueness of Tryptophan in Biology: Properties, Metabolism, Interactions and Localization in Proteins. *Int J Mol Sci* **21** (Nov. 2020).

- 40. Bartok, O. *et al.* Anti-tumour immunity induces aberrant peptide presentation in m elanoma. *Nature* **590**, 332–337 (Feb. 2021).
- Holecek, M. Why Are Branched-Chain Amino Acids Increased in Starvation and Diabetes? *Nutrients* 12 (Oct. 2020).
- 42. Mayers, J. *et al.* Elevation of circulating branched-chain amino acids is an early event in human pancreatic adenocarcinoma development. *Nat Med* **20**, 1193–1198 (Oct. 2014).
- Jiang, Z., Zheng, J., Liu, J., Yang, X. & Chen, K. Novel Branched-Chain Amino Acid-Catabolism Related Gene Signature for Overall Survival Prediction of Pancreatic Carcinoma. *J Proteome Res* 21, 740–746 (Mar. 2022).
- 44. Ghosh, K. & Dill, K. Computing protein stabilities from their chain lengths. *Proc Natl Acad Sci U S A* **106**, 10649–10654 (June 2009).
- Chen, S., Francioli, L. C., Goodrich, J. K. & et al. A genomic mutational constraint map using variation in 76,156 human genomes. *Nature* 625. Published correction appears in Nature. 2024 Feb;626(7997):E1. doi: 10.1038/s41586-024-07050-7, 92–100 (2024).
- 46. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancerassociated genes. *Nature* **499**, 214–218. ISSN: 1476-4687. https://doi.org/10. 1038/nature12213 (2013).
- Arango, D. *et al.* Acetylation of cytidine in mRNA promotes translation efficiency. *Cell* 175, 1872–1886 (2018).
- 48. Eyler, D. E. *et al.* Pseudouridinylation of mRNA coding sequences alters translation. *Proceedings of the National Academy of Sciences* **116**, 23068–23074 (2019).
- Berry, M. J., Banu, L., Harney, J. W. & Larsen, P. R. Functional characterization of the eukaryotic SECIS elements which direct selenocysteine insertion at UGA codons. *The EMBO Journal* 12, 3315–3322 (1993).
- Pan, T. Modifications and functional genomics of human transfer RNA. *Cell research* 28, 395–404 (2018).

- 51. Slavov, N., Semrau, S., Airoldi, E., Budnik, B. & van Oudenaarden, A. Differential stoichiometry among core ribosomal proteins. *Cell Reports* 13, 865–873. https://doi. org/10.1016/j.celrep.2015.09.056 (5 2015).
- Emmott, E., Jovanovic, M. & Slavov, N. Ribosome stoichiometry: from form to function. *Trends in biochemical sciences* 44, 95–109 (2019).
- 53. Dever, T. E., Dinman, J. D. & Green, R. Translation Elongation and Recoding in Eukaryotes. *Cold Spring Harbor Perspectives in Biology* **10** (2018).
- 54. Griss, J. *et al.* Recognizing millions of consistently unidentified spectra across hundreds of shotgun proteomics datasets. *Nature methods* **13**, 651–656 (2016).
- 55. Sinitcyn, P. *et al.* Global detection of human variants and isoforms by deep proteome sequencing. *Nature biotechnology* **41**, 1776–1786 (2023).
- 56. Slavov, N. Driving Single Cell Proteomics Forward with Innovation. Journal of Proteome Research 20, 4915–4918. https://doi.org/10.1021/acs.jproteome.1c00639 (2021).
- 57. MacCoss, M. J. *et al.* Sampling the proteome by emerging single-molecule and mass spectrometry methods. en. *Nat. Methods* 20, 339–346. https://www.nature.com/articles/s41592-023-01802-5 (Mar. 2023).
- 58. Gatto, L. *et al.* Initial recommendations for performing, benchmarking, and reporting singlecell proteomics experiments. *Nat. Methods* 20, 375–386. https://doi.org/10.1038/ s41592-023-01785-3 (2023).
- 59. Leduc, A., Khoury, L., Cantlon, J., Khan, S. & Slavov, N. Massively parallel sample preparation for multiplexed single-cell proteomics using nPOP. *Nature Protocols*, 1–27. https://doi.org/10.1038/s41596-024-01033-8 (2024).
- Huffman, R. G. *et al.* Prioritized mass spectrometry increases the depth, sensitivity and data completeness of single-cell proteomics. *Nat. Methods.* http://dx.doi.org/10.1038/s41592-023-01830-1 (2023).

