

Supporting Information

The Discovery of Human sORF-Encoded Polypeptides (SEPs) in Cell Lines and Tissue

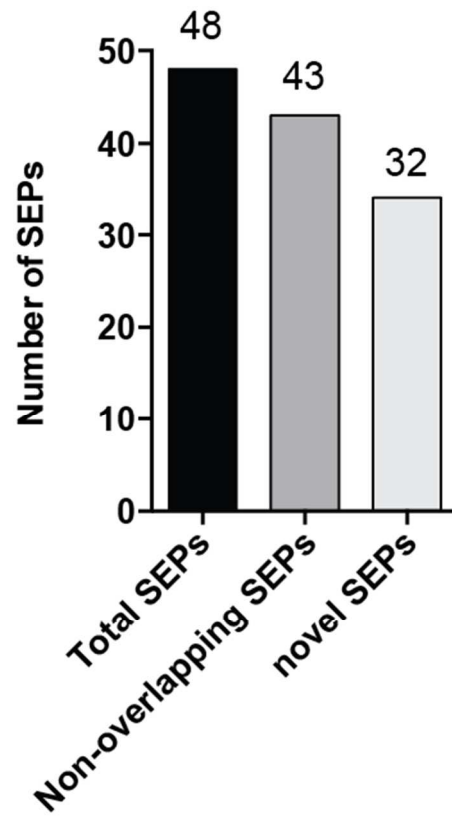
*Jiao Ma¹, Carl Ward¹, Irwin Jungreis^{3,4}, Sarah Slavoff¹, Adam Schwaid¹, John Neveu²,
Bogdan A. Budnik², Manolis Kellis^{3,4} and Alan Saghatelian¹*

1 Department of Chemistry and Chemical Biology, Harvard University, 12 Oxford Street,
Cambridge, Massachusetts 02138, USA

