

moPepGen: Rapid and Comprehensive Identification of Non-canonical Peptides

Chenghao Zhu^{a,b,c,d,t,#}, Lydia Y. Liu^{a,b,e,f,g,†}, Annie Ha^{e,f}, Takafumi N. Yamaguchi^{a,b,c}, Helen Zhu^{e,f,g}, Rupert Hugh-White^{a,b,c}, Julie Livingstone^{a,b,c}, Yash Patel^{a,b,c}, Thomas Kislinger^{e,f,#}, Paul C. Boutros^{a,b,c,d,e,#}

^a Department of Human Genetics, University of California, Los Angeles, CA, USA

^b Jonsson Comprehensive Cancer Center, University of California, Los Angeles, CA, USA

^c Institute for Precision Health, University of California, Los Angeles, CA, USA

^d Department of Urology, University of California, Los Angeles, CA, USA

^e Department of Medical Biophysics, University of Toronto, Toronto, Canada

^f Princess Margaret Cancer Centre, University Health Network, Toronto, Canada

^g Vector Institute for Artificial Intelligence, Toronto, Canada

[†] These authors contributed equally to this work

[#] Correspondence: chenghaozhu@mednet.ucla.edu, thomas.kislinger@utoronto.ca, pboutros@mednet.ucla.edu

Abstract

Gene expression is a multi-step transformation of biological information from its storage form (DNA) into functional forms (protein and some RNAs). Regulatory activities at each step of this transformation multiply a single gene into a myriad of proteoforms. Proteogenomics is the study of how genomic and transcriptomic variation creates this proteomic diversity, and is limited by the challenges of modeling the complexities of gene-expression. We therefore created moPepGen, a graph-based algorithm that comprehensively generates non-canonical peptides in linear time. moPepGen works with multiple technologies, in multiple species and on all types of genetic and transcriptomic data. In human cancer proteomes, it enumerates previously unobservable noncanonical peptides arising from germline and somatic genomic variants, noncoding open reading frames, RNA fusions and RNA circularization. By enabling efficient detection and quantitation of previously hidden proteins in both existing and new proteomic data, moPepGen facilitates all proteogenomics applications. It is available at:

<https://github.com/uclahs-cds/package-moPepGen>.

Main Text

A single stretch of DNA can give rise to multiple protein products through genetic variation and through transcriptional, post-transcriptional and post-translational processes (e.g., RNA editing, alternative splicing and RNA circularization)¹⁻⁴. The number of potential proteoforms rises combinatorially with the number of possibilities at each of these levels, so despite advances in proteomics technologies^{5,6}, much of the proteome is undetected in high-throughput studies⁷.

The most common strategies to detect peptide sequences absent from canonical reference databases⁷⁻⁹ (i.e., non-canonical peptides) are *de novo* sequencing and open search. Despite continued algorithmic improvements, these strategies are computationally expensive, have elevated false-negative rates, and lead to variant identification issues and difficult data interpretation^{10,11}. As a result, the vast majority of proteogenomic studies use non-canonical peptide databases that have incorporated DNA and RNA alterations⁷. These databases are often generated using DNA and RNA sequencing of the same sample, and this improves error rates relative to community-based databases (e.g., UniProt¹², neXtProt¹³ and the Protein Mutant Database¹⁴) by focusing the search space^{7,15}.

This type of sample-specific proteogenomics relies on the ability to predict all potential protein products generated by the complexity of gene expression. Modeling transcription, translation and peptide cleavage to fully enumerate the combinatorial diversity of peptides is computationally demanding. To simplify the search-space, existing methods have focused on generating peptides caused by individual variants or variant types¹⁶⁻³³, greatly increasing false negative rates and even potentially resulting in false positive detections if the correct peptide is absent from the database (**Extended Data Table 1**). To fill this gap, we created a graph-based algorithm for the exhaustive elucidation of protein sequence variations and subsequent *in silico* non-canonical peptide generation. This method is moPepGen (multi-omics Peptide Generator; **Figure 1a**).

moPepGen captures peptides that harbour any combination of small variants (e.g., single nucleotide polymorphisms [SNPs], small insertions and deletions [indels], RNA editing sites) occurring on canonical coding transcripts, as well as on non-canonical transcript backbones resulting from novel open reading frames (ORFs), transcript fusion, alternative splicing and

RNA circularization (**Supplementary Figure 1**). It performs variant integration, *in silico* translation and peptide cleavage in a series of three graphs for every transcript, enabling systematic traversal across every variant combination (**Online Methods; Extended Data Figure 1a-d**). All three reading frames are explicitly modeled for both canonical coding transcripts and non-canonical transcript backbones to efficiently capture frameshift variants and facilitate three-frame ORF search (**Extended Data Figure 2a**). Alternative splicing events (e.g., retained introns, etc.) and transcript fusions are modeled as subgraphs with additional small variants (**Extended Data Figure 2b**). Graphs are replicated four times to fully cover peptides of back-splicing junction read-through in circular RNAs (circRNAs; **Extended Data Figure 2c-d**). moPepGen outputs non-canonical peptides that cannot be produced by the chosen canonical proteome database. It documents all possible sources of each peptide to eliminate redundancy - for example where different combinations of genetic and transcriptomic events can produce the same peptide.

We first validated moPepGen using 1,000,000 iterations of fuzz testing (**Supplementary Figure 2**). For each iteration, a transcript model, its nucleotide sequence, and a set of variants composed of all supported variant types were simulated. Then non-canonical peptides generated by moPepGen were compared to those from a ground-truth brute-force algorithm. moPepGen demonstrated perfect accuracy and linear runtime complexity (4.7×10^{-3} seconds per variant) compared to exponential runtime complexity for the brute-force method (**Figure 1b-c**). A comprehensive non-canonical peptide database of human germline polymorphisms was generated with 15 GB of memory in 3.2 hours on a 16-core compute node; the brute-force method was unable to complete this task.

Having established the accuracy of moPepGen, we next compared it to two popular custom database generators, customProDBJ¹⁸ and pyQUILTS²². We tested all three methods on five prostate tumours with extensive multi-omics characterization³⁴⁻³⁶. We first evaluated the simple case of germline and somatic point mutations and indels. Most peptides ($84.0 \pm 0.9\%$ [median \pm MAD {median absolute deviation}]) were predicted by all three methods, with moPepGen being modestly more sensitive (**Extended Data Figure 3a**). Next, we considered the biological complexity of alternative splicing, RNA editing, RNA circularization and transcript fusion. Only moPepGen was able to evaluate peptides generated by all four of these processes, and therefore $80.2 \pm 2.1\%$ (median \pm MAD) of peptides were uniquely

predicted by moPepGen (**Extended Data Figure 3b**). By contrast only 3.2% of peptides were not predicted by moPepGen, and these corresponded to specific assumptions around the biology of transcription and translation made by other methods (**Extended Data Figure 3c; Online Methods**). By generating a more comprehensive database, moPepGen enabled the unique detection of $53.7 \pm 12.2\%$ (median \pm MAD) of peptide hits from matched proteomic data (**Extended Data Figure 3d**). An example of a complex variant peptide identified only by moPepGen is the combination of a germline in-frame deletion followed by a substitution in *SYNPO2* (**Figure 1d**). moPepGen's clear variant annotation system also readily enables verification across the central dogma. For example, the somatic mutation D1249N in *AHNAK* was detected in about 30% of both DNA and RNA reads and was detected by mass spectrometry (MS; **Figure 1e-i**), confirmed by three search engines. Taken together, these benchmarking results demonstrate the robust and comprehensive nature of moPepGen.

To illustrate the utility of moPepGen for proteogenomic studies, we first evaluated it across multiple proteases (**Extended Data Figure 4a**). Using independent conservative control of false discovery rate (FDR) across canonical and custom databases (**Online Methods; Supplementary Figure 3**)^{7,36}, we focused on detection of novel ORFs (*i.e.*, polypeptides from transcripts canonically annotated as noncoding) across seven proteases³⁷ in a deeply fractionated human tonsil sample³⁸ (**Supplementary Table 1**). moPepGen enabled the detection of peptides from 1,787 distinct ORFs previously thought to be noncoding, and these peptides were most easily detected with the Arg-C protease (**Extended Data Figure 4b**), suggesting alternative proteases may enhance noncoding ORF detection (**Extended Data Figure 4c**). In total 184 noncoding ORFs were detected by proteomics of four or more preparation methods in this single sample, demonstrating that moPepGen can identify novel proteins (**Extended Data Figure 4d-e**).

We next sought to demonstrate that moPepGen can benefit analyses in different species by studying germline variation in the C57BL/6N mouse^{39,40}. DNA sequencing of the related C57BL/6J strain was used to predict 5,481 non-canonical peptides arising from germline variants in protein-coding genes and 15,475 peptides from noncoding transcript novel ORFs (**Extended Data Figure 5a**). Across the proteomes of three bulk tissues (cerebellum, liver and uterus), we detected 18 non-canonical peptides in protein-coding genes and 343 from noncoding ORFs (**Extended Data Figure 5b-d; Supplementary Table 2**). Thus moPepGen

can support proteogenomics in non-human studies to identify variants of protein-coding genes and novel proteins.