- 61. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *GigaScience* **10.** giab008. ISSN: 2047-217X. https://doi.org/10.1093/gigascience/giab008 (Feb. 2021).
- Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor. Bioinformatics 34, i884–i890. ISSN: 1367-4803 (Sept. 2018).
- Kim, D., Pagg, J. M., Park, C., Bennett, C. & Salzberg, S. L. Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology* 37, 907–915. ISSN: 1546-1696 (2019).
- 64. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nature Biotechnology* 33, 290–295. ISSN: 1546-1696. https://doi.org/10.1038/nbt.
  3122 (2015).
- 65. Pertea, G. & Pertea, M. GFF Utilities: GffRead and GffCompare. *F10000Research* 9, 304 (2020).
- 66. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. *arXiv preprint* (2012).
- 67. Wang, X. & Zhang, B. customProDB: an R package to generate customized protein databases from RNA-Seq data for proteomics search. *Bioinformatics* **29**, 3235–3237 (2013).
- 68. Gessulat, S. *et al.* Prosit: proteome-wide prediction of peptide tandem mass spectra by deep learning. *Nature Methods* **16**, 509–518 (2019).
- 69. Lautenbacher, L. *et al.* Koina: Democratizing machine learning for proteomics research. *bioRxiv.* eprint: https://www.biorxiv.org/content/early/2024/06/03/ 2024.06.01.596953.full.pdf.https://www.biorxiv.org/content/ early/2024/06/03/2024.06.01.596953 (2024).
- Gabriel, W. *et al.* Prosit-TMT: Deep Learning Boosts Identification of TMT-Labeled Peptides. *Analytical Chemistry* 94, 7181–7190 (2022).
- Halloran, J. T. & Rocke, D. M. A Matter of Time: Faster Percolator Analysis via Efficient SVM Learning for Large-Scale Proteomics. *Journal of Proteome Research* 17, 1978–1982. ISSN: 1535-3907 (2018).

- 72. Nettling, M. *et al.* DiffLogo: a comparative visualization of sequence motifs. *BMC Bioinformatics* **16**, 387 (Nov. 2015).
- Behle, A. *et al.* Manipulation of topoisomerase expression inhibits cell division but not growth and reveals a distinctive promoter structure in *Synechocystis. Nucleic Acids Res* 50, 12790–12808 (Dec. 2022).
- Erdös, G., Pajkos, M. & Dosztanyi, Z. IUPred3: prediction of protein disorder enhanced with unambiguous experimental annotation and visualization of evolutionary conservation. *Nucleic Acids Res* 49, W297–W303 (July 2021).
- 75. Eddy, S. Accelerated Profile HMM Searches. *PLoS Comput Biol* 7, e1002195 (Oct. 2011).
- 76. Moffat, L. & Jones, D. Increasing the accuracy of single sequence prediction methods using a deep semi-supervised learning framework. *Bioinformatics* **37**, 3744–3751 (Nov. 2021).
- 77. Hu, G. *et al.* flDPnn: Accurate intrinsic disorder prediction with putative propensities of disorder functions. *Nat Commun* **12**, 4438 (July 2021).
- 78. Peng, Z. & Kurgan, L. High-throughput prediction of RNA, DNA and protein binding regions mediated by intrinsic disorder. *Nucleic Acids Res* **43**, e121 (Oct. 2015).
- 79. Steinegger, M. & Soding, J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat Biotechnol* **35**, 1026–1028 (Nov. 2017).
- Zhao, B. *et al.* DescribePROT: database of amino acid-level protein structure and function predictions. *Nucleic Acids Res* 49, D298–D308 (Jan. 2021).
- Meng, E. *et al.* UCSF ChimeraX: Tools for structure building and analysis. *Protein Sci* 32, e4792 (Nov. 2023).
- Chen, T. & Guestrin, C. XGBoost: A Scalable Tree Boosting System in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (Association for Computing Machinery, San Francisco, California, USA, 2016), 785–794. ISBN: 9781450342322. https://doi.org/10.1145/2939672.2939785.
- 83. Seplyarskiy, V. *et al.* A mutation rate model at the basepair resolution identifies the mutagenic effect of polymerase III transcription. *Nature Genetics* **55**, 2235–2242 (2023).