2 MSPRL, Center for Systems Biology, Harvard University, 52 Oxford Street,
Cambridge, Massachusetts 02138, USA

3 MIT Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute
of Technology, Cambridge, Massachusetts 02139, USA

4 The Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02139, USA

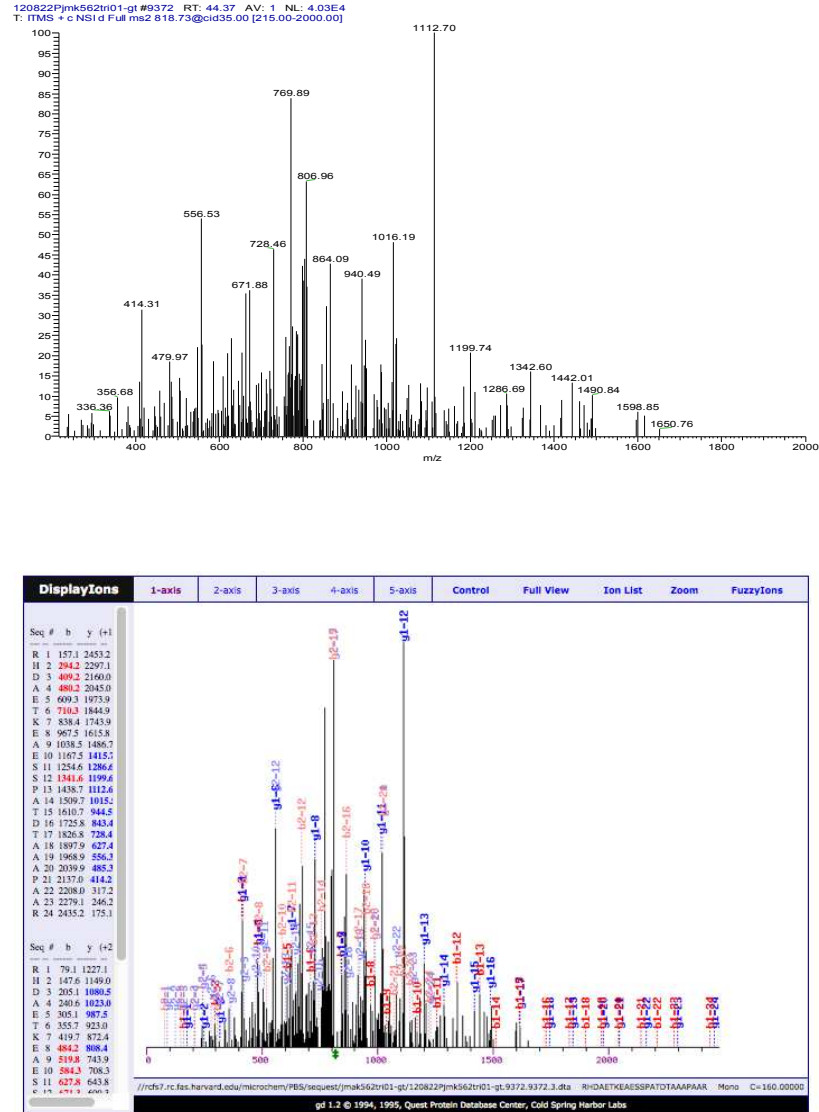


Supplementary Figure 1. Total number of SEPs detected in K562 cells using PAGE + LC-MS/MS workflow after performing an additional six technical replicates.

Supplementary Figure 2. Alignment for ASNSD1-SEP shows protein-coding signature.

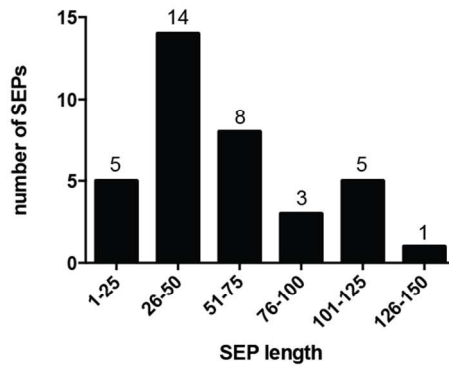
Alignment for ASNSD1-SEP across 29 eutherian mammals, color-coded by CodAlignView ("CodAlignView: a tool for visualizing protein-coding constraint", I Jungreis, M Lin, M Kellis, in preparation). The amino acid sequences of the four tryptic peptides detected, NILDELKK, IVVDELSNLKK, QQQNSNIFFLADR, and EYQEIENLDKTK, are highlighted in yellow. The high concentrations of synonymous substitutions (light green) and conservative amino acid changes (dark green), and relatively low concentrations of radical amino acid changes (red) and frame-shifted regions (orange) is characteristic of protein-coding regions. The region's evolutionary coding potential as measured by per-codon PhyloCSF score, 4.315, is higher than 99.97% of non-coding regions, implying that it has been functional at the amino acid level in much of the eutherian mammal tree.

RHDAETKEAESSPATDTAAAPAAR, sf: .94, PRR3-SEP

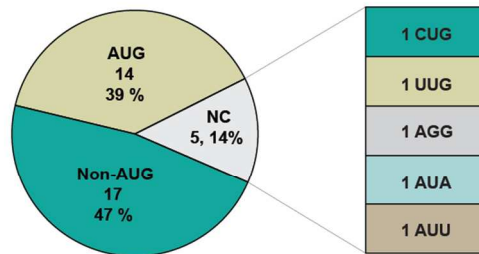


Supplementary Figure 3. MS/MS spectra (raw spectrum: top, SEQUEST annotated spectrum: bottom) for the detected peptide: RHDAETKEAESSPATDTAAAPAAR for PRR3-SEP, with SF score of 0.94.

