

Supplementary Materials for Manuscript

Comparative Proteomic Profiling of Unannotated Microproteins and Alternative Proteins in Human Cell Lines

Xiongwen Cao^{1,2,4}, Alexandra Khitun^{1,2,4}, Zhenkun Na^{1,2}, Daniel G. Dumitrescu¹,
Marcelina Kubica^{2,5}, Elizabeth Olatunji^{2,5}, Sarah A. Slavoff^{1,2,3*}

¹Department of Chemistry, Yale University, New Haven, Connecticut 06520, United States

²Chemical Biology Institute, Yale University, West Haven, Connecticut 06516, United States

³Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut 06529, United States

⁴These authors contributed equally

⁵These authors contributed equally

*Correspondence: sarah.slavoff@yale.edu

.

Table of Contents

Supplementary Figures

Figure S1. mRNA-Seq support for novel SCRIB transcript variant.

Figure S2. Confirmation of differential expression of three additional detected alt-proteins with extracted ion chromatograms (EICs).

Figure S3. Predicted start codons used by the six selected alt-ORFs.

Figure S4. Results of TargetP bioinformatic analysis.

Figure S5. Nuclear localization signals predicted using cNLS-mapper.

Figure S6. Amino acid-level conservation of alt-proteins in other species.

Figure S7. Structure prediction of two nuclear alt-proteins.

Figure S8. Uncropped Western blots for confirmation of alt-protein expression (related to Figure 5).

Figure S9. K562 MS/MS spectra.

Figure S10. MOLT4 MS/MS spectra.

Figure S11. MS/MS spectra of peptides found in both cell lines.

Supplementary Tables (all in separate Excel sheets)

Table S1. Full list of peptides and proteins detected in each sample. (XLSX)

Table S2. PepQuery results. (XLSX)

Table S3. Full list of peptides detected from unannotated alt-proteins. (XLSX)

Table S4. TargetP 2.0 prediction results. (XLSX)

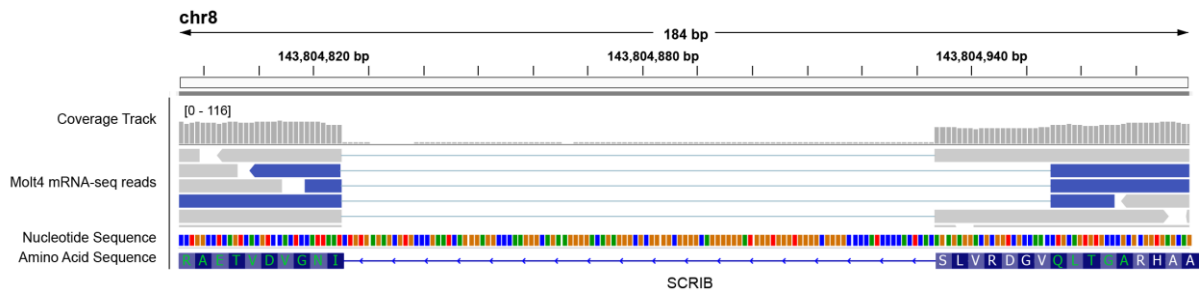


Figure S1. mRNA-Seq support for novel SCRIB transcript variant. Overall mRNA-seq coverage is shown in the top coverage track. mRNA-seq reads from MOLT4 cells are shown in the second track from the top. mRNA-seq reads supporting the novel splice junction between exons 20 and 21 of the SCRIB gene are highlighted in blue. Reads mapping to the annotated splice junction are shown in gray. Note that there is a small deletion in the novel variant, compared with the annotated transcripts. The genomic nucleotide sequence is displayed on the nucleotide sequence track with green, red, blue, and orange tiles representing A,T,C and G, respectively. The bottom track represents the annotated SCRIB protein sequence with amino acids highlighted in green corresponding to the peptide (sequence, AGTLQINGVDVTEAR) detected in our proteomics data which aligns exclusively to the alternatively spliced variant. Note that seven amino acids (VGDRVLS) are excluded in the protein translated by the novel variant, compared with the annotated SCRIB protein.