84. Cheng, J. *et al.* Accurate proteome-wide missense variant effect prediction with AlphaMissense. *Science* 381, eadg7492. https://www.science.org/doi/abs/10.1126/science.adg7492 (2023).

# **Extended Data Figures**



**Extended Data Fig. 1** | Systematic identification and validation of amino acid substitutions (a) Number of tumor and normal samples analyzed from each CPTAC dataset. (b) Number of samples analyzed for each healthy tissue from the label-free dataset. (c) Distribution of the percentage of each transcript with a read that is included in the patient-specific databases. (d) Distribution of the number of transcripts with 100% sequence coverage included in each patient-specific protein database. (e) Non-substitution modifications identified in the dependent peptide search are majorly comprised of post-translational modifications, and include artifacts and chemical derivatives from MS analysis. (f)– (j) (Continued on the next page)

(Continued) (f) Butterfly plots showing a systematic mass shift in MS2 spectra between SAAP and BP for a representative SAAP with median RAAS=1.2 in pericentriolar material 1 protein isoform 1 (PCM1). The fragmentation spectra were predicted by the Prosit TMT model<sup>70</sup>. (g) Cumulative density distributions of p-values (MaxQuant) and FDR-controlled q-values computed using only SAAP. Red dashed line indicates confidence threshold for SAAP inclusion in further analysis. (h) Over 80% of substitutions identified from lysine (K) or arginine (R) are at sites of missed cleavage or are substitutions between K and R. (i) Mass error distributions for SAAP and all peptides identified in the database search show no significant differences. (j) Observed and predicted (DeepRT+,<sup>28</sup>) retention times show strong agreement for all main peptides identified in standard database search and for SAAP.



**Extended Data Fig. 2** | **Establishing confidence in AAS abundance** (a) SAAP with high RAAS $\geq$ 0 are identified with the same FDR-controlled confidence as SAAP with low RAAS<0. (b) SAAP with high RAAS $\geq$ 0 are identified with as many fragment ions providing evidence at the site of alternate translation as SAAP with low RAAS<0. (c) SAAP with high RAAS $\geq$ 0 have similar mass error distributions as SAAP with low RAAS<0. (d-f) SAAP abundance estimates are unlikely to be affected by differences in ionization efficiency between base peptides and alternatively translated peptides. The ionization efficiency distributions for SAAP and BP are indistinguishable (d), correlate strongly (e), and there is negligible fold change between them (f). (g) SAAP abundances (RAAS) computed for the same site of alternate translation from peptides in different enzymatic digests of tonsil are consistent for peptides across the range of peptide abundances. (h)– (i) (Continued on the next page)

(Continued) (h) Base peptides with missed cleavages are generally an order of magnitude more lowly abundant than their fully cleaved counterparts. (i) Correlation of BP (top panel) or SAAP (middle panel) abundance with shared peptide abundance across samples. Shared peptides are peptides found in both the encoded and alternatively translated proteoforms. BP correlation to shared peptides decreases with increasing RAAS, while SAAP correlation to shared peptides tends to increase, especially at RAAS>1, as indicated by the difference in these correlations (bottom panel), and in support of the hypothesis presented in Fig. 2a. Abundances are computed with MS2-level intensities.