To evaluate the utility of moPepGen for somatic variation, we analyzed 375 human cancer cell line proteomes with matched somatic mutations and transcript fusions^{41,42} (**Supplementary Data**). moPepGen processed each cell line in 2:58 minutes (median \pm 1:20 minutes, MAD), generating $2,683 \pm 2,513$ (median \pm MAD) potential non-canonical variant peptides per cell line. The number of predicted variant peptides varied strongly with tissue of origin, ranging from 838 - 16,255 (**Figure 2a**) and was driven largely by somatic mutations in protein-coding genes and by fusion events in noncoding genes (**Extended Data Figure 6a-c**). Searching the cell line proteomes identified 39 ± 27 (median \pm MAD) non-canonical peptides per cell line (**Online Methods; Supplementary Figure 4**). The majority of these were derived from noncoding transcript ORFs (**Extended Data Figure 6d; Supplementary Table 3**). Variant peptides from coding somatic mutations were more easily detected than those from transcript fusion events (**Extended Data Figure 6e-f**). 26 genes had variant peptides detected in cell lines from three or more tissues of origin, including the cancer driver genes *TP53*, *KRAS* and *HRAS* (**Figure 2b**). Peptide evidence was also found for fusion transcripts involving cancer driver genes like *MET* and *STK11* (**Extended Data Figure 6g-h**). We validated non-canonical peptide-spectrum matches (PSMs) by predicting tandem mass (MS2) spectra using Prosit⁴³ and verifying that variant peptide MS2 spectra correlated better with predictions based on the matched non-canonical peptide sequences than predictions based on their canonical peptide counterparts (**Online Methods; Extended Data Figure 6i**). Coding variant peptide PSMs also showed high cross-correlations with Prosit-predicted MS2 variant spectra, on par with those of canonical PSMs and their canonical spectra (**Extended Data Figure 6j**). Thus moPepGen can effectively and rapidly detect variant peptides arising from somatic variation. These variant peptides may prove to harbour functional consequences in future studies. Genes with non-canonical peptide hits, such as *KRAS*, trended towards greater essentiality for cell growth in multiple corresponding cell lines, and the effects may be independent of gene dosage (**Extended Data Figure 7a-c**). Across cell lines, detected variant peptides were also predicted to give rise to 416 putative neoantigens (3.0 ± 1.5 [median \pm MAD] per cell line; **Extended Data Figure 7d; Supplementary Table 4**), including recurrent neoantigens in *KRAS*, *TP53* and *FUPBP3* (**Extended Data Figure 7e**).

We next sought to demonstrate the utility of moPepGen in data-independent acquisition (DIA) MS using eight clear cell renal cell carcinoma tumours with matched whole-exome sequencing, RNA-sequencing and DIA proteomics⁴⁴. In each tumour, moPepGen predicted $157,016 \pm 34,215$ (median \pm MAD) unique variant peptides from protein-coding genes (**Extended Data Figure 8a**). Using a ProSight-generated spectral library, we detected 307 ± 112 (median \pm MAD) variant peptides in each tumour using DIA-NN⁴⁵ (**Extended Data Figure 8b**; **Supplementary Table 5**). Germline-SNP and alternative splicing were the most common sources of detected variant peptides (**Extended Data Figure 8c-d**). Non-canonical peptides derived from RNA editing events were detected in 21 genes (**Extended Data Figure 8e-i**). Thus moPepGen can enable detection of variant peptides from DIA proteomics.

Finally, to demonstrate the utility of moPepGen on the most complex gene-expression data, we analyzed five primary prostate cancer samples with matched DNA whole-genome sequencing, ultra-deep ribosomal-RNA-depleted RNA-sequencing and mass-spectrometry-based proteomics³⁴⁻³⁶. moPepGen generated $1,382,666 \pm 64,281$ (median \pm MAD) unique variant peptides per sample, spanning 115 variant combination categories (**Figure 2c**). Searching this database resulted in the detection of 206 ± 56 (median \pm MAD) non-canonical peptides per sample. The distribution of intensities and Comet expectation scores of non-canonical target peptide hits closely resembles that of canonical target peptides and is distinct from all decoy hits (**Supplementary Figure 5**), lending confidence in our non-canonical peptide detection. 138 ± 28 (median \pm MAD) detected non-canonical peptides were derived from protein-coding genes (**Extended Data Figure 9a**; **Supplementary Table 6**). All samples harboured proteins containing multiple variant peptides (9 ± 1.5 [median \pm MAD] proteins per tumour; range of 2-6 variant peptides per protein; **Figure 2d**). Some detected peptides harboured multiple variants, including two from prostate-specific antigen (PSA from the *KLK3* gene; **Extended Data Figure 9b**). Germline SNPs were the major common cause of variant peptides on coding transcripts and alternative splicing events were the most common cause on noncoding transcripts (**Extended Data Figure 9c-e**). Nine genes showed recurrent detection of peptides caused by circRNA back-splicing (**Extended Data Figure 9f-g**), with 36/78 circRNA PSMs validated by *de novo* sequencing (**Supplementary Table 7**)⁴⁶. These recurrent circRNA-derived peptides were verified in five additional prostate tumours (**Supplementary Figure 6**). We also detected four peptides from noncoding transcripts with

the recently reported tryptophan-to-phenylalanine substituents⁴⁷. Thus moPepGen can identify peptides resulting from highly complex layers of gene-expression regulation.

moPepGen is a computationally efficient algorithm that enumerates transcriptome and proteome diversity across arbitrary variant types. It enables the detection of variant and novel ORF peptides across species, proteases and technologies. moPepGen integrates into existing proteomic analysis workflows, and can broadly enhance proteogenomic analyses for many applications.

References

1. Zhang, B. *et al.* Proteogenomic characterization of human colon and rectal cancer. *Nature* 513, 382–387 (2014).
2. Sinitcyn, P. *et al.* Global detection of human variants and isoforms by deep proteome sequencing. *Nat Biotechnol* (2023) doi:10.1038/S41587-023-01714-X.
3. Nilsen, T. W. & Graveley, B. R. Expansion of the eukaryotic proteome by alternative splicing. *Nature* 463, 457–463 (2010).
4. Peng, X. *et al.* A-to-I RNA Editing Contributes to Proteomic Diversity in Cancer. *Cancer Cell* 33, 817-828.e7 (2018).
5. Creighton, C. J. Clinical proteomics towards multiomics in cancer. *Mass Spectrom Rev* (2022) doi:10.1002/MAS.21827.
6. Edwards, N. J. *et al.* The CPTAC data portal: A resource for cancer proteomics research. *J Proteome Res* 14, 2707–2713 (2015).
7. Nesvizhskii, A. I. Proteogenomics: concepts, applications and computational strategies. *Nature Methods* 2014 11:11 11, 1114–1125 (2014).
8. Chick, J. M. *et al.* A mass-tolerant database search identifies a large proportion of unassigned spectra in shotgun proteomics as modified peptides. *Nat Biotechnol* 33, 743–749 (2015).
9. Rodriguez, H., Zenklusen, J. C., Staudt, L. M., Doroshow, J. H. & Lowy, D. R. The next horizon in precision oncology: Proteogenomics to inform cancer diagnosis and treatment. *Cell* 184, 1661–1670 (2021).
10. Ma, B. & Johnson, R. De novo sequencing and homology searching. *Mol Cell Proteomics* 11, (2012).
11. Fu, Y. Data analysis strategies for protein modification identification. *Methods in Molecular Biology* 1362, 265–275 (2016).
12. Bateman, A. UniProt: A worldwide hub of protein knowledge. *Nucleic Acids Res* 47, D506–D515 (2019).
13. Lane, L. *et al.* neXtProt: a knowledge platform for human proteins. *Nucleic Acids Res* 40, (2012).
14. Kawabata, T., Ota, M. & Nishikawa, K. The Protein Mutant Database. *Nucleic Acids Res* 27, 355–357 (1999).

15. Salz, R. *et al.* Personalized Proteome: Comparing Proteogenomics and Open Variant Search Approaches for Single Amino Acid Variant Detection. *J Proteome Res* 20, 3353–3364 (2021).
16. Wang, X. *et al.* Protein identification using customized protein sequence databases derived from RNA-seq data. *J Proteome Res* 11, 1009–1017 (2012).
17. Wang, X., Zhang, B. & Wren, J. customProDB: an R package to generate customized protein databases from RNA-Seq data for proteomics search. *Bioinformatics* 29, 3235–3237 (2013).
18. Wen, B., Li, K., Zhang, Y. & Zhang, B. Cancer neoantigen prioritization through sensitive and reliable proteogenomics analysis. *Nature Communications* 2020 11:1 11, 1–14 (2020).
19. Sinitcyn, P., Gerwien, M. & Cox, J. MaxQuant Module for the Identification of Genomic Variants Propagated into Peptides. *Methods in Molecular Biology* 2456, 339–347 (2022).
20. Van De Geer, W. S., Van Riet, J. & Van De Werken, H. J. G. ProteoDisco: a flexible R approach to generate customized protein databases for extended search space of novel and variant proteins in proteogenomic studies. *Bioinformatics* 38, 1437–1439 (2022).
21. Umer, H. M. *et al.* Generation of ENSEMBL-based proteogenomics databases boosts the identification of non-canonical peptides. *Bioinformatics* 38, 1470–1472 (2022).
22. Ruggles, K. V. *et al.* An Analysis of the Sensitivity of Proteogenomic Mapping of Somatic Mutations and Novel Splicing Events in Cancer. *Molecular & Cellular Proteomics* 15, 1060–1071 (2016).
23. Sheynkman, G. M., Shortreed, M. R., Frey, B. L. & Smith, L. M. Discovery and Mass Spectrometric Analysis of Novel Splice-junction Peptides Using RNA-Seq. *Molecular & Cellular Proteomics* 12, 2341–2353 (2013).
24. Sheynkman, G. M. *et al.* Using Galaxy-P to leverage RNA-Seq for the discovery of novel protein variations. *BMC Genomics* 15, 1–9 (2014).
25. Sheynkman, G. M., Shortreed, M. R., Frey, B. L., Scalf, M. & Smith, L. M. Large-scale mass spectrometric detection of variant peptides resulting from nonsynonymous nucleotide differences. *J Proteome Res* 13, 228–240 (2014).
26. Wen, B. *et al.* sapFinder: an R/Bioconductor package for detection of variant peptides in shotgun proteomics experiments. *Bioinformatics* 30, 3136–3138 (2014).
27. Wen, B. *et al.* PGA: An R/Bioconductor package for identification of novel peptides using a customized database derived from RNA-Seq. *BMC Bioinformatics* 17, 1–6 (2016).