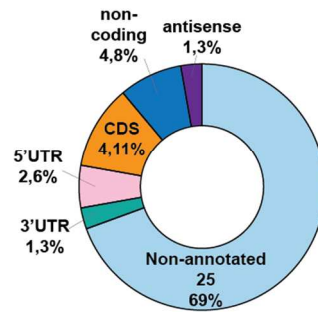
A



B

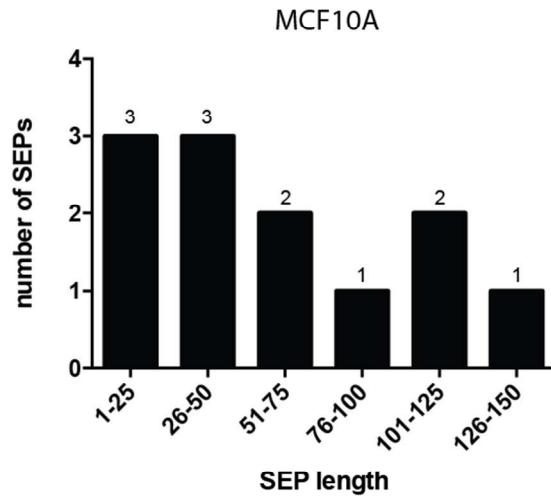


C

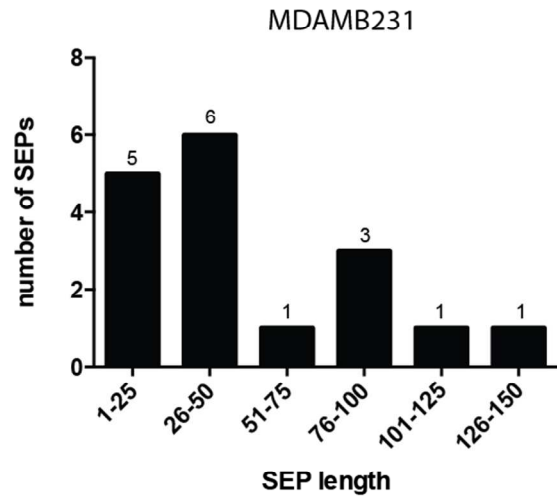


Supplementary Figure 4. **The characteristics of the 36 SEPs in K562 cells validated by Skyline-MRM.** (A) The length distribution of the SEPs, (B) the start codon usage of the SEPs, (C) the SEPs mRNA annotation by RefSeq.

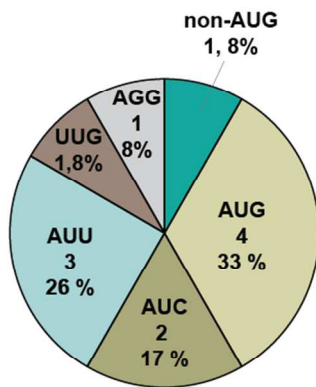
A



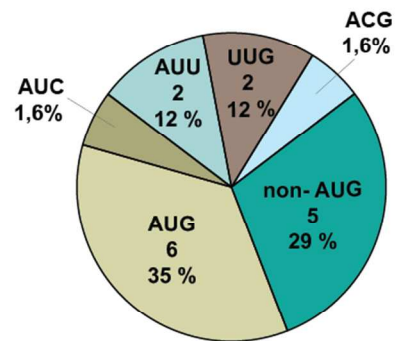
D



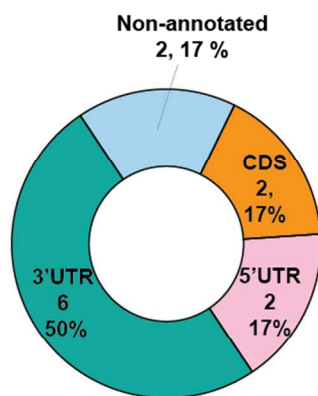
B



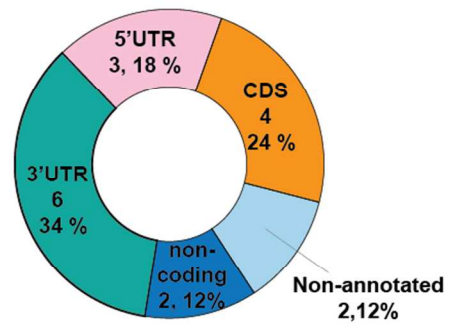
E



C



F



Supplementary Figure 5. The characteristics of SEPs detected in MCF10A and MDAMB231 cell lines. (A) The length distribution of the SEPs, (B) the start codon usage of the SEPs, (C) the SEPs mRNA annotation by RefSeq in MCF10A cells. Similarly (D) The length distribution of the SEPs, (E) the start codon usage of the SEPs, (F) the SEPs mRNA annotation by RefSeq in MDAMB231 cells.