Extended Data Fig. 3 | Quantification of substituted amino acid peptides (SAAP) (a)– (k) (Continued on the next page)

(Continued) (a) Distributions of substitution ratios for all SAAP identified in each dataset computed for each experiment (TMT set, CPTAC data) or sample (label-free data), using MS1 precursor ion intensities. N indicates the number of RAAS computed at the MS1 level in each dataset. (b) Median RAAS were computed for each unique SAAP-BP pair using reporter ion intensities (MS2, CPTAC data) or precursor ion intensities (MS1, label-free data). Distributions of median RAAS across all SAAP in a dataset are shown. N indicates number of unique SAAP-BP pairs identified in each dataset. (c) Median RAAS across all SAAP identified in each sample were computed using reporter ion intensities (MS2, CPTAC data) or precursor ion intensities (MS1, label-free data). Distributions shown are of these medians across all samples in a dataset. N indicates the number of samples in each dataset. (d) Substitution ratio distributions shown in (a), (b), (c) have consistent medians, highlighting variability in RAAS across datasets. (e) Upset plot showing overlap in unique SAAP identified across all datasets. Dataset combinations require at least 10 shared SAAP to be included in visualization. (f) Heatmap displaying the percentage of samples in each dataset in which SAAP identified in 6+ datasets are found. Hierarchical clustering shows a cluster of shared SAAP that are commonly identified across majority of samples in addition to 6+ datasets. (g) To confirm variability in RAAS across datasets, we looked at the subset of SAAP that were identified in at least 1 sample in at least 6 datasets. LUAD and LSCC substitutions consistently have the lowest RAAS, while PDAC substitutions have the highest RAAS. N indicated the number of RAAS computed for shared SAAP in each dataset. (h) Boxplots highlighting the difference between RAAS in CPTAC datasets relative to RAAS computed in LUAD. Only SAAP shared between LUAD and the compared dataset are used. Each data point is a  $log_{10}(RAAS)$ difference computed for a unique SAAP-BP pair. (i) RAAS as a function of the minimum number of codonanticodon mismatches needed for incorporating the detected amino acid across all datasets. (i) An example of a substitution that can be partially explained by synthesis errors arising from significantly (t-test) higher abundance of the amino acyl-tRNA ligase supplying the alternatively translated amino acid relative to the abundance of the amino acyl-tRNA ligase supplying the encoded amino acid. \*: q-value <  $10^{-3}$ , \*\*: qvalue <  $10^{-5}$ , \*\*\*: q-value <  $10^{-20}$ . (k) RAAS negatively correlates to the codon stability coefficient, an empirical measure of codon usage. n denotes number of codons, r is Pearson correlation, p is correlation p-value, and the red line is the ordinary least squares fit. (I) RAAS distributions for SAAP identified and validated in human hepatocytes. (m) The stability of SAAP relative to BP in primary human B cells is inversely proportional to their RAAS (Pearson correlation). (n) Same as (m) but in NK cells.



**Extended Data Fig. 4** | Associations between codons, incorporated amino acids and RAAS (a) Relative codon frequencies in the full transcripts of the set of all proteins with identified substitution sites. The total count for each codon was divided by the total count of all codons for the same amino acid. Codon groups (per amino acid, separated by vertical lines) were sorted by amino acid property groups. Within each codon group, codons were sorted by their relative frequencies. All other panels are aligned with this sorting, see (e) for x-axis labels. These relative frequencies (bar heights) were also used in Fig. 2j (x axis). (b) RAAS dotplot for codons (columns) and incorporated amino acids (rows), sorted and color-coded by amino acid property groups. (c) RAAS dotplot for codons by datasets, i.e., cancer types and healthy tissues. (d) RAAS dotplot for all codons without further subsetting. Note, that the median RAAS values (colors) correspond to the y-axis values in Fig. 2j. (e) RAAS distributions for each codon.



Extended Data Fig. 5 | Substitution ratios depend on substitution and tissue types (a) RAAS dotplot for all encoded and incorporated amino acids. (b)– (h) (Continued on the next page)

(Continued) (**b**) Violin plots of RAAS medians for each substitution type in every dataset. N indicates number of substitution identified in each dataset. (**c**) RAAS dotplot as in (**a**) but by chemical properties of the encoded and incorporated amino acids. (**d**) Heatmap of median RAAS by substitution type for substitution types with variance <10% across datasets. (**e**) Heatmap of median RAAS by substitution type for substitution types with variance >50% across datasets. (**f**) RAAS dotplots for encoded (left panel) and incorporated (right panel) amino acids. (**g**) Median RAAS values for SAAP grouped based on the encoded amino acid (left panel) or incorporated amino acid (right panel) correlate strongly and significantly across all datasets (Pearson correlation). (**h**) Substitution types with significantly higher RAAS in a given tissue type (cancer and healthy samples) relative to all other tissues analyzed (t-test, Benjamini-Hochberg FDR-corrected). (**i**) RAAS distributions (boxplots) and number of SAAP identified (barplot) for substitution types that are significantly higher in a given tissue (colored) relative to all other tissues (gray) analyzed. Colors indicate the same tissue types as in (**h**).