28. Woo, S. *et al.* Proteogenomic database construction driven from large scale RNA-seq data. *J Proteome Res* 13, 21–28 (2014).
29. Woo, S. *et al.* Proteogenomic strategies for identification of aberrant cancer peptides using large-scale next-generation sequencing data. *Proteomics* 14, 2719–2730 (2014).
30. Cesnik, A. J. *et al.* Spritz: A Proteogenomic Database Engine. *J Proteome Res* 20, 1826–1834 (2021).
31. Cifani, P. *et al.* ProteomeGenerator: A Framework for Comprehensive Proteomics Based on de Novo Transcriptome Assembly and High-Accuracy Peptide Mass Spectral Matching. *J Proteome Res* 17, 3681–3692 (2018).
32. Huber, F. *et al.* A comprehensive proteogenomic pipeline for neoantigen discovery to advance personalized cancer immunotherapy. *Nat Biotechnol* (2024) doi:10.1038/s41587-024-02420-y.
33. Kwok, N. *et al.* Integrative Proteogenomics Using ProteomeGenerator2. *J Proteome Res* 22, 2750–2764 (2023).
34. Fraser, M. *et al.* Genomic hallmarks of localized, non-indolent prostate cancer. *Nature* 541, 359–364 (2017).
35. Chen, S. *et al.* Widespread and Functional RNA Circularization in Localized Prostate Cancer. *Cell* 176, 831-843.e22 (2019).
36. Sinha, A. *et al.* The Proteogenomic Landscape of Curable Prostate Cancer. *Cancer Cell* 35, 414-427.e6 (2019).
37. Giansanti, P., Tsiatsiani, L., Low, T. Y. & Heck, A. J. R. Six alternative proteases for mass spectrometry–based proteomics beyond trypsin. *Nature Protocols* 2016 11:5 11, 993–1006 (2016).
38. Wang, D. *et al.* A deep proteome and transcriptome abundance atlas of 29 healthy human tissues. *Mol Syst Biol* 15, e8503 (2019).
39. Keane, T. M. *et al.* Mouse genomic variation and its effect on phenotypes and gene regulation. *Nature* 477, 289–294 (2011).
40. Giansanti, P. *et al.* Mass spectrometry-based draft of the mouse proteome. *Nature Methods* 2022 19:7 19, 803–811 (2022).
41. Ghandi, M. *et al.* Next-generation characterization of the Cancer Cell Line Encyclopedia. *Nature* 2019 569:7757 569, 503–508 (2019).
42. Nusinow, D. P. *et al.* Quantitative Proteomics of the Cancer Cell Line Encyclopedia. *Cell* 180, 387-402.e16 (2020).

43. Gessulat, S. *et al.* Prosit: proteome-wide prediction of peptide tandem mass spectra by deep learning. *Nature Methods* 2019 16:6 16, 509–518 (2019).
44. Li, Y. *et al.* Histopathologic and proteogenomic heterogeneity reveals features of clear cell renal cell carcinoma aggressiveness. *Cancer Cell* 41, 139-163.e17 (2023).
45. Demichev, V., Messner, C. B., Vernardis, S. I., Lilley, K. S. & Ralser, M. DIA-NN: neural networks and interference correction enable deep proteome coverage in high throughput. *Nat Methods* 17, 41–44 (2020).
46. Ma, B. Novor: real-time peptide de novo sequencing software. *J Am Soc Mass Spectrom* 26, 1885–1894 (2015).
47. Pataskar, A. *et al.* Tryptophan depletion results in tryptophan-to-phenylalanine substituents. *Nature* 2022 603:7902 603, 721–727 (2022).

Online Methods

Transcript Variant Graph

A transcript variant graph (TVG) is instantiated for each transcript, incorporating all associated variants. In a TVG, nodes are transcript fragments with reference or alternative nucleotide sequences, while edges are the opening or closing of variant nodes connecting them to the reference sequence, or simply the elongation of reference sequences. The TVG starts with three linear nodes of the entire transcript sequence representing the three reading frames, with x number of nucleotides removed from the transcript N-terminus, where x equals 0, 1, or 2. A variant is incorporated into the graph by cutting nodes in the graph at the variant's start and end positions and attaching a new node with the alternative sequence to the new upstream and downstream nodes. An in-frame variant is represented as a node that has incoming and outgoing nodes in the same reading frame subgraph, while frameshifting variants have incoming nodes and outgoing nodes in different subgraphs representing different reading frames. The outgoing reading frame index equals to $(S_{ref} - S_{alt}) \bmod 3$, where S_{ref} is the length of the reference sequence and S_{alt} is the length of the alternative sequence. For transcripts with an annotated known canonical open reading frame (ORF), variants are only incorporated into the subgraph with the appropriate reading frame index (**Extended Data Figure 2a**). If frameshifting variants are present, variants are also incorporated into the subgraphs with the reading frame indices of the outgoing frameshift nodes. For transcripts without an annotated ORF, all variants are incorporated into all three subgraphs (**Extended Data Figure 2a**). Large insertions and substitutions as the result of alternative splicing events (e.g., retained introns, alternative 3'/5' splicing, etc.) are represented as subgraphs that can carry additional variants (**Extended Data Figure 2b**).

Variant Bubbles and Peptide Variant Graph

After the TVG has been populated with all variants, nodes that overlap with each other in transcriptional coordinates are aligned to create variant bubbles within which all nodes point to the same upstream and downstream nodes (**Extended Data Figure 1b**). This is done by first finding connection nodes in the TVG. A connection node is a reference node without any variants in its transcriptional coordinates that connects two variant bubbles after they are aligned. The root node is the first connection node, and the next connection node is found by

looking for the first commonly connected downstream node with length of five or more nucleotides that is outbound to more than one node (**Supplementary Note 1; Supplementary Figure 7**). Nodes between the two connection nodes are then aligned and a variant bubble is formed by generating all combinations of merged nodes so that they all point to the same upstream and downstream nodes (**Extended Data Figure 1b**). Overlapping variants in the variant bubble are eliminated because they are disjoint. The sequence lengths of nodes in the variant bubble are also adjusted by taking nucleotides from the commonly connected upstream and downstream nodes to ensure that they are multiples of three. A peptide variant graph (PVG) is then instantiated by translating the nucleotide sequence of each TVG node into amino acid sequences.

Peptide Cleavage Graph

The peptide variant graph is converted into the peptide cleavage graph (PCG), where each edge represents an enzymatic cleavage site (**Extended Data Figure 1c**). For connection nodes, all enzymatic cleavage sites are first identified, and the node is cleaved at each cleavage site. Because enzymatic cleavage sites can span over multiple nodes (for example, the trypsin exception of not cutting given K/P but cutting given WK/P), connection nodes are also merged with all downstream and/or upstream nodes and cut at additional cleavage sites if found. To optimize run time, different *merge-and-cleave* operations are used depending on the number of incoming and outgoing nodes, and the number of cleavage sites in a node (**Supplementary Figure 8**). Hypermutated regions where variant bubbles contain many variants and/or the lack of cleavage sites in connection nodes can result in an exponential increase in the number of nodes in the aligned variant bubble. To deal with this, we used a *pop-and-collapse* strategy, such that when *merge-and-cleave* is applied to a connection node, x number of amino acids are popped from the end of each node in the variant bubble. The popped nodes are collapsed if they share the same sequence. The pop-and-collapse operation is only applied when the number of nodes in a variant bubble exceed a user-defined cutoff.

Calling Variant Peptides

Variant peptides with the permitted number of miscleavages are called by traversing through the peptide cleavage graph. A node in a PCG can be in a different ORF when arriving at the node by traversing from a different incoming node. Thus, we use a *stage-and-call* approach

that first visits all incoming nodes to determine the valid ORFs of a peptide node (**Supplementary Note 2**). *Stage-and-call* also allows cleavage-gain mutations and upstream frameshift mutations to be carried over to the downstream peptide nodes. Peptide nodes are then extended by merging with downstream nodes to call variant peptides with miscleavages (**Supplementary Note 3**). For noncoding transcripts, novel ORF start sites, including those caused by start-gain mutations, are found by looking for any methionine (M) in all three subgraphs.

Fusion and Circular Transcripts

Most fusion transcript callers detect the fusion events between genes, thus causing ambiguity as to which transcripts of the genes are involved in a particular fusion event. We took the most comprehensive approach and endeavored to capture all possible variant peptides by assuming that the fusion event could happen between any transcript of the donor and acceptor genes. Fusion transcripts are considered as novel backbones in graph instantiation, with an individual graph instantiated for each donor and acceptor transcript pair. Single nucleotide variants (SNV) and small insertion/deletions (indels) of both donor and acceptor transcripts are incorporated into the TVG. The translated and cleaved PCG is then traversed to call variant peptides, identical to a canonical transcript backbone. If the fusion breakpoint is in an intronic region, the intronic nucleotide sequence leading up to or following the breakpoint is retained as unspliced, and any SNVs or indels that it carries are also incorporated into the graph. The ORF start site of the donor transcript, if exists, is used when calling variant peptides. The fusion transcript is treated as a noncoding transcript if the donor transcript is annotated as noncoding.

Similar to fusion transcripts, circular RNAs (circRNAs) are treated as novel backbones, with an individual graph instantiated for each circRNA (**Extended Data Figure 2c**). A circular variant graph (CVG, a counterpart to TVG) is instantiated by connecting the linear sequence of the circRNA onto itself at the back-splice junction and incorporating SNVs and indels. Novel peptides can theoretically be translated from circRNAs if a start codon is present, by ribosome readthrough across the back-splicing junction site. If the circRNA length is not a multiple of three nucleotides, translation across the back-splicing site induces a frameshift. Without a stop codon, the ribosome may traverse the circRNA up to three times before the amino acid sequence repeats. Therefore, moPepGen extends the circular graph linearly by

appending three copies of each reading frame as a subgraph to account for frameshifts. The extended graph is then translated to a PVG and converted to a PCG. Variant peptides are called by treating every circRNA as a noncoding transcript and exploring all novel start codons in all three reading frames.