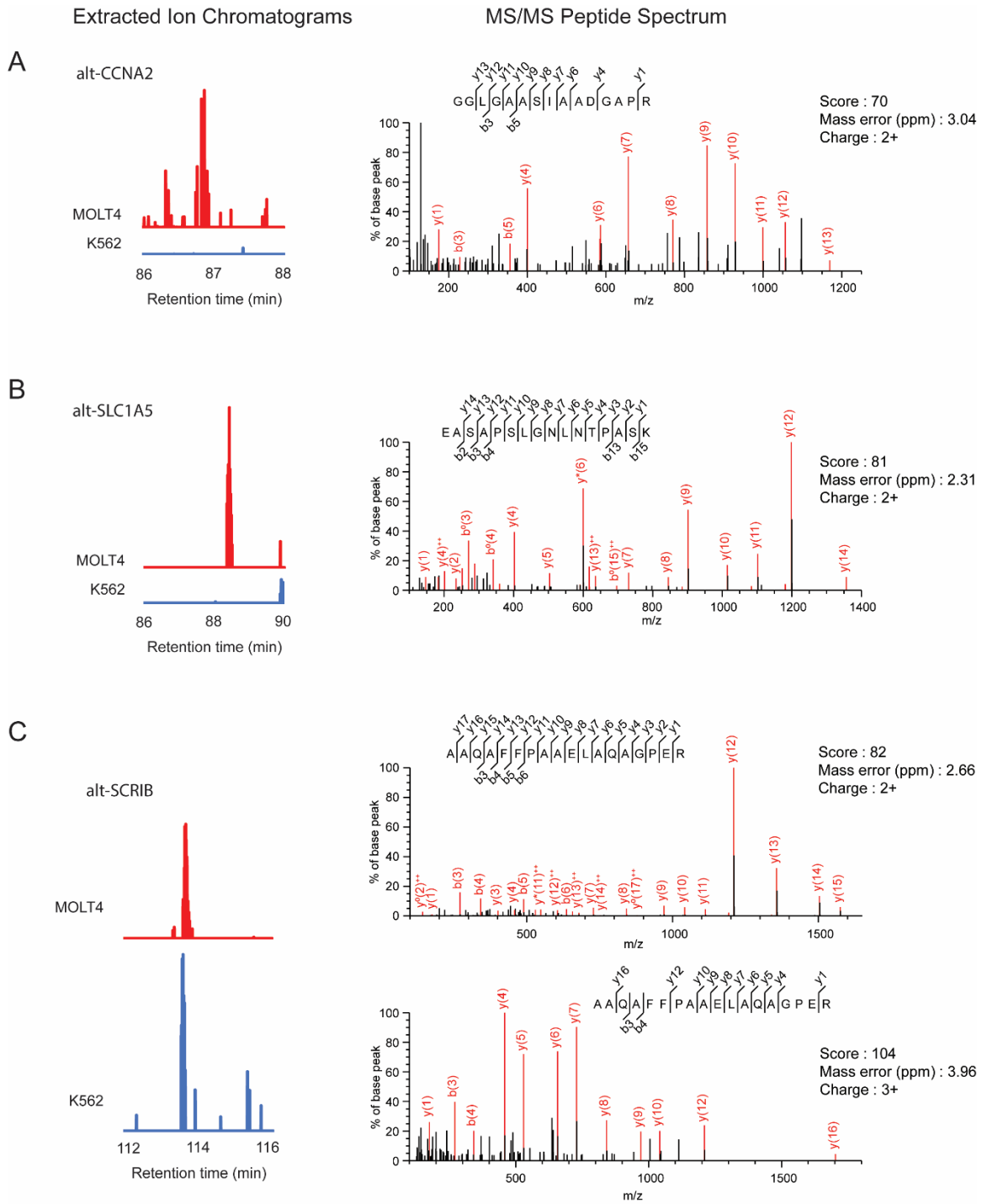


Figure S2. Confirmation of differential expression of three additional detected alt-proteins with extracted ion chromatograms (EICs). Shown are EICs (left) from MS₁ spectra corresponding to MS/MS spectra (right) of three tryptic peptides identified in MOLT4 cells only (A-B) or both cell lines (C). The same y-axis scale is used for each matched EIC pair. The EIC intensity at the same retention time was compared for the paired samples.