**Extended Data Fig. 6** | Associations of substitutions with cancer (a) RAAS fold change between tumor and patient-matched normal adjacent tissue samples with median of distribution shown in red. N indicates the number of patient-specific RAAS values compared. (b) Patient-level RAAS distributions stratified by clinical tumor stage. No significant associations were measured between RAAS and tumor stage. (c) RAAS for the  $N \rightarrow G$  substitution in serine/threonine-protein phosphatase PP1-beta catalytic subunit (PP1CB) is significantly higher in tumor samples than in patient-matched normal adjacent tissue, for the majority of patients in LUAD and LSCC (t-test).



Extended Data Fig. 7 | Proteins with high RAAS organized by functional groups RAAS dotplots for all proteins with significantly high median RAAS ( $p \le 10^{-5}$ ) from each functional category shown in Fig. 4c.



**Extended Data Fig. 8** | Amino acid sequence context around substitution sites (a) Counts (text) and enrichments (gray scale: p-values of cumulative hypergeometric distribution tests) of amino acids surrounding the amino acid substitution sites. Lysine (K) and arginine (R) are enriched directly upstream of the substitution sites. Tryptophan (W), methionine (M), glycine (G) and cysteine (C) are enriched directly adjacent to substitution sites. (b)– (c) (Continued on the next page)

(b) As (a) but for substitution types (encoded:incorporated) vs. their position in identified peptides. Substitutions by glycine ( $\rightarrow$ G) or alanine ( $\rightarrow$ A) are enriched within the 3 N-terminal amino acids of base peptides (N1 to N3), i.e., directly after the trypsin cleavage sites (K or R). Various substitutions involving arginine (N), methionine (M) or glutamate (E) as either the encoded or the incorporated amino acid are enriched distant from the N- and C-termini (>9). Only substitution types with at least one significant enrichment ( $p < 10^{-10}$ ) are shown in (b). (c) Sequence difference logos were calculated for all unique sequences surrounding substitution sites, subset for all observed substitution types (encoded-incorporated amino acids), and plots were only generated if any of the positions -3 to +3 around a substitution site showed a significant enrichment with  $p \le 10^{-10}$  (\*:  $p \le 10^{-3}$ , \*\*:  $p \le 10^{-5}$ , \*\*\*:  $p \le 10^{-10}$ ), and all resulting logos are shown. The logos were grouped by common patterns (rows from top to bottom): (i) Substitutions by glycine or alanine ( $Q \rightarrow A$ .  $Q \rightarrow G, M \rightarrow G, L \rightarrow G$ ) are enriched directly upstream with lysine (K) or arginine (R), i.e. they are preferentially observed at the N-terminus of base peptides, next to the trypsin cleavage sites (K or R). (ii) Substitutions of glutamine (Q $\rightarrow$ A) or arginine (N $\rightarrow$ G) are flanked by cysteine (C) enrichments. (iii) substitutions E $\rightarrow$ N,  $T \rightarrow V$  and  $N \rightarrow M$  are flanked by methionine (M) enrichments. (iv) Substitutions  $E \rightarrow C$ ,  $I \rightarrow Q$ ,  $L \rightarrow Q$  and  $P \rightarrow T$ are flanked by tryptophan (W) enrichments. To generalize these grouped observations, we defined four sequence classes used for the difference logos in Fig. 4d, where sequence difference logos (each against all other sequences in our set) were calculated for all substitutions to G or A within the 3 N-terminal sites of base peptides (motif KRAQ), and for all substitutions that had at least one W (WWXWW), M (MMXMM), C (CCxCC) or G (GGxGG) within 2 positions of the substitution site. The selection GGxGG showed no specific enrichments at the substitution site and is not shown.