Biological Assumptions for Edge Cases

moPepGen applies various assumptions to selectively include or exclude certain variant events or peptides (**Extended Data Figure 3c**). Start-codon-altering variants are excluded due to the uncertainty around whether and where translation will still occur. Similarly, splice-site-altering variants are omitted due to the complexity of splicing determinants, which can result in skipping to the next canonical or non-canonical splice site. We terminate translation at the last complete peptide when stop codons are unknown, as incomplete transcript annotations create ambiguity in downstream sequences, obscuring enzymatic cleavage sites. Stop-codon-altering variants do not extend translation beyond the transcript, as the downstream genomic region is not assumed to be part of the RNA transcript.

GVF File Format and Parsers

Genomic (single nucleotide polymorphisms [SNP], SNV, indel) and transcriptomic variants (fusion transcripts, RNA editing sites, alternative splicing transcripts, circRNAs) are first converted into gene-centric entries for each transcript that they impact. We defined the GVF (genetic variant format) file format, derived from the VCF (variant calling format) file format, to store all relevant information for each variant. The location of each variant is represented by gene ID along with offset from the start of the gene. Parsers have been implemented as part of moPepGen to parse the outputs of a variety of variant calling tools and convert them into the GVF file format. VCF files containing SNPs, SNVs and indels require annotation via Variant Effect Predictor (VEP) for compatibility with the parseVEP module. Native output files for fusion transcripts, alternative splicing events, RNA editing sites and circRNAs are directly processed by their respective moPepGen parsers. moPepGen was implemented in Python and supports easy extension and addition of new parsers. The corresponding Nextflow pipeline, available at <https://github.com/uclahs-cds/pipeline-call-NonCanonicalPeptide>, orchestrates the input data processing, non-canonical peptide prediction, database stratification into tiers, and other ancillary functions, including filtering based on transcript abundance^{48,49}. Additionally, our DNA data processing pipeline, designed for calling germline

and somatic SNVs and indels from whole-genome sequencing (WGS) and whole-exome sequencing (WXS) data with recalibration and adjustment, is accessible at <https://github.com/uclahs-cds/metapipeline-DNA>⁵⁰.

Fuzz Testing and Brute Force Algorithm

A simulation-based fuzz testing framework is used to ensure the validity of moPepGen. In every fuzz test case, a transcript is simulated with a transcript model of various parameters (e.g., coding or noncoding, positive or negative strand, selenoprotein or not, unknown start or stop codon position) and artificial sequence. Variant records associated with the transcript are then simulated, covering all supported variant types, including SNV, indel, fusion, alternative splicing, RNA editing and circRNA. Variant peptides produced by moPepGen are then compared to the output of the brute force algorithm, which iterates through all possible variant combinations to identify non-canonical peptides. The brute force algorithm also performs three-frame translation for noncoding transcripts. Software encoding this algorithm is available in the same repository as moPepGen.

Datasets

Cancer Cell Line Encyclopedia proteome

Proteomics characterization of 375 cell lines from the Cancer Cell Line Encyclopedia (CCLE) was obtained from Nusinow *et al.*, 2020⁴². Fractionated raw mass spectrometry (MS) data were downloaded from MassIVE (project ID: MSV000085836). Somatic SNVs and indels, and fusion transcript calls from the CCLE project were downloaded from the DepMap portal (<https://depmap.org/portal>, 22Q1). Somatic mutations were converted to GRCh38 coordinates from hg19 using CrossMap (v0.5.2)⁵¹. Gene and transcript IDs were assigned to each SNV/indel using the Variant Effect Predictor (VEP)⁵² (v104) with genomic annotation GTF downloaded from GENCODE (v34)⁵³. Fusion results were aligned to the GENCODE v34 reference by first lifting over the fusion coordinates to GRCh38 using CrossMap (v0.5.2). After lift-over, the records were removed if the donor or acceptor breakpoint location was no longer associated with the gene, if either breakpoint dinucleotides did not match with the reference, or if either gene ID was not present in GENCODE (v34).

Mouse proteome

Mass spectrometry-based characterization of the proteome of mouse strain C57BL/6N was obtained from Giansanti *et al.*, 2022⁴⁰. Fractionated raw mass spectrometry data of the liver, uterus and cerebellum proteomes was downloaded from the PRIDE repository (project ID: PXD030983). Germline single nucleotide polymorphisms (SNPs) and indels were obtained from the Mouse Genomes Project³⁹ with GRCm38 VCFs downloaded from the European Variation Archive (accession: PRJEB43298). Germline SNPs and indels were mapped using VEP (v102) to Ensembl GRCm38 GTF (v102)⁵⁴.

Alternative protease and fragmentation proteome

A human tonsil tissue processed using ten different combinations of proteases and peptide fragmentation methods (ArgC_HCD, AspN_HCD, Chymotrypsin_CID, Chymotrypsin_HCD, GluC_HCD, LysC_HCD, LysN_HCD, Trypsin_CID, Trypsin_ETD, Trypsin_HCD) was obtained from Wang *et al.*, 2019³⁸. Fractionated raw mass spectrometry data were downloaded from the PRIDE repository (project ID: PXD010154).

DIA Proteome

Data-independent acquisition (DIA) proteomic data from eight clear cell renal cell carcinoma samples were obtained from Li *et al.*, 2023⁴⁴. Raw mass spectrometry data were retrieved from the Proteomic Data Commons (PDC) under accession number PDC000411. WXS and RNA-seq BAM files were obtained from Genomic Data Commons (GDC, Project: CPTAC-3, Primary Site: Kidney). WXS data was processed using a standardized pipeline to identify germline SNPs, somatic SNVs and indels⁵⁰. BAM files were first reverted to FASTQ format using Picard toolkit (v2.27.4) and SAMtools (v1.15.1)⁵⁵ and subsequently re-aligned to the human reference genome GRCh38 using BWA-MEM2 (v2.2.1)⁵⁶. The re-aligned BAM files were then calibrated using BQSR and IndelRealignment from GATK (v4.2.4.1)⁵⁷. Germline SNPs and indels were called following GATK (v4.2.4.1) best practices^{57,58}, while somatic SNVs and indels were called using Mutect2 (from GATK v4.5.0.0). Gene and transcript IDs from GENCODE (v34) were assigned to variants by VEP (v104)⁵⁹. Similarly, RNA-seq BAM files were converted back to FASTQ format using Picard toolkit (v2.27.4) and SAMtools (v1.15.1) and re-aligned to human genome GRCh38.p13 with GENCODE v34 GTF using STAR (2.7.10b)⁶⁰. Transcript fusion events were called using STAR-Fusion (v1.9.1)⁶¹,

alternative splicing events were called using rMATS (v4.1.1)⁶², and RNA editing sites were called using REDIttools2 (v1.0.0)⁶³ using paired RNA and DNA BAMs.

Prostate cancer proteome

The proteomics characterization of five prostate cancer tissues were obtained from Sinha *et al.*, 2019³⁶. Raw mass spectrometry data were downloaded from MassIVE (project ID: MSV000081552). Germline SNPs and indels, as well as somatic SNVs and indels were obtained from the ICGC Data Portal (Project code: PRAD-CA). Variants were indexed using VCFtools (v0.1.16)⁶⁴ and lifted over to GRCh38 using Picard toolkit (v2.19.0), followed by chromosome name mapping from the Ensembl to the GENCODE system using BCFtools (v1.9-1)⁶⁵. Gene and transcript IDs were mapped by VEP (v104)⁵⁹ to the GENCODE (v34) GTF. Raw mRNA sequencing data were obtained from Gene Expression Omnibus (accession: GSE84043). Transcriptome alignment was performed using STAR (v2.7.2) to reference genome GRCh38.p13 with GENCODE (v34) GTF and junctions were identified by setting the parameter `--chimSegmentMin 10`⁶⁶. CIRCexplorer2 (v.2.3.8) was used to parse and annotate junctions for circular RNA detection⁶⁷. Fusion transcripts were called using STAR-Fusion (v1.9.1)⁶¹. RNA editing sites were called using REDIttools2 using paired RNA and DNA BAMs (v1.0.0)⁶³. Alternative splicing transcripts were called using rMATS (v4.1.1)⁶².

Terminology

Non-canonical Peptides

Non-canonical peptides are absent from the canonical reference protein sequence database. They can arise from genomic variants (*e.g.*, germline SNPs, somatic SNVs, indels), transcriptomic events (*e.g.*, alternative splicing, transcript fusion, RNA editing, circRNAs), and sequence modifications during translation (*e.g.*, selenocysteine terminations, tryptophan-to-phenylalanine [W>F] substitutants) in protein-coding transcripts. Peptides translated from novel ORFs in transcripts annotated as noncoding (*e.g.*, lncRNAs, pseudogenes) are also classified as non-canonical peptides when derived from valid biological processes, with possible further additions of genomic or transcriptomic variants.

Variant Peptides

Variant peptides are a subset of non-canonical peptides harbouring genomic and/or transcriptomic variants translated from protein-coding transcripts.

Proteoform

A proteoform is a specific molecular form of a protein with a unique amino acid sequence resulting from genetic variations (e.g., germline SNPs, somatic SNVs, indels), transcriptional processes (e.g., alternative splicing, transcript fusion, RNA editing), translation-level sequence modifications (e.g., selenocysteine terminations, W>F substituents) or post-translational/co-translational modifications generated by enzymatic (e.g., phosphorylation, acetylation, N/O-glycosylation) or non-enzymatic (e.g., glycation, oxidation, etc.) attachment of specific chemical moieties to the side chains of specific amino acids.