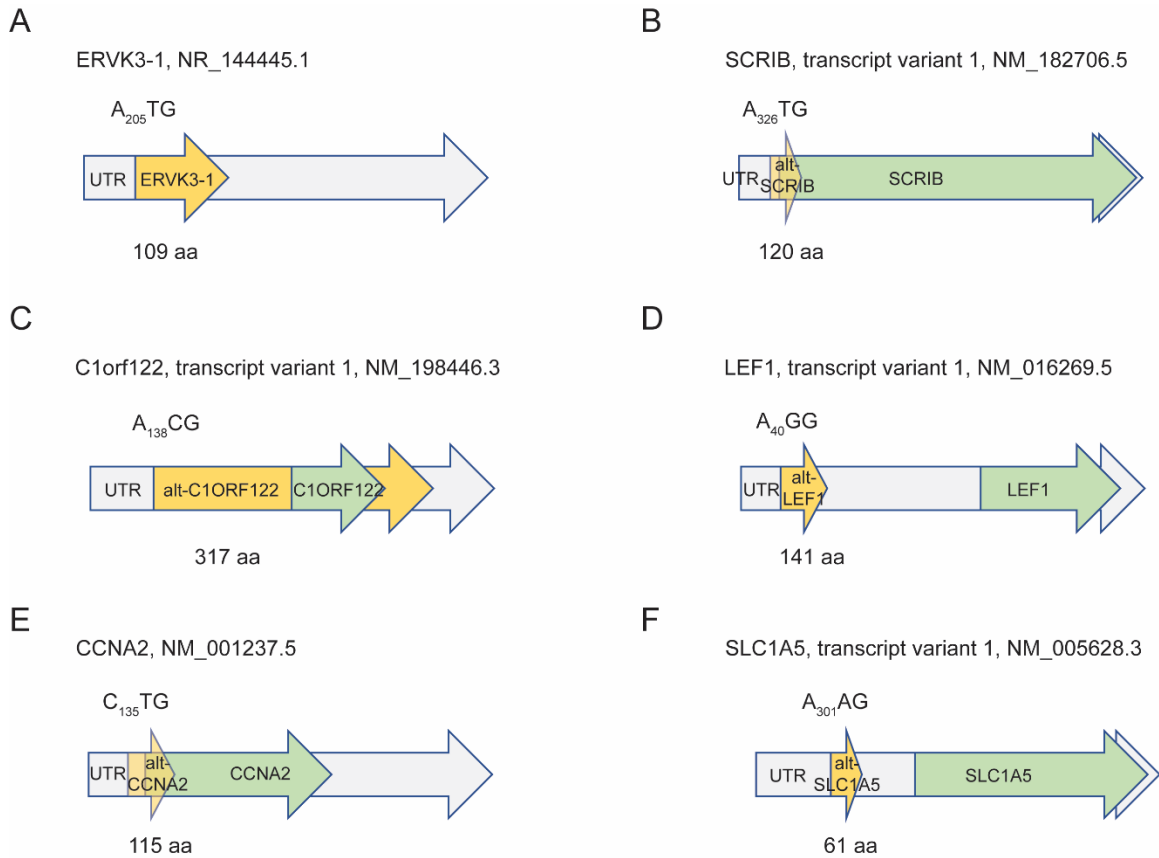


Figure S3. Predicted start codons used by the six selected alt-ORFs. Shown are schematic representations of the transcripts encoding the six selected alt-proteins. Top, transcript name and the accession number; light gray arrow, untranslated regions (UTR); yellow arrow, alt-ORF; green arrow, annotated coding sequence; the predicted start codon used by the alt-ORF is indicated above the transcript, numbered relative to the first nucleotide of the corresponding cDNA, and the predicted full length of the alt-protein is indicated at the bottom. aa, amino acid. (A-B) Two alt-proteins detected in both K562 and MOLT4 cells, (C) one detected in K562 only, and (D-F) three detected in MOLT4 only.

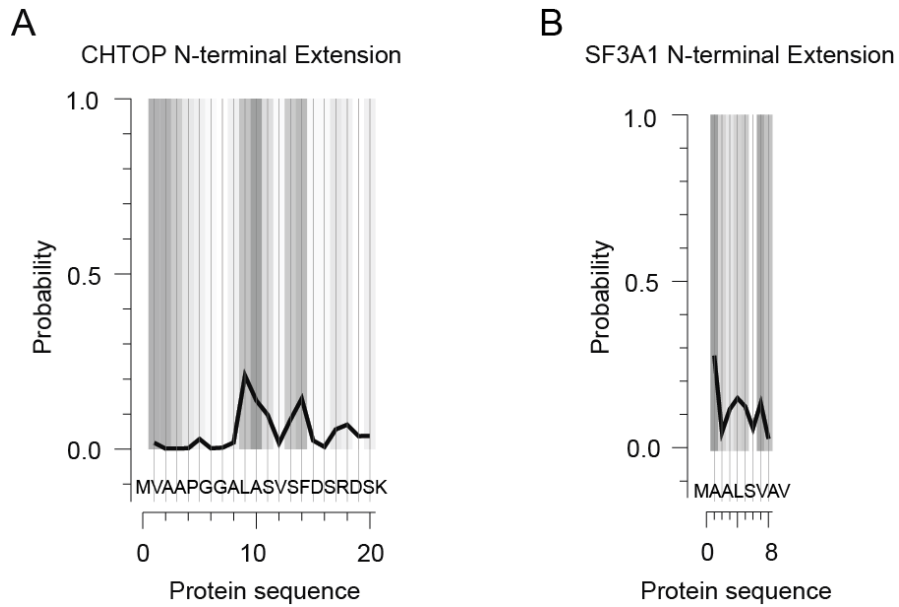


Figure S4. Results of TargetP bioinformatic analysis. (A) N-terminal extension of CHTOP protein is predicted to encode a mitochondrial transfer peptide (B) TargetP prediction result for SF3A1 N-terminal extension where the peak corresponds to a predicted signal peptide. The protein names are indicated at the top of each diagram. TargetP 2.0 probabilistically scores localization signals using a deep-learning model trained on proteins with known localization within three classes 1) mitochondrial transfer peptide 2) signal peptide and 3) other. Only sequences predicted to encode mitochondrial or signal peptides with >40% probability are represented. See Table S4 for full TargetP results.

A

Alt-C1orf122, cNLS-mapper score = 6

MATTAPARSSATAARSASVQPARSGLAACTAPLPGSRNSRAPGAAGA
GLGGRKRRLRE^{PAGPSLNPTLAATAALIPLHRRAGDIRPGPLPGGSD}
APSQLPVQETRPRLCAGQLKGTTPRRIEGDRWGRGGAQPCEALTQ
RWGGGRPKGGRWSGRASGDGMGPGLRLVTGGCRRGPREGGAG
ARREATPTAAPGGARPAAGRGAAAAAPGHHRSLRGDVTAGPP
APGAGWWRERLSQTWSAPPAGCLRQRRLSKGCW^{RWSCGALTIPEQ}
NTLTSLPPQGRWLDSEQLPSVKGPVFTGSGWYLALSLEWQLC