Supplementary Table 1. **A list of 36 SEPs detected in K562 cells that were validated by Skyline-MRM.** The peptides detected by shotgun proteomics were shown in blue, the additional peptides validated in Skyline-MRM were shown in red, the overlapping peptides detected in both shotgun method and Skyline-MRM were highlighted in purple.

Detected peptide by shotgun	Detected peptide by Skyline-MRM	Predicted SEP length	SEP sequence
IVYGDIRK KIVYGDIR	IVYGDIR	41	LYQRREELKLQWRSFNLDKIV YGDIRKEGILNALETMGSSF
ADAPAVSPESPQK	KPPFQPR ADAPAVSPESPQK	101	MSRGGQFGQGQEPLDMFFWV NEISGEITYPPQKADAPAVSPE SPQKPPFQPRSVQEAPCSPQ GPPAQRPALAPPSKPSLKDSG SRNPCPSAPTWARPKPEE
KNFFISK	ALAYIGAR	20	MPKNFFISKIKALAYIGARC
EAESSPATDTAAAPAAR RHDAETKEAESSPATDTAAA PAAR	APSSLPSK EAESSPATDTAAAPAAR	93	MPGRDGGRLAPECGRRRCS PQPSGSGMAGASALLPTDPP EAETESPRAPSSLPSKSRCSHC RRHDAETKEAESSPATDTAAA PAARAGR DWR
FSPLELAAGGVR	FSPLELAAGGVR LLLPCPCCFG	127	MAGGDPAAGRAVGAAKAKR AAGICGDTAAPALSGHPPDCL PESPGRGGLTCSPSQSRFSPLE LAAGGVRGCGAQLVTVGCSE PSWSPGEGAGPGEQQSLDVE AAEALQAGPASPRCRLLLPCP CCFG
YLDLHER	QIFHSNK	33	IRKYLDLHERQIFHSNKHLLNII NCKAQTTRVK
QGSLVQQVALR	EVLEQLMK	66	MMKNQWMEIGHHQLHQPS RFLHLKHPEILPTVQGSLVQQ VALRLEVLLRKEVLEQLMKMIL
RNLVVVINL	DPMNAK TLEVVIS	71	MNVRNVRNLVVVINLFYITG FMSLRDPMNAKSVGRTFVVA INLLYIKDFILVRNPMNVQNV GRTLEVVIS

IYNRLYFLEK	VNFCNK	24	VNFCNKASWDLK KIYNRLYFL EKF
EGGWRQVEGTGTPK	DSAKPIR QVEGTGTPK	75	MTTGER DSAKPIR ATATRQED RSP EGGWRQVEGTGTPK SKQ GSRVLAAQEETQHPEAVPQR ADPKGASASPLRRQ
ELSQYLK	MDELSQYLK	42	MDELSQYLK VILPSTVVLDVILL QLSFLLYIANFLFSLGSLP
SRQVDQEVRSR	QVDQEV	24	SQHFG SRQVDQEVRSR TA WPRW
ENIPDITK	DFVFNLSK	32	DFVFNLSK ILVENRPAFVNEN IPDITKPKHF
VAEIIER	TLYTCLR LVSHGINLALIFSIWK	116	RVAEIIERLVSHGINLALIFSIW K CLKENHFHCRKSFFKYLPR EISLYLPPQAVISCFREWNPPC PSIFWFLGLNSSLVKSPWLGIL SWEQILSCSLMCLHSP TLYTC LLRA
AIVVARVVTIPK	MVAIVVAR	43	MVAIVVARVVTIPK IMHQPVL SFLNFHVPLYTFMSVYVDLSLV
GDFLNLR	AGDFLNLR	41	AGDFLNLRIGISYQFCK FSPINY FFFLFSPCLLYGILLDIS
AGDFLNLR	IGISYQFCK		
GFLAGYVVAK TLRDYLQLLR NQLESLQR RVEDEVNSGVGQDGSLLSSP FLK	NQLESLQR NQLESLQRR GFLAGYVVAK	103	MADDKDSLPLKDLAFL KNQL ESLQRRVEDEVNSGVGQDGS LLSSPFLKGFLAGYVVAK LRAS AVLGFVAVGTCTGIYAAQAYAV PNVEK TLRDYLQLLR KGPD
RLLFAGK IRLLFAGK	GTTFSWVIR	22	GTTFSWVIR LLFAGKLNYS S