Canonical Database

A canonical database is a reference protein database containing sequences derived exclusively from well-annotated, protein-coding genes and transcripts without variants, modifications, or alternative splicing events. It represents the standard, widely accepted protein sequences in a given organism, excluding non-canonical or variant sequences that arise from genomic or transcriptomic alterations.

Non-canonical Database

A non-canonical database is a protein database that contains only sequences absent from the canonical database. It includes peptides or proteins derived from genomic or transcriptomic variants, alternative splicing events, or previously unannotated and misannotated regions, ensuring no overlap with standard, well-annotated protein sequences. This database enables the detection of unique, rare, or condition-specific peptides that cannot be identified using canonical databases alone.

Canonical Database Search

All mass spectrometry raw files (.raw) were converted to the open format mzML using ProteoWizard (3.0.21258)⁶⁸. The human GRCh38 canonical proteome database was obtained from GENCODE (v34), concatenated with common contaminants⁶⁹, and appended with reversed sequences to enable target-decoy false discovery rate (FDR) control. Mouse GRCm38 canonical proteome database was obtained from Ensembl (v102) and similarly processed. Database search was performed using Comet (v2019.01r5)⁷⁰. All searches were performed with static modifications of cysteine carbamidomethylation, and up to three variable modifications of methionine oxidation, protein N-terminus acetylation and peptide N-

terminus pyroglutamate formation. All searches were performed with fully specific digestion on both peptide ends for peptide lengths of 7-35 amino acids. Except for the tonsil samples processed with alternative enzymes, searches were performed with trypsin digestion and up to two miscleavages. For CCLE, static modification of tandem mass tag (TMT; 10plex) on the peptide N-terminus and lysine residues and variable modification of TMT on serine residues were additionally included, in accordance with the original study. For CCLE, searches were set to low resolution with 20 ppm precursor mass tolerance, 0.5025 Da fragment mass tolerance, and clear m/z range corresponding to TMT10plex, in accordance with the original publication. All other searches were of high-resolution label-free quantification (LFQ), with precursor mass tolerance of 20 ppm, 10 ppm, 30 ppm for the mouse proteome, tonsil proteome and prostate cancer proteome, respectively, and fragment mass tolerances of 0.025 Da for the tonsil proteome and 0.01 Da otherwise, in accordance with original publications. Tonsil proteomes were searched with the appropriate protease used in sample preparation, with maximum two miscleavages for Lys-C and Arg-C, three miscleavages for Glu-C and Asp-N and four miscleavages for chymotrypsin, as in the original publication³⁸.

Peptide level target-decoy FDR calculation was performed using the FalseDiscoveryRate module from OpenMS (v3.0.0-1f903c0)⁷¹ using the formula $(D+1)/(T+D)$, where D is the number of decoy peptide-spectrum matches (PSMs) and T is the number of target PSMs. Peptides were filtered at 1% FDR, and all PSMs were removed from the corresponding mzML for subsequent non-canonical database search. *Post-hoc* cohort level FDR was calculated to verify an FDR cutoff smaller than 1%. Peptide quantification was performed using FeatureFinderIdentification as part of OpenMS (v3.0.0-1f903c0)⁷², using the “internal IDs only” strategy and with adjusted precursor mass tolerances as above, and otherwise default parameters for LFQ proteomics. The IsobaricAnalyzer module from OpenMS (v3.0.0-1f903c0) was used for the quantification of channel intensities for TMT proteomics, with no isotope correction due to the lack of correction matrix.

Non-canonical Database Generation

Human GRCh38 proteome reference files were obtained from GENCODE (v34) while mouse GRCm38 proteome reference files were obtained from Ensembl (v102). All non-canonical peptide databases were generated with trypsin digestion of up to two miscleavages and peptide lengths 7-25, except for databases used for alternative protease samples. Peptides

from alternative translation were generated using *callAltTranslation*, including those with selenocysteine termination⁷³ or W>F substituents⁴⁷. Peptides from noncoding ORFs were generated using *callNovelORF* with ORF order as min and with or without alternative translation. Noncoding ORF peptide databases were also generated for each alternative protease used in processing of the tonsil proteome, with appropriate number of maximum miscleavages as outlined above.

Non-canonical peptide databases were generated for 376 cell lines from CCLE, 375 of which have non-reference channel proteomics characterization. This included all 10 cell lines in the bridge line and 366 non-reference cell lines with mutation data. Of the 378 non-reference channels across 42 plexes, three cell lines were duplicated, seven were in the bridge line, two didn't have mutation or fusion information and additional eight didn't have fusion information. Variant databases from all cell lines in a TMT plex, including the ten cell lines in the reference channel, were merged along with noncoding ORF peptides to generate plex-level databases. Plex-level databases were split into "Coding" (coding point mutations and transcript fusions), "Noncoding" (noncoding transcript novel ORFs) and "Noncoding Variant" (point mutations and transcript fusions on noncoding transcript ORFs) databases for tiered database search.

Non-canonical peptide databases for the proteome of mouse strain C57BL/6N was generated by calling variant peptides based on germline SNPs and indels, followed by merging with the noncoding ORF peptides database, and splitting into the three databases. The "Germline" database contained coding germline variations, "Noncoding" the novel ORFs, and "Noncoding-Germline" the peptides from germline variations on noncoding transcript ORFs. Genomic and transcriptomic variants (*i.e.*, germline SNPs, germline indels, somatic SNVs, somatic indels, RNA editing sites, transcript fusions, alternative splicing and circRNA) from five prostate tumour samples were similarly used to call variant peptides, which were merged with the noncoding ORF peptide database and alternative translation peptide database and split into five databases. The "Variant" database included non-canonical peptides from coding transcripts with SNVs, indels, RNA editing, alternative splicing or transcript fusions. The "Noncoding" database included all peptides from noncoding transcript novel ORFs and noncoding peptides with any variants were in the "Noncoding Variant" database. The "Circular RNA" database included all peptides representing circRNA open reading frames

(ORFs) with or without other sequence changes. The “Alt Translation” database included any peptides with selenocysteine insertion or W>F substituents⁴⁷.

Non-canonical Database Search

Non-canonical database searches were performed in nearly identical fashion as canonical proteome searches for each dataset, as described in detail above. Custom databases of peptide sequences were concatenated with the reverse sequence for FDR control. Non-canonical peptide searches with Comet (v2019.01r5) were set to “no cleavage” and did not permit protein N-terminus modifications or clipping of N-terminus methionine. Peptide-level FDR was set to 1% independently for each tier of non-canonical database, and PSMs of peptides that passed FDR were removed from the mzML for subsequent searches. *Post-hoc* cohort level FDR was calculated to verify an FDR cutoff smaller than 1%. Each tier of database thus had independent FDR control using database-specific decoy peptides, and a spectra is excluded from subsequent searches after finding its most probable match. This strategy has been shown to both alleviate the detection of false-positive non-canonical peptides due to joint FDR calculation with canonical peptides and to enable a highly conservative approach to non-canonical peptide detection^{7,74}. For CCLE specifically, peptide detection and quantification were only considered for a cell line when the non-canonical peptide exists in the sample-specific database. For prostate tumors, additional searches were conducted with the same non-canonical databases using MSFragger (v3.3)⁷⁵ and X!Tandem (v2015.12.15)⁷⁶ with equivalent parameters for verification. For all datasets, quantified peptides were distinguished by charge and variable modifications, and detected but not quantified peptides were excluded from subsequent analysis.

DIA Non-canonical Spectral Library Search

Raw files were converted to .mzML files using ProteoWizard (3.0.21258)⁶⁸. Sample-specific variant peptide FASTA databases were generated using the aforementioned non-canonical database generation pipeline, with individual spectral libraries .msp files generated by ProSIT⁷⁷. ProSIT was configured with the instrument type of LUMAS, collision energy of 34, and fragmentation method of HCD, with all other parameters set to default. Each sample was searched against the sample-specific predicted variant peptide spectral library using DIA-NN (v.1.8.1)⁴⁵. Spectral library searches were performed using the sample-specific spectral

library, with protein inference disabled. All searches were performed with q-value cutoff of 0.01 and quantified using the “high precision” mode.

Neoantigen Prediction

Neoantigens were predicted from non-canonical peptide hits detected in the CCLE proteomes. First, cell line-specific *HLA* genotype was predicted using OptiType (v1.3.5)⁷⁸ from whole-genome or whole-exome sequencing data. Non-canonical peptide hits in the “Coding” database tier were converted to FASTA format using a custom script. Neoantigen predictions were subsequently performed MHCflurry (v2.0.6)⁷⁹, with default parameters and cell line-specific *HLA* genotypes.

Statistical Analysis and Data Visualization

All statistical analysis and data visualization were performed in the R statistical environment (v4.0.3), with visualization using BoutrosLab.plotting.general (v6.0.2)⁸⁰. All boxplots show all data points, the median (center line), upper and lower quartiles (box limits), and whiskers extend to the minimum and maximum values within 1.5 times the interquartile range. Schematics were created in Inkscape (v1.0) and Adobe Illustrator (27.8.1), and figures were assembled using Inkscape (v1.0).

Gene Dependency Association Analysis

Gene dependency data from the CRISPR screens in the CCLE project were downloaded from the DepMap data portal (<https://depmap.org/portal>, 24Q2). Twelve cell lines with non-canonical peptide detections in proteomic data from at least ten genes were selected. The CERES scores⁸¹ of genes with non-canonical peptide hits were compared to those without, using the Mann-Whitney U-test. Additionally, pooled CERES scores across all genes and cell lines were compared between the two groups using the same test. For *KRAS*, CERES scores and RNA abundance in cell lines with non-canonical peptide detections in proteomic data were compared to those without using the Mann-Whitney U-test.