C

Alt-LEF1, cNLS-mapper score = 5

MARRRDPGCALSGRGTRTRPEPRTRPSSGRPLPSAGRTQGAQLFAL
TELAGGGVQSGEPASQAEKLEPGTKRGRTE^{CVCRLELRAEAFGPEA}
PAVTPRDSAVPSTAESPRLPAKTLFLANFSFSSPPPRPPSSAPPSLA
D

B

A0A499FIZ0, cNLS-mapper score = 14.5

MASLGEETLASASSSSSDSDTGGASPPPRKKPRASAAEGVGEPEGASA
GRAGLSPPSSSSSSSSSSSSVVVVVGLPPAAAPPA^{AAAAVPHRSSGH}
SLVSGSI

Figure S5. Three alt-proteins are predicted to encode nuclear localization signals (red) by cNLS-mapper. (A-C) The cNLS-mapper prediction results. The name of alt-proteins and the cNLS mapper score are indicated above each sequence. The default cutoff score of cNLS-mapper is 5.

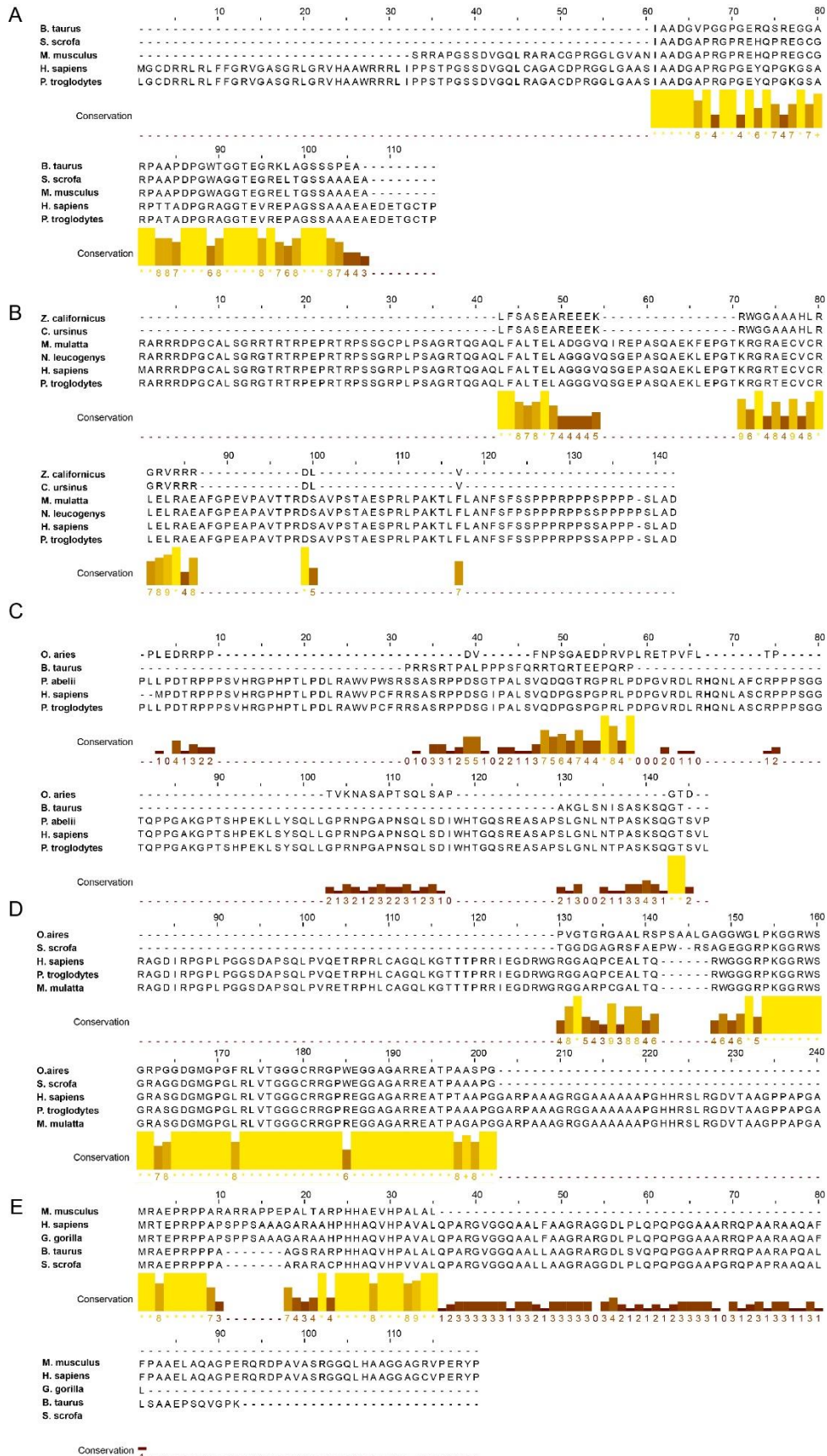
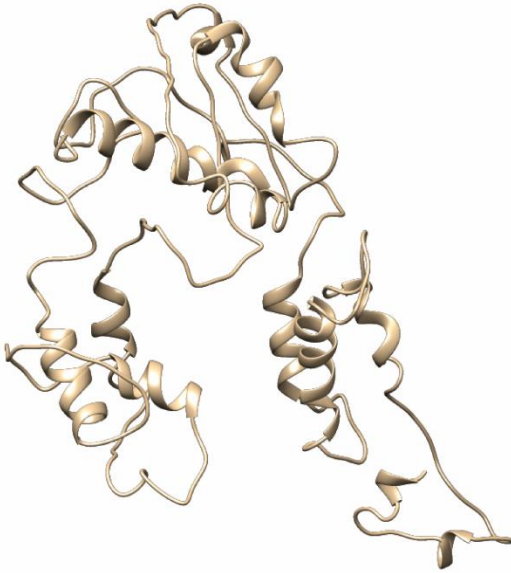