**Extended Data Fig. 9** | **Substitution proximity in domains, 1D and 3D protein structures** (a) Pfam domains significantly enriched in substituted peptides (FDR-adjusted p-values < 0.05, see Methods for details). (b) The number of base peptides decreases exponentially with the number of distinct SAAP detected per base peptide. This reflects the abundance bias of AAS detectability, such that many distinct SAAPs can be detected for highly abundant peptides. (c) High density region of substitutions in marginal zone B- and B1-cell specific protein (MZB1). A high RAAS substitutions cluster in  $\beta$ -sheet and unstructured regions is immediately followed by a lower RAAS cluster in an  $\alpha$ -helix region. (d) Three 1D clusters of substitutions are found in the glycolytic region of aldolase B (ALDOB). (e) Many substitutions are identified in the ribosome complex, some of which cluster across complex subunits in the 3D protein structure. (f) High and low RAAS substitutions cluster in the 3D protein structure of the proteasome complex.



**Extended Data Fig. 10** | Associations of substitutions with protein properties (a) Protein RAAS is negatively correlated to protein abundance. The intensity was calculated as the median over all "leading razor protein" intensities of all data sets (CPTAC and tissues) where a given base peptide was identified. (b) RAAS is positively correlated to the disordered score (IUPred3 prediction). RAAS was computed as the median RAAS for each unique BP/SAAP pair across all datasets (CPTAC and tissues) and the disordered score is the score of the protein at the AAS site. (c) as (b) but for the mean conservation score (MMseqs2 score via the DescribePROT database). (d) Allele frequency in the gnomAD database for missense variants coding for identified substitution in alternatively translated codons. Sites with allele variation frequency  $\geq 10^{-3}$  correspond to 131 SAAPs. (e) Observed / expected ratios for missense variants coding for identified substitution in alternatively translated codons, determined from analysis of gnomAD database (see Methods). Missense variants are less constrained with increasing RAAS ( $p < 10^{-6}$ ). Data point colors correspond to RAAS quartiles as in (d) or low or high AlphaMissense (AM) controls. (f) Upset plot showing overlap in unique SAAP identified between human and mouse tissues.

# **Supplemental Information**

# **Description of Supplemental Data tables**

- Supplemental\_Data\_1.PTMs.csv A table of peptides identified through the dependent peptide search as having known post-translational or chemical modifications. This table includes peptides from all 6 CPTAC datasets and the label-free healthy human tissue dataset. Types and locations of modifications are specified for each peptide.
- Supplemental\_Data\_2.SAAP\_proteins.xlsx A table containing 1 row per unique SAAP-BP pair per dataset (CPTAC and label-free). Each SAAP-BP pair is listed along with RAAS summary stats across the dataset and protein mapping information, including Pfam domains. The first tab of the file provides a description of the columns in the data table.
- **Supplemental\_Data\_3.SAAP\_precursor\_quant.xlsx** A table containing 1 row per unique SAAP-BP pair per TMT set (CPTAC) or healthy tissue type (label-free). Each SAAP-BP pair is listed along with their precursor ion level abundances and RAAS values computed from precursor ion intensities. The first tab of the file provides a description of the columns in the data table.
- Supplemental\_Data\_4.SAAP\_reporter\_quant.xlsx A table containing 1 row per unique SAAP-BP pair per patient sample (CPTAC only). Each SAAP-BP pair is listed along with their reporter ion level abundances and RAAS values computed from reporter ion intensities. The first tab of the file provides a description of the columns in the data table.
- **Supplemental\_Data\_5.SAAP\_SILAC\_quant.xlsx** Table of SAAP-BP pairs identified in SILAClabeled liver cells<sup>35</sup> containing one row per unique SAAP-BP pair per sample. Precursor ion level abundances of the light and heavy peptides in each sample are listed, along with their sums, ratios and RAAS values. The first tab of the file provides a description of the columns in the data table.