Spectrum Visualization and Validation

Target PSM experimental spectra were extracted from mzML files using custom scripts based on pyOpenMS⁸² and visualized in R. Theoretical spectra were generated from the target peptide sequences using the TheoreticalSpectrumGenerator module of OpenMS, and hyperscores between the experimental and theoretical spectra were calculated using the

HyperScore module with the same parameters (e.g., fragment mass tolerance) used during database search. Fragment ion matching between the experimental and theoretical spectra was performed using a similar approach to IPSA⁸³. Theoretical spectra with predicted fragment peak intensities were generated using Prosit through the Oktoberfest⁸⁴ Python package using parameters (e.g., fragmentation method and energy) in accordance with the original publication⁴³. Similarities between the experimental spectra and the Prosit predicted spectra were estimated using cross-correlation⁸⁵, using the same parameters (e.g., fragment_bin_offset) during database search with Comet. To assess the distribution of cross-correlation values for variant peptide PSMs, we randomly selected 1,000 PSMs from the canonical database search of each of the 42 TMT-plexes as control. Spectra matched to peptides derived from circRNA events were validated using the Novor algorithm through app.novor.cloud, using parameters (e.g., fragmentation method, MS2 analyzer, enzyme, precursor and fragment mass tolerance) consistent with database searches⁴⁶.

Cohort Level FDR

A *post-hoc* approach was employed to estimate the FDR threshold at the cohort level for each database tier. Within each sample and database tier, we first identified the target hit with the highest FDR value under the 1% threshold, denoted as FDR_i .

$$FDR_i = \max_{j \in \{1, 2, \dots, n\}} (FDR_{j,k} | FDR_{j,k} < 0.01)$$

The number of decoy and target hits with FDR values less than FDR_i for each sample were tallied. The equivalent cohort-level FDR threshold was then calculated by dividing the total number of decoy hits by the total number of target and decoy hits across the cohort.

$$\text{Cohort Level FDR Cutoff} = \frac{\sum 1(\text{Decoy Hits } FDR_j < FDR_i)}{\sum 1(\text{Target \& Decoy Hits } FDR_j \leq FDR_i)}$$

Online Methods References

48. Di Tommaso, P. *et al.* Nextflow enables reproducible computational workflows. *Nat Biotechnol* 35, 316–319 (2017).
49. Patel, Y. *et al.* NFTTest: automated testing of Nextflow pipelines. *Bioinformatics* 40, (2024).
50. Patel, Y. *et al.* Metapipeline-DNA: A Comprehensive Germline & Somatic Genomics Nextflow Pipeline. *bioRxiv* 2024.09.04.611267 (2024) doi:10.1101/2024.09.04.611267.
51. Zhao, H. *et al.* CrossMap: a versatile tool for coordinate conversion between genome assemblies. *Bioinformatics* 30, 1006–1007 (2014).
52. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* 17, 122 (2016).
53. Frankish, A. *et al.* GENCODE 2021. *Nucleic Acids Res* 49, D916–D923 (2021).
54. Cunningham, F. *et al.* Ensembl 2022. *Nucleic Acids Res* 50, D988–D995 (2022).
55. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* 10, (2021).
56. Vasimuddin, Md., Misra, S., Li, H. & Aluru, S. Efficient Architecture-Aware Acceleration of BWA-MEM for Multicore Systems. in *2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS)* 314–324 (IEEE, 2019). doi:10.1109/IPDPS.2019.00041.
57. Poplin, R. *et al.* Scaling accurate genetic variant discovery to tens of thousands of samples. *bioRxiv* 201178 (2018) doi:10.1101/201178.
58. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet* 43, 491–8 (2011).
59. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* 17, 122 (2016).
60. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29, 15–21 (2013).
61. Haas, B. J. *et al.* Accuracy assessment of fusion transcript detection via read-mapping and de novo fusion transcript assembly-based methods. *Genome Biol* 20, 213 (2019).
62. Shen, S. *et al.* rMATS: robust and flexible detection of differential alternative splicing from replicate RNA-Seq data. *Proc Natl Acad Sci U S A* 111, E5593-601 (2014).
63. Lo Giudice, C., Tangaro, M. A., Pesole, G. & Picardi, E. Investigating RNA editing in deep transcriptome datasets with REDIttools and REDIportal. *Nat Protoc* 15, 1098–1131 (2020).

64. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* 27, 2156–2158 (2011).
65. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* 10, 1–4 (2021).
66. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29, 15–21 (2013).
67. Zhang, X. O. *et al.* Diverse alternative back-splicing and alternative splicing landscape of circular RNAs. *Genome Res* 26, 1277–1287 (2016).
68. Chambers, M. C. *et al.* A cross-platform toolkit for mass spectrometry and proteomics. *Nat Biotechnol* 30, 918–20 (2012).
69. Mellacheruvu, D. *et al.* The CRAPome: a Contaminant Repository for Affinity Purification Mass Spectrometry Data. *Nat Methods* 10, 730 (2013).
70. Eng, J. K., Jahan, T. A. & Hoopmann, M. R. Comet: An open-source MS/MS sequence database search tool. *Proteomics* 13, 22–24 (2013).
71. Röst, H. L. *et al.* OpenMS: a flexible open-source software platform for mass spectrometry data analysis. *Nature Methods* 2016 13:9 13, 741–748 (2016).
72. Weisser, H. & Choudhary, J. S. Targeted Feature Detection for Data-Dependent Shotgun Proteomics. *J Proteome Res* 16, 2964–2974 (2017).
73. Berry, M. J., Harney, J. W., Ohama, T. & Lhatfield, D. Selenocysteine insertion or termination: factors affecting UGA codon fate and complementary anticodon:codon mutations. *Nucleic Acids Res* 22, 3753–3759 (1994).
74. Wen, B., Li, K., Zhang, Y. & Zhang, B. Cancer neoantigen prioritization through sensitive and reliable proteogenomics analysis. *Nature Communications* 2020 11:1 11, 1–14 (2020).
75. Kong, A. T., Leprevost, F. V., Avtonomov, D. M., Mellacheruvu, D. & Nesvizhskii, A. I. MSFragger: ultrafast and comprehensive peptide identification in mass spectrometry-based proteomics. *Nat Methods* 14, 513–520 (2017).
76. Craig, R. & Beavis, R. C. TANDEM: matching proteins with tandem mass spectra. *Bioinformatics* 20, 1466–7 (2004).
77. Gessulat, S. *et al.* Prosit: proteome-wide prediction of peptide tandem mass spectra by deep learning. *Nat Methods* 16, 509–518 (2019).
78. Szolek, A. *et al.* OptiType: precision HLA typing from next-generation sequencing data. *Bioinformatics* 30, 3310–6 (2014).

79. O'Donnell, T. J., Rubinsteyn, A. & Laserson, U. MHCflurry 2.0: Improved Pan-Allele Prediction of MHC Class I-Presented Peptides by Incorporating Antigen Processing. *Cell Syst* 11, 42-48.e7 (2020).
80. P'ng, C. *et al.* BPG: Seamless, automated and interactive visualization of scientific data. *BMC Bioinformatics* 20, 42 (2019).
81. Meyers, R. M. *et al.* Computational correction of copy number effect improves specificity of CRISPR-Cas9 essentiality screens in cancer cells. *Nat Genet* 49, 1779–1784 (2017).
82. Röst, H. L., Schmitt, U., Aebersold, R. & Malmström, L. pyOpenMS: a Python-based interface to the OpenMS mass-spectrometry algorithm library. *Proteomics* 14, 74–77 (2014).
83. Brademan, D. R., Riley, N. M., Kwiecien, N. W. & Coon, J. J. Interactive Peptide Spectral Annotator: A Versatile Web-based Tool for Proteomic Applications. *Mol Cell Proteomics* 18, S193–S201 (2019).
84. Picciani, M. *et al.* Oktoberfest: Open-source spectral library generation and rescoring pipeline based on Prosit. *Proteomics* 2300112 (2023) doi:10.1002/PMIC.202300112.
85. Eng, J. K., Fischer, B., Grossmann, J. & MacCoss, M. J. A fast SEQUEST cross correlation algorithm. *J Proteome Res* 7, 4598–4602 (2008).

Acknowledgments

The authors thank members of the Boutros and Kislinger labs for their continued support, particularly Dr. Amanda Khoo, Meinusha Govindarajan and Dr. Matthew Waas. The authors also thank Dr. James Wohlschlegel from UCLA Proteome Research Center and Dr. Mathias Wilhelm from Technical University of Munich. This work was supported by the NIH *via* awards P30CA016042, U01CA214194, U2CCA271894, U24CA248265, P50CA092131 and R01CA244729, by the Canadian Cancer Society *via* an Impact Grant (705649) and by the Canadian Institute of Health Research *via* a Project Grant (PJT156357). CZ was supported by the UCLA Jonsson Comprehensive Cancer Center Fellowship Award. LYL was supported by a CIHR Vanier Fellowship and Ontario Graduate Scholarship. HZ was supported by a CIHR Doctoral Award. TK is supported through the Canadian Research Chair program. University Health Network was supported by the Ontario Ministry of Health and Long-Term care.

Data Availability

Data supporting the conclusions of this paper is included within it and its supplementary files. The processed CCLE data are available at the DepMap portal (<http://www.depmap.org>). The raw WGS and WXS cell lines sequencing data are available at Sequence Read Archive (SRA) and European Genome-Phenome Archive (EGA) under access number PRJNA523380 and EGAD00001001039. The raw mass spectrometry proteomic data are publicly available without restrictions at the ProteomeXchange *via* the PRIDE partner repository under accession number PXD030304 for cell lines, PXD030983 for mouse strain C57BL/6N, and PXD010154 for alternative protease and fragmentation analyses. The proteomic data for the five prostate tumour samples are freely available at UCSD's MassIVE database under accession number MSV000081552, whereas their raw WGS and RNA-seq data are available at EGA under accession EGAS00001000900. Proteomic data for the eight kidney tumour samples are freely available at Proteomic Data Commons (PDC) under accession number PDC000411, whereas the genomic and transcriptomic data are available at Genomic Data Commons (GDC, Project: CPTAC-3, Primary Site: Kidney) with dbGaP accession number phs001287, generated by the National Cancer Institute's Clinical Proteomic Tumor Analysis Consortium (CPTAC).