Figure S6. Amino acid-level conservation of alt-proteins in other species. (A) alt-CCNA2, (B) alt-LEF1, (C) alt-SLC1A5, (D) alt-C1orf122 (partial sequence) (E) alt-SCRIB. Alignments from Clustal Omega were quantified on the basis of physico-chemical properties using JalView. * indicates sequence identity.

A

Alt-C1orf122

**B**

Alt-LEF1

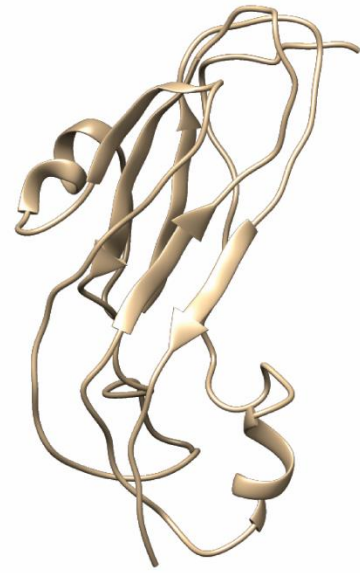


Figure S7. Structure prediction of two nuclear alt-proteins. Predicted structural models for alt-C1orf122, C-score = -1.79 (A) or alt-LEF1, C-score = -3.03 (B) using I-TASSER. UCSF chimera was used to visualize and generate the figure. C-score is a confidence score for estimating the quality of predicted models by I-TASSER, which is typically in the range of [-5, 2].

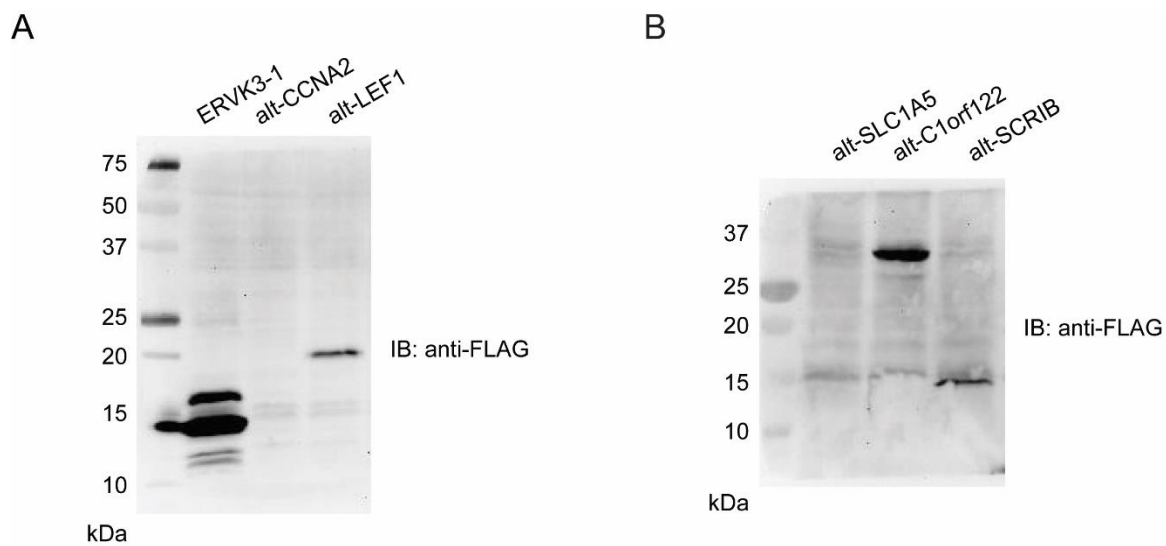
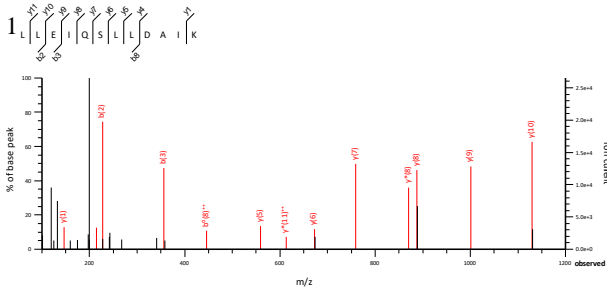
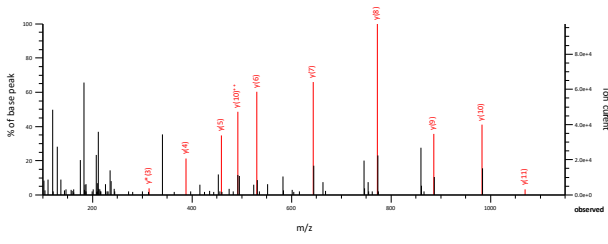