GLIENPALIR	LMQEGK QGLIENPALIR	100	QRVQAERLAIRARLKREYLLQY NDPNRQGLIENPALIRWAYAR TTNVYPNFRPTPKNSLMGALC GFGPLIFIYYIIKTERDRKEKLM QEGKLDRTVHLSY
QNIKGLENILQK	GLENILQK VVTTLQSSSENQR	81	IWSRVVTLTLQSSSENQRQNIK GLENILQKEAATCVDNGLFMP LLSVDLVQETCSGDGCEGG MRIDIDTPVSQTCLFITLL
SWLTPVAGK	MAPLGLK	26	MAPLGLKDPLSSWLTPVAGKL VMAVS
HALPLLK	QEFHALPLLK NSTNFFLLIK	52	NSTNFFLLIKQRSFGGFIPIADK RGKDGKCSRFLSFHKQEFHAL PLLKQRKE
GAGILLR	TGAGILLR	44	TGAGILLRWLTHWLLAGSLR SSPGVPLHVLLHGLMMWHEP HSV
KQNSLIANMEK	SGYINR SCLHSIK GEEAAEEK MQIEATR GQTLFSSTK QNSLIANMEK	114	LIKQNSLIANMEKVLVWVM EDQTSJNIPLSQSLIQSKGQTL FSSTKNEKGEEAAEEKFEASRV WLMRFKERSCLSIKMQIEAT RADEEGTASDPEDPAKLIDKS GYINRFTM
MKNFLAVTITGK	TSVQGITTVILK	110	IDWRRKRRKKIEKRKSFRAEV NTKNISPLPHLPPPPPLLLRL QKAVVRVRVTIKKKYKGRKE RKTSVQGITTVILKRRSLRRES

			FMKNFLAVTITGKPRKSPGS
KDLHLSWEPK	FEFFPK NGLPSVLLVK	43	KQQPPLFSLYKFEFFPKLDLH LSWEPKEKNGLPSVLLVKEIL
KNEFLK	DHVLFFK	27	IIFKNEFLKDHVLFFKSIFSSYF CYC
AEIILK	VFDLQDF MAEIIILK	17	MAEIIILKAKVFDLQDF
TPLLAYIQ	TPLLAYIQPDTSAF	49	MNLEMEKKAGLFQRVDLSEL DSTIELCCIFCGSSKTPLLAYIQP DTSAF
HAFLNLR	HALFLNLR NLQTPGAVGEDK	54	HAFLNLRRAIPSPQSNLNERPQ VQLLHSPDLLSTPRNLQTPG AVGEDKKGSGVA
EVEGAVSR	QSEVMSQK IFNNHTLIK	42	TQVEGAVSRDCITALQPGKQ SEVMSQKQTTKIFNNHTLIK
RKPLYTIGWNL	DFTSHQLER	64	SSGKGKNSQRDFTSHQLERL SSKRQNIKRVGKNAEKRKPLY TIGWNLNWYSHYKKQHGGSS KN
KINALLK	GNILLSNK	50	SQPPLKCLCIKINALLKGNILLS NKCVCVYHTSILRKCWTSEY HKTGN
FQPPHHVQSSPDVK	GLSFQPPHHVQSSPDVK	32	ESCEPTEQKGLSFQPPHHVQ SSPDVKSQFWF
ARDQYGHLIPTK	KPSFSPR DQYGHLIPTK GSCHFLSQVGGWGI	61	MCAEIEEGAEGVTARDQYGH LIPTKVASGPQGLSGARKPSFP SPRLRGSCHFLSQVGGWGI
LAFIFLPDR	NDLAFIFLPDR	65	AKIVPLHSSLGDRVRPCLKTKQ TKEFRNDLAFIFLPDRQCIHQD GTLTGNQVLAPLLAGKEHEVF

Supplementary Table 2. A list of all novel SEPs detected in this study.