Code Availability

moPepGen is publicly available at: <https://github.com/uclahs-cds/package-moPepGen>. Data processing, analysis and visualization scripts are available upon request.

Conflicts of Interest

PCB sits on the Scientific Advisory Boards of Intersect Diagnostics Inc., BioSymetrics Inc. and previously sat on that of Sage Bionetworks. All other authors declare no conflicts of interest.

Figure Legends

Figure 1: moPepGen is a graph-based algorithm that uncovers non-canonical peptides with variant combinations

a) moPepGen algorithm schematic. moPepGen is a graph-based algorithm that generates databases of non-canonical peptides that harbour genomic and transcriptomic variants (e.g., single nucleotide variants [SNVs], small insertions and deletions [indels], RNA editing, alternative splicing, gene fusion and circular RNAs [circRNAs]) from coding transcripts, as well as from novel open reading frames of noncoding transcripts. **b)** and **c)** moPepGen achieves linear runtime complexity when fuzz testing with SNVs only (**b**) and with SNVs and indels (**c**). **d)** A variant peptide from *SYNPO2* that harbours a small deletion and an SNV. Fragment ion mass spectrum from peptide-spectrum match (PSM) of the non-canonical peptide harbouring two variants (top, both) is compared against the canonical peptide theoretical spectra (left, theoretical spectra at the bottom) and against the variant peptide theoretical spectra (right, bottom). Fragment ion matches are colored, with b-ions in blue and y-ions in red. **e-g)** A somatic SNV D1249N in *AHNAK* was detected in DNA sequencing at chr11:62530672 (**e**), in RNA-sequencing (**f**) and as the non-canonical peptide MDIDAPDVEVQGP~~N~~WHLK (**g**). **h-i)**: Fragment ion mass spectrum from PSM of the canonical peptide MDIDAPDVEVQGP~~D~~WHLK (**h**) and the non-canonical peptide (**i**).

Figure 2: moPepGen generates comprehensive non-canonical databases that support for proteogenomic analysis

a) Sizes of variant peptide databases generated by moPepGen using somatic single nucleotide variants, small insertions and deletions and transcript fusions for 376 cell lines from the Cancer Cell Line Encyclopedia project. Color indicates cell line tissue of origin. **b)** Genes with variant peptides detected in cell lines across three or more tissues of origin (bottom covariate). The barplot shows number of recurrences across tissues and color of heatmap indicates number of cell lines. **c)** Number of non-canonical peptides from different variant combinations (bottom heatmap) generated using genomic and transcriptomic data from five primary prostate tumours, shown across four tiers of custom databases and grouped by the number of variant sources in combination. Alternative translation (Alt Translation) sources with ≥ 10 peptides are visualized. **d)** Five variant peptides detected in

one prostate tumour (CPCG0183) from the protein plectin (PLEC). Fragment ion matches are colored, with b-ions in blue and y-ions in red.

Extended Data Figure Legends

Extended Data Figure 1: Core graph algorithm of moPepGen

The graph algorithm of moPepGen implements the following key steps: **a)** A transcript variant graph (TVG) is generated from the transcript sequence with all variants associated. All three reading frames are explicitly generated to efficiently handle frameshift variants. **b)** Variant bubbles of the TVG are aligned and expanded to ensure the sequence length of each node is a multiple of three. **c)** Peptide variant graph (PVG) is generated by translating the sequence of each node of the TVG. **d)** Peptide cleavage graph is generated from the PVG in such a way that each node is an enzymatically cleaved peptide.

Extended Data Figure 2: Differential handling of noncoding transcripts, subgraphs and circular RNAs

a) For coding transcripts, variants are only incorporated into the effective reading frames. For transcripts that are canonically annotated as noncoding, variants are added to all three reading frames to perform comprehensive three-frame translation. **b)** Subgraphs are created for variant types that involve the insertion of large segments of the genome, which can carry additional variants. **c)** The graph of a circular RNA is extended four times to capture all possible peptides that span the back-splicing junction site in all three reading frames. In the bottom panel, the nodes in rose-red harbour the variant 130-A/T and the nodes in yellow harbour 165-A/AC. **d)** Illustration of a circRNA molecule with a novel open reading frame. Each translation across the back-splicing site may shift the reading frame. If no stop codon is encountered, the original reading frame is restored after the fourth crossing.

Extended Data Figure 3: moPepGen demonstrates comprehensive results and deliberate biological assumptions

a) and **b)** Non-canonical peptide generation results from benchmarking of moPepGen, pyQUILTS and customProDBJ using only point mutations and small insertions and deletions (indels; **a)**, and with inputs from point mutations, indels, RNA editing, transcript fusion, alternative splicing and circular RNAs (**b)**. Top boxplot shows the number of peptides in each set intersection and right barplot shows the total number of non-canonical peptides generated by each algorithm in five primary prostate tumour samples. **c)** Assumptions made by moPepGen for handling edge cases that differ from other algorithms. Start-codon-altering

and splice-site-altering variants are omitted due to the uncertainty of the resulting translation and splicing outcomes. Transcripts with unknown stop codons do not have trailing peptide outputs because of the uncertainty of the trailing enzymatic cleavage site. Stop-codon-altering variants do not result in translation beyond the transcript end, adhering to central dogma. **d)** Non-canonical database search results from benchmarking of moPepGen, pyQUILTS and customProDBJ using point mutations, indels, RNA editing, transcript fusion, alternative splicing and circular RNAs.

Extended Data Figure 4: Detection of novel open reading frame peptides across proteases

a) Peptide length distributions after *in silico* digestion with seven enzymes, as indicated by color, of the canonical human proteome and three-frame translated noncoding transcript open reading frames (ORFs). The dotted lines indicate the 7-35 amino acids peptide length range used for database search. **b)** Noncoding peptide detection across ten enzyme-fragmentation methods. The top barplot shows the number of peptides in each set intersection and the right barplot shows the total number of non-canonical peptides from noncoding ORFs detected in each enzyme-fragmentation method, as indicated by covariate color. **c)** Optimal combinations of one to ten enzyme-fragmentation methods for maximizing the number of transcripts detected from the canonical proteome, or the number of ORFs detected from noncoding transcripts. The bottom covariate indicates the optimal combinations of enzyme-fragmentation methods from combinations of one to ten, with color indicating enzyme-fragmentation method. **d)** Noncoding transcript ORFs with peptides detected across four or more enzyme-fragmentation methods, with recurrence count shown in the right barplot. The color of the heatmap indicates the number of peptides detected per ORF per enzyme-fragmentation method. **e)** Example ORFs with coverage by multiple proteases are shown, with peptides tiled according to detection in each enzyme-fragmentation method, as indicated by covariate color. Representative fragment ion mass spectra of peptide-spectrum matches are shown, with theoretical spectra at the bottom and fragment ion matches colored (blue: b-ions, red: y-ions in).

Extended Data Figure 5: Germline non-canonical peptide detection in mouse strain C57BL/6N

a) Comparison of canonical and custom database sizes for the C57NL/6N mouse. Germline database includes single nucleotide polymorphisms (SNPs) and short insertions and deletions. **b)** Number of non-canonical peptides detected from each database in each tissue, with database indicated by color. **c)** Comparison of a variant peptide-spectrum match (PSM) spectra (top, both) with the theoretical spectra of the canonical peptide counterpart (left, bottom) as well as the theoretical spectra of the variant peptide harbouring a SNP (right, bottom). Fragment ion matches are colored, with b-ions in blue and y-ions in red. **d)** Noncoding transcripts with open reading frames yielding two or more non-canonical peptides recurrently detected across tissues, with color indicating the number of peptides detected in each tissue.

Extended Data Figure 6: Proteogenomic investigation of the Cancer Cell Line Encyclopedia

a) Number of non-canonical peptides generated per cell line, with color indicating peptide source. Bottom covariate indicates tissue of origin. **b)** and **c)** Number of variant peptides per cell line with given number of variants in coding (**b**) and noncoding transcripts (**c**). **d)** Number of non-canonical peptides detected in each cell line, with color representing peptide source. Bottom covariate indicates tissue of origin. **e)** Per cell line, the number of variant effect predictor (VEP) annotated intragenic coding mutations, mutations predicted to produce detectable non-canonical peptides and mutations detected through proteomics. **f)** Per cell line, number of transcript fusions, fusions theoretically able to produce detectable non-canonical peptides and fusions with detected peptide products. Color indicates tissue of origin. **g)** Fusion transcripts (upstream transcript gene symbol – downstream transcript gene symbol) with detected peptide products, with number of peptides shown across cell lines. Color indicates whether the upstream fusion transcript was coding or noncoding. **h)** Fragment ion mass spectrum from peptide-spectrum match (PSM) of the non-canonical peptide at the junction of the *FLNB-SLMAP* fusion transcript. The peptide theoretical spectrum is shown at the bottom and fragment ion matches are colored (blue: b-ions, red: y-ions in). **i)** Comparison of mass spectrum (top, both) from PSM of a non-canonical peptide with a single nucleotide variant against Prosit-predicted MS2 mass spectra based on the canonical counterpart

peptide sequence (left, bottom) and the detected variant peptide sequence (right, bottom). Fragment ion matches are colored, with b-ions in blue and y-ions in red. **j**) Cross-correlation (Xcorr) distribution of PSMs of coding variant non-canonical peptides against ProSight-predicted MS2 mass spectra (solid lines, color indicate charge), in comparison with Xcorr of control canonical PSMs against ProSight-predicted mass spectra (dotted lines).