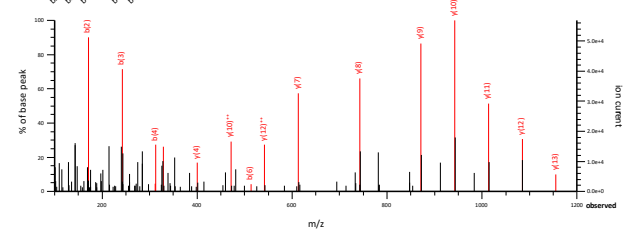
Figure S8. Uncropped Western blots for confirmation of alt-protein expression (related to Figure 5).



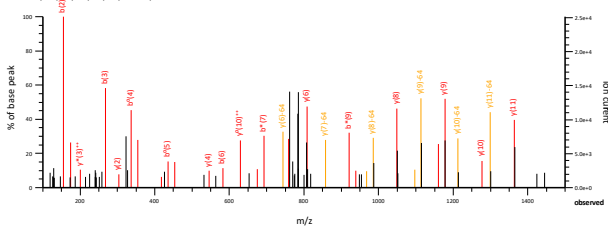
11
 F S P L E L A A G V R

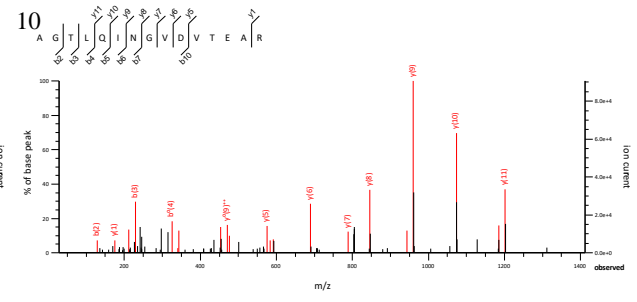
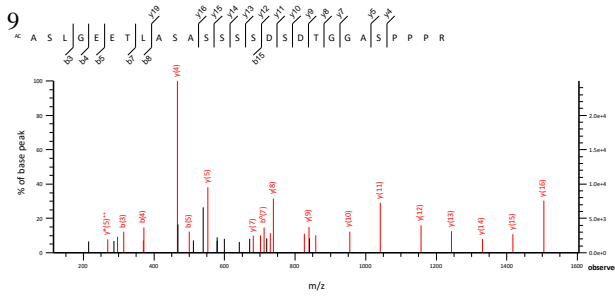
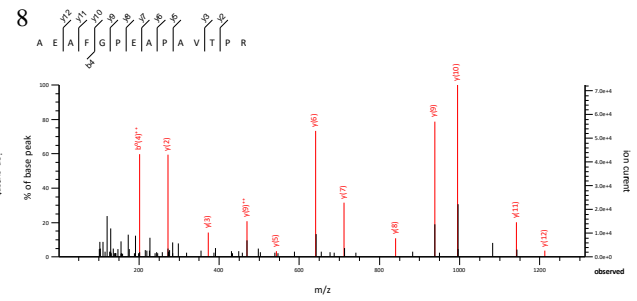
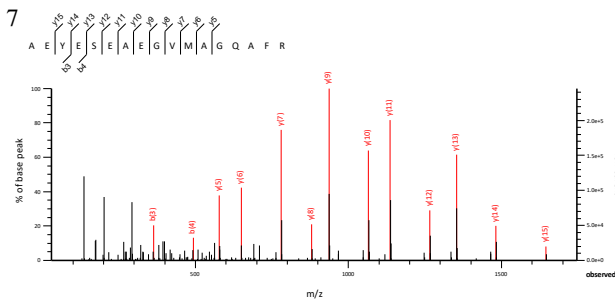
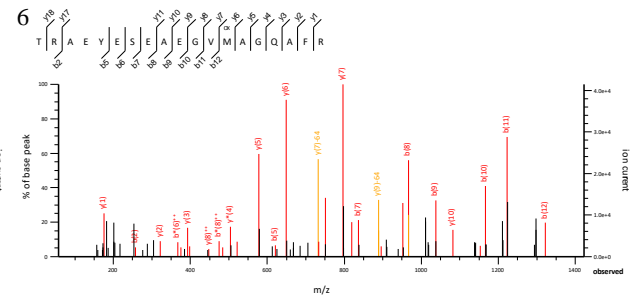
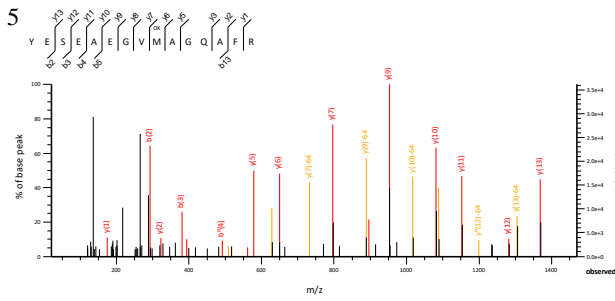
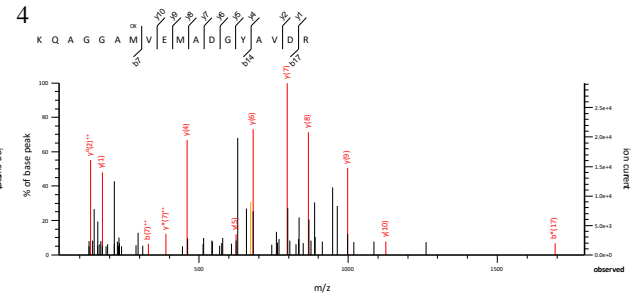
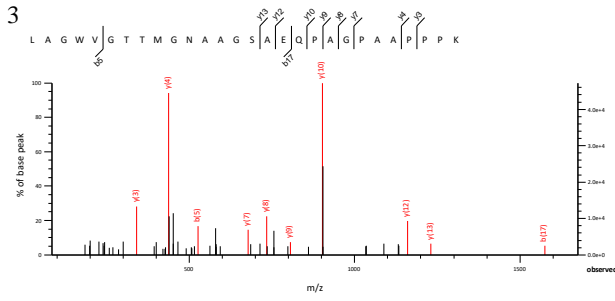
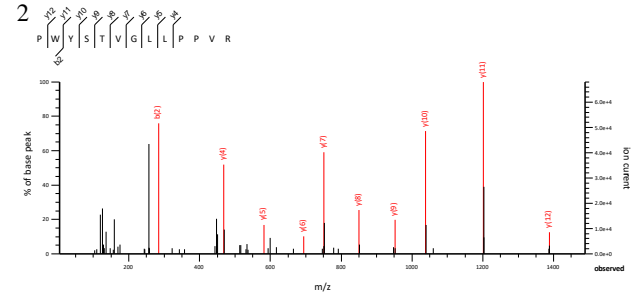
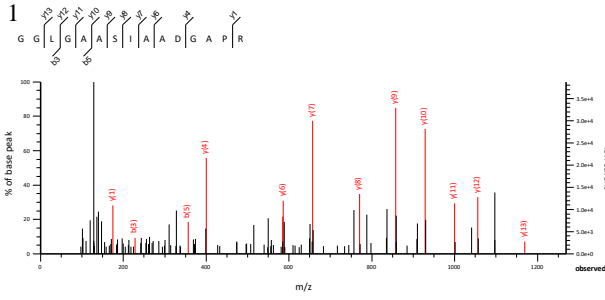