Supplemental\_Data\_6.SAAP\_neurodegenerative.xlsx A subset of the substitutions from

Supplemental\_Data\_2.SAAP\_proteins.xlsx corresponding to SAAP, BP and RAAS data pertaining to proteins that have been implicated in neurodegeneration and dementia (Uniprot).

- Supplemental\_Data\_7.SAAP\_coordinates.tsv Mapping of the amino acid substitution sites of unique pairs of BP/SAAP to proteins, their transcripts and genome coordinates, as defined in the Ensembl genome release GRCh38.110.
- **Supplemental\_Data\_8.High\_confidence\_SAAP\_precursor\_quant.xlsx** Supplemental\_Data\_3 filtered for SAAP with RAAS>0.1 and positional probability>0.9.
- Supplemental\_Data\_9.gnomAD.xlsx Data tables with gene locus, allele frequency and constraint data for all observed substitutions.
- **Supplemental\_Data\_10.Aligned\_reads.txt** Aligned reads to the transcripts coding for the base peptides (and corresponding SAAP) shoon in Supplemental Fig. 1 and 2.



Supplemental Fig. 1 | RNA-seq and fragment ion evidence for SAAP with RAAS>1 (part I). Mirror plots display the empirical and predicted fragment ion spectra for 1 label-free (a) and 4 TMT-labeled SAAPs (b-e) and their corresponding BP, with substitution sites bolded in the peptide sequences. Genome browser tracks display the DNA sequence corresponding to the detected RNA molecules and the corresponding amino acid sequence predicted by the genetic code. The aligned reads are available in Supplemental\_Data\_10.Aligned\_reads.txt.



Supplemental Fig. 2 | RNA-seq and fragment ion evidence for SAAP with high RAAS (part II). Additional examples of RNA data and mass spectra (TMT-labeled) supporting high RAAS substitutions. All notation is as in Supplemental Fig. 1



Supplemental Fig. 3 | Localizing mass shifts and distinguishing between alternative hypotheses for the mass shifts (a) Distribution of the positional probabilities, which quantifies the confidence with which a mass shift can be assigned to an amino acid residue, for all substituted peptides identified by the validation search of the CPTAC and healthy tissues data. **b-f** Mirror plots of predicted and empirical fragment ion spectra for TMT-labeled SAAP with  $Q \rightarrow G$  substitutions with uncertain localization of the mass shift. Mirror plots for an alternate hypothesis (alanine cleavage at the N-terminus of the peptide) that is consistent with the mass shift are also shown. The empirical spectra are better represented by the predicted spectra for the  $Q \rightarrow G$  substituted peptides than for the alternate hypothesis peptides. This conclusion is supported by cosine similarity scores and the presence of peaks (annotated with a star) that match the b2 fragment ion in the spectra predicted for the  $Q \rightarrow G$  substituted peptides. See Methods for details.



Supplemental Fig. 4 | Overlap of amino acid substitution with RNA modification sites (a) Modified uracil nucleotides (defined by nanopore sequencing<sup>34</sup>) overlap with codons of amino acid substitution sites. To evaluate the significance of the 180 overlapping sites, we calculated a p-value by a cumulative hypergeometric distribution test (see Methods for details) and here indicate this p-value (red x) in the context of the p-value distribution over different counts of overlaps with otherwise unchanged background numbers. (b) The measured fraction of modified U ( $\psi$  fraction corresponds to the column mm.DirectMINUSmm.IVT in the original data set) is not correlated to the median RAAS of these sites globally. (c) The  $\psi$  fraction at sites measured in the A549 cell line and overlapping with codon position 1 of amino acid substitution sites is positively correlated to the median RAAS at these sites. (d) The  $\psi$  fraction at sites measured in the NTERA cell line and overlapping 2 of amino acid substitution sites is negatively correlated to the median RAAS at these sites.



**Supplemental Fig. 5** | **SAAP detection and validation in metabolic pulse data** (a) SAAP detected in the CPTAC/label-free healthy data analysis that were also detected in the metabolic pulse data from Savtitski, *et al.*<sup>35</sup> are majorly derived from the analysis of healthy liver tissue. (b) Percentage of candidate SAAP from dependent peptide search of Savtitski, *et al.*<sup>35</sup> data that were validated by MSFragger.