Extended Data Figure 7: Functional investigation of non-canonical peptide detection in Cancer Cell Line Encyclopedia

a) Gene dependency CERES scores for genes with detected non-canonical peptides (orange), detected canonical peptides only (pink) and no detected peptides (gray). A lower CERES score indicates higher gene dependency. Cell lines were selected based on the detection of non-canonical peptides in more than 10 genes. P-values were calculated using the Mann-Whitney U-test. The red vertical line indicates $\alpha = 0.05$. The bottom panel represents data pooled across all genes and cell lines. **b**) Gene dependency CERES score and **c**) mRNA abundance of *KRAS* in cell lines with only canonical peptides detected compared to those with detected non-canonical peptides. P-values were calculated using the Mann-Whitney U-test. **d**) Number of putative neoantigens predicted based on detected non-canonical peptides in cell lines with more than two neoantigens. The color indicates cell line tissue of origin. **e**) Recurrent neoantigens observed across multiple cell lines, along with their associated gene, variant, *HLA* genotype and the full peptide sequence as detected by trypsin-digested whole cell lysate mass spectrometry. The color in the left heatmap represents neoantigen affinity.

Extended Data Figure 8: Detection of non-canonical peptides from DIA proteomics

a) Number of variant peptides from different variant combinations generated using genomic and transcriptomic data from eight clear cell renal cell carcinoma (ccRCC) tumours, grouped by the number of variant sources in combination. **b**) Number of detected variant peptides in the data-independent acquisition (DIA) proteome of eight ccRCC tumours. **c-e**) Detection of non-canonical peptides harbouring germline single nucleotide polymorphisms (SNPs; **c**), alternative splicing (**d**) and RNA editing sites (**e**) across genes. Heatmap colors indicate the number of peptides detected per gene per sample. The barplot indicates recurrence across samples. **f**) Illustration of non-canonical peptides derived from the canonical sequence

FSGSNSGNTATLTISR in gene *IGLV3-21* caused by RNA editing events. **g-i)** Extracted ion chromatograms of the canonical peptide (**g**) and non-canonical peptides derived from *IGLV3-21* caused by RNA editing events: chr22:22713097 G-to-C (**h**) and chr22:22713111 A-to-G (**i**).

Extended Data Figure 9: Detection of non-canonical peptides from genomic variants, alternative splicing and circular RNAs

a) Number of detected non-canonical peptides in five primary prostate tumour samples per database tier (colored by database). **b)** Peptides as the result of a combination of two variants, with variant type indicated in left covariate and gene on the right. The heatmap shows presence of peptide across samples. **c-f)** Non-canonical peptide detection results across genes, with color of heatmap representing the number of peptides detected per gene per sample. The barplot indicates recurrence across samples, and when colored indicates variant type associated with the gene entry. The Variant database includes non-canonical peptides from coding transcripts with single nucleotide variants, small insertion and deletion, RNA editing, alternative splicing or transcript fusion (**c**). Noncoding database includes all peptides from noncoding transcript three-frame translation open reading frames (**d**) and noncoding peptides with any variants are included in the Noncoding Variant database (**e**). The Circular RNA database includes all peptides representing circular RNA open reading frames (ORFs) with or without other variants (**f**). The bottom covariate indicates prostate cancer sample. **g)** Mass spectrum from peptide-spectrum match of a non-canonical peptide spanning the back-splicing junction between exon 29 and exon 24 of *MYH10*, reflective of circular RNA translation. The peptide theoretical spectrum is shown at the bottom and fragment ion matches are colored (blue: b-ions, red: y-ions in).

Extended Data Tables

Extended Data Table 1: Feature comparison of custom database generation algorithms

Algorithm / Feature	moPepGen	customProDBj ¹⁶⁻¹⁸	MaxQuant module ¹⁹	ProteoDisco ²⁰	ProteomeGenerator ^{31,33}	pypgatk ²¹	pyQUILTS ²²	samplespecificDBGenerator ²³	sapFinderPGA ^{26,27}	spliceDB ^{28,29}	Spritz ³⁰	NeoDisc ³²
DNA Single Nucleotide Variants	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
DNA Small Insertions / Deletions	✓	✓		✓	✓	✓	✓		✓	✓	✓	✓
RNA Point Variants ^a	✓										✓	
Alternative Splice Junctions	✓			✓	✓		✓	✓	✓	✓		
Fusion Transcripts	✓			✓	✓		✓	✓	✓	✓		
Circular RNAs	✓											
Supported Variant Combinations	Any combinations ^b	None	All SNVs only ^c	All SNVs only ^c	All variants only ^c	None	Single SNV with splicing or fusion ^d	None	None	All variants only ^c	None	All variants only ^c
Modular Support for New Input Formats	✓							✓				
Noncoding Transcript Three-Frame Translation	✓				✓	✓			✓			✓
Noncoding Transcript with Variants ^e	✓				✓							✓
Alternative Translation W>F	✓											
Export Only Non-canonical Proteotypic Peptides	✓			✓				✓	✓			
Summary of Database Generation	✓	✓			✓				✓			✓
Visualization of Database Generation	✓								✓			✓
Database Splitting for Tiered False Discovery Rate Control	✓											
Filter FASTA by RNA Abundance	✓				✓			✓				✓

^aRNA Point Variants: variants at single nucleotide positions within RNA sequences (e.g., RNA editing).

^bAny variant combinations: generates peptides with any combination of variants.

^cAll SNVs only / All variants only: generates peptides containing all variants simultaneously, but not separate combinations of individual variants.

^dSingle SNV with splicing or fusion: generates peptides with a single SNV, with or without additional alternative splicing or transcript fusion events.

^eNoncoding Transcript with Variants: non-canonical peptides derived from novel open reading frames in noncoding transcripts harbouring additional DNA and/or RNA variants.

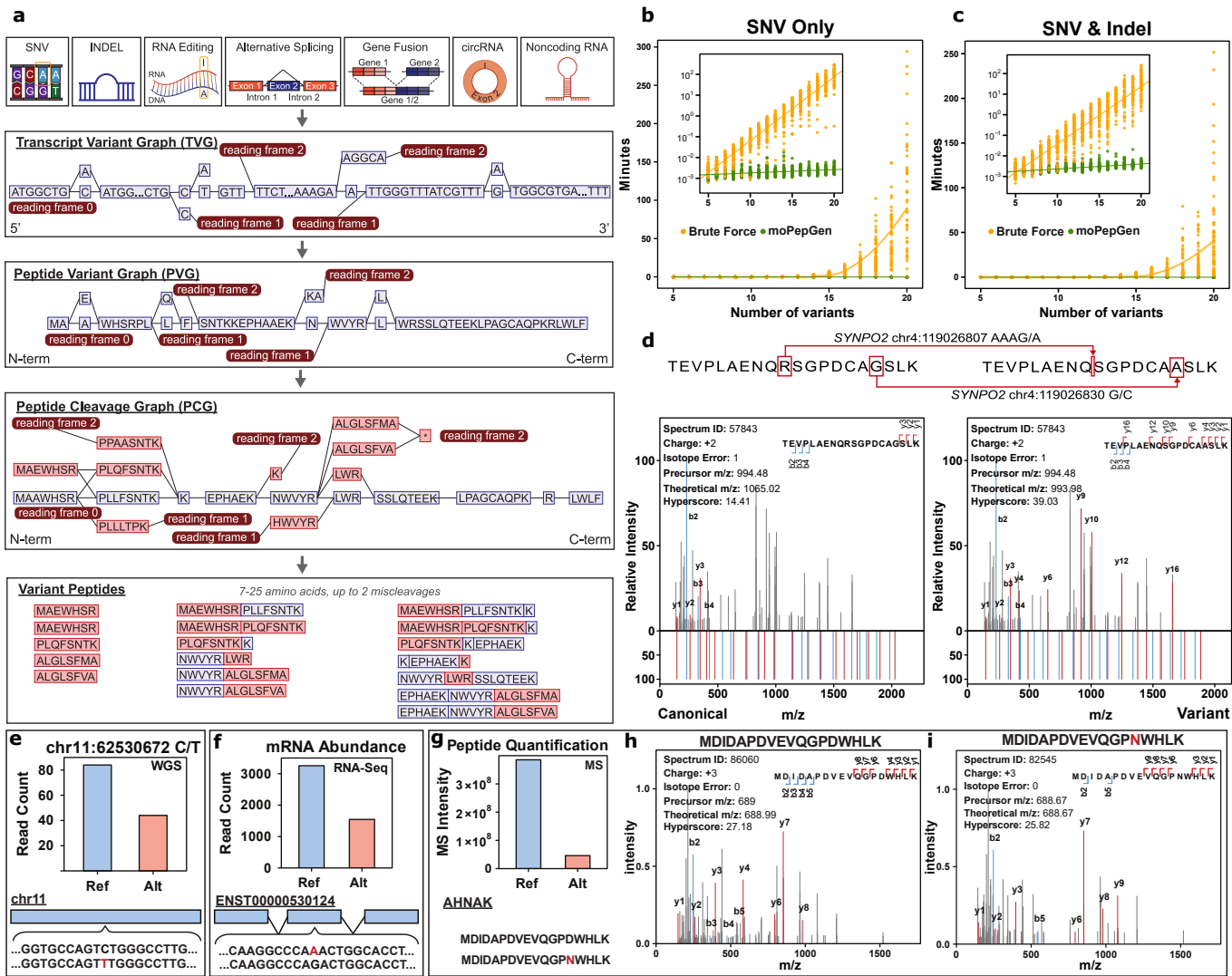
Figure 1

Figure 2