12
 V A A A E A A A G P R



13
 P G S V Q L D M E N Q R





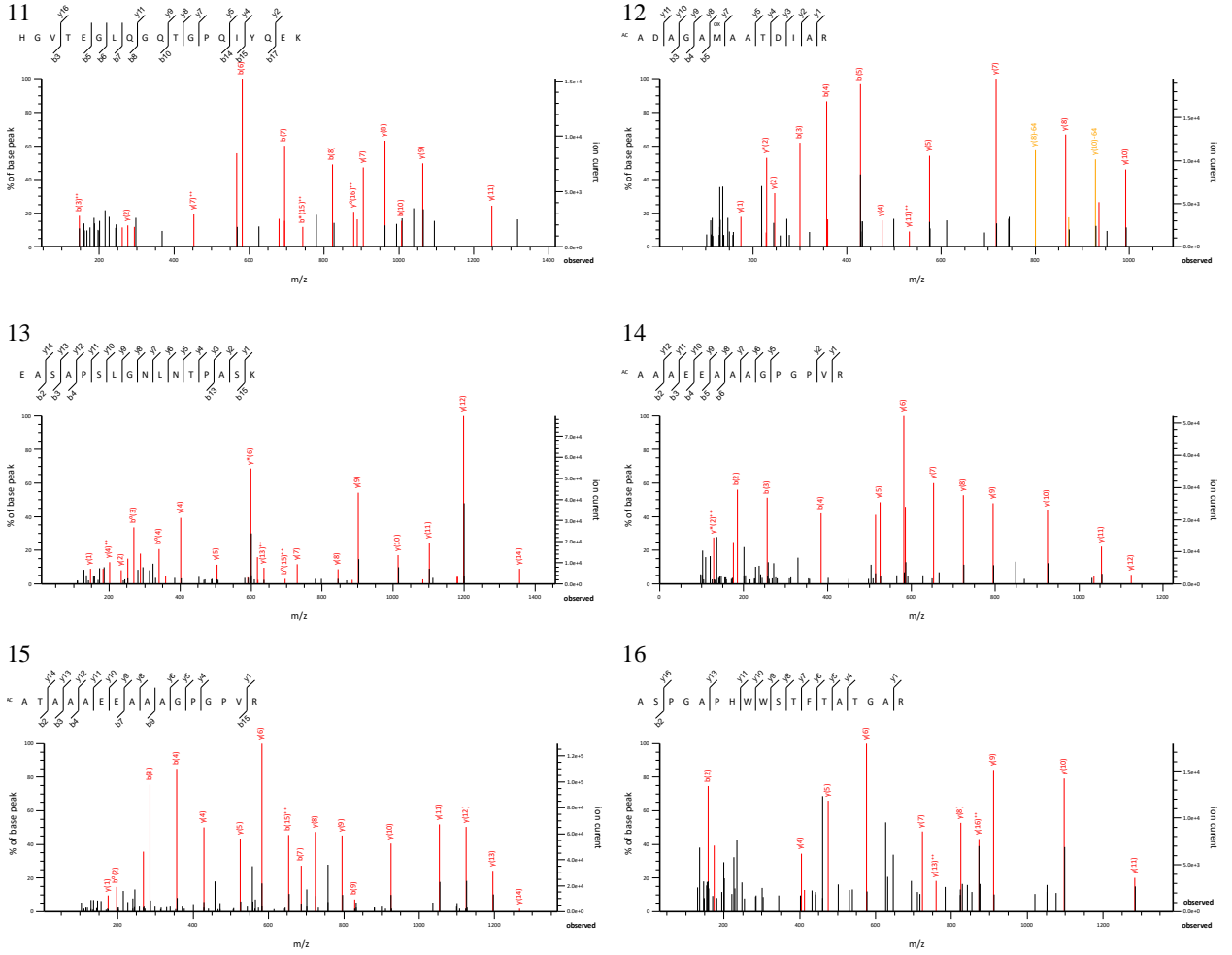
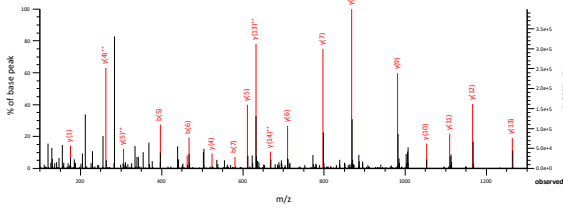
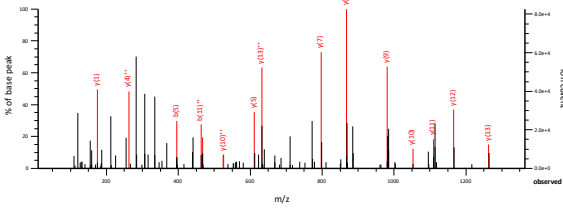


Figure S10. MOLT4 MS/MS spectra corresponding to peptides listed in Table S3.

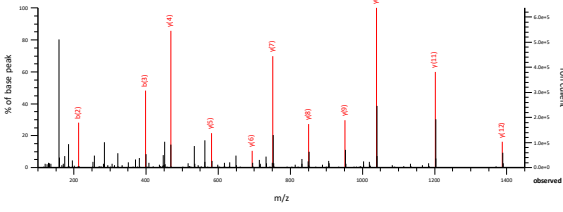
1-K562



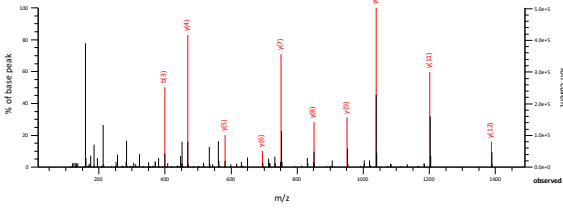
1-MOLT4



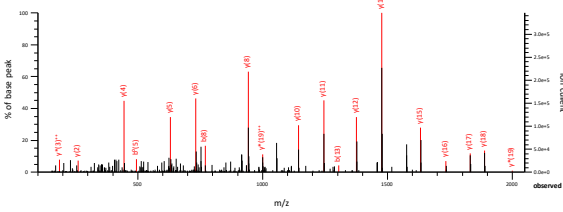
2-K562



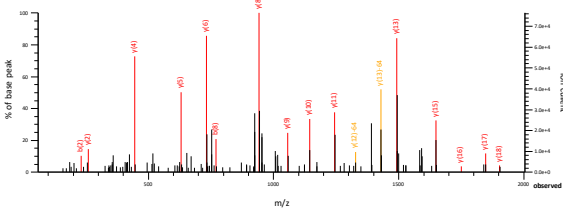
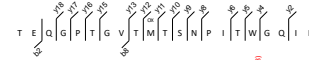
2-MOLT4



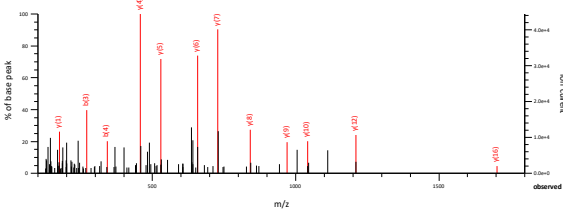
3-K562



3-MOLT4



4-K562



4-MOLT4

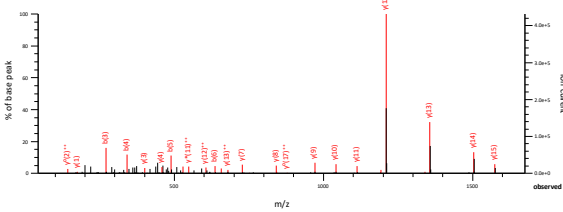


Figure S11. MS/MS spectra corresponding to peptides listed in Table S3 which are found in both K562 and MOLT4 cells.