

Supplemental information

**A pan-cancer transcriptome analysis of exitron
splicing identifies novel cancer driver
genes and neoepitopes**

Ting-You Wang, Qi Liu, Yanan Ren, Sk. Kayum Alam, Li Wang, Zhu Zhu, Luke H. Hoeppner, Scott M. Dehm, Qi Cao, and Rendong Yang

A

12 EIS events confirmed in SKBR3
RNA-Seq data were selected for validation

**B**

$$\text{PCC}=0.88; p=0.02$$

SETDB2 V189Ifs*63

NOD1 G195Afs*56

MOGS V277Afs*19

MICALL2 V316Pfs*98

SSFA2 Δ593-761

CHD2 G1575Dfs*8

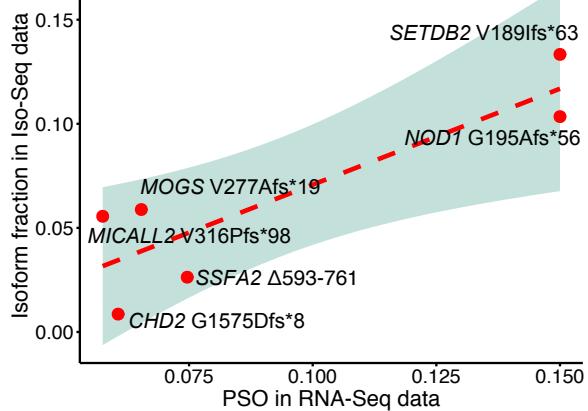
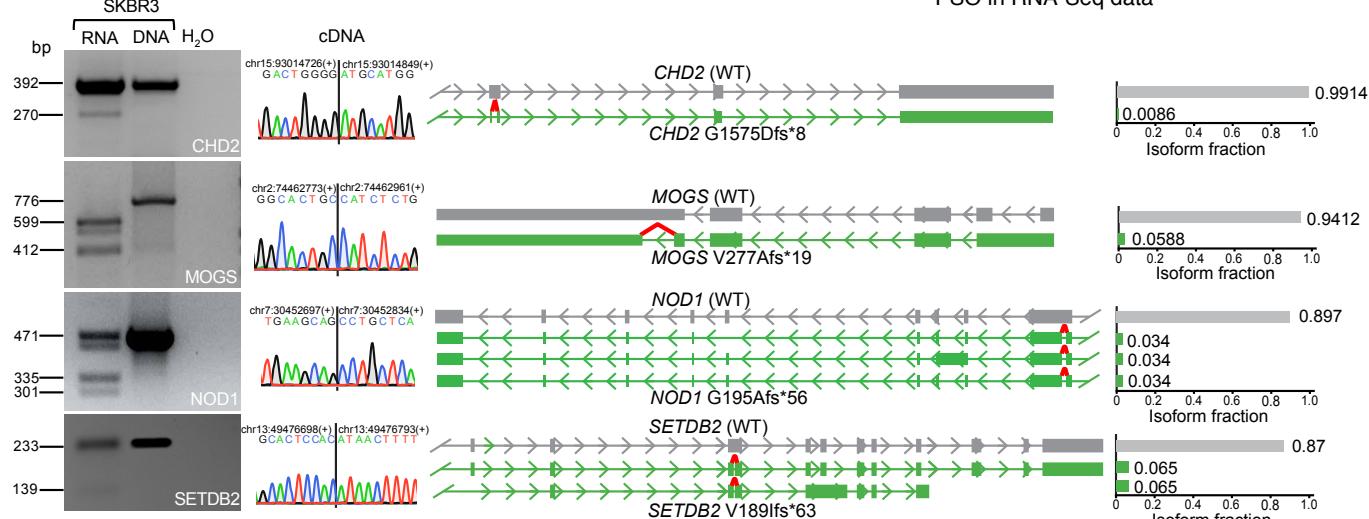
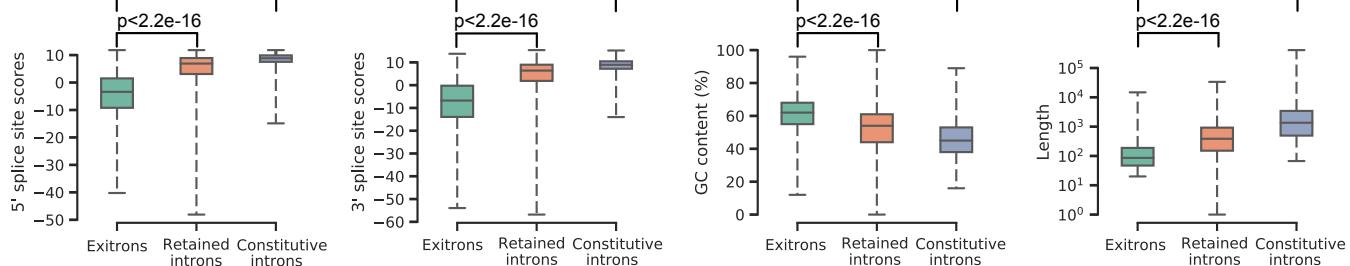
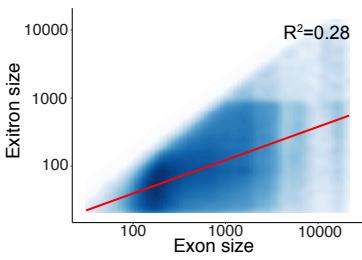
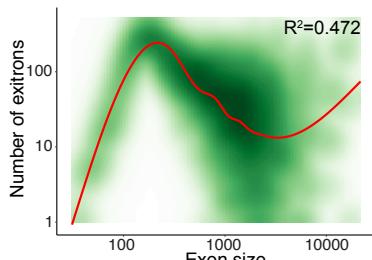
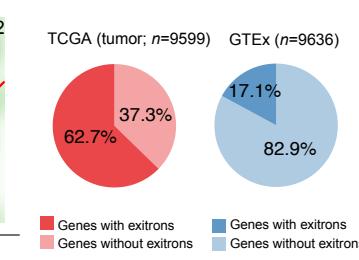
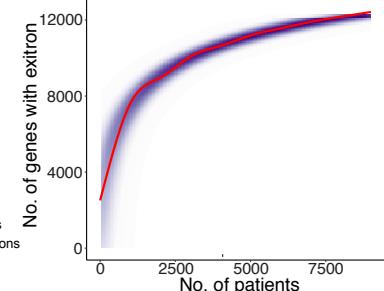
**C****D****E****F****G****H**

Figure S1. Exitron splicing (EIS) events detection in cancer, Related to Figure 1.

- (A) Information flow during the validation of EIS events identified in TCGA-BRCA.
- (B) The fractions of exitron-spliced isoforms from PacBio long read Iso-Seq data are highly correlated with their PSO values predicted from the RNA-Seq data. Pearson's correlation coefficient (PCC) was used to determine the correlation.
- (C) Four EIS events (*CHD2* G1575Dfs*8, *MOGS* V277Afs*19, *NOD1* G195Afs*56 and *SETDB2* V189Ifs*63) were validated by RT-PCR and Sanger sequencing. Isoform fraction of wild-type (grey color) and exitron-spliced transcripts (green color) are counted according to long read Iso-Seq.
- (D) Comparisons of 5' and 3' splice site strength (measured using maximum entropy; see Methods), GC content and size distribution of exitrons, retained introns and constitutively spliced introns. All p values were calculated by Mann-Whitney U test.
- (E) Exitron size and exitron-containing exon size are correlated ($R^2=0.28$). Exitron size ranges from 21bp to 14,765bp and exon size ranges 31bp to 21,693bp. Both exon size and exitron size are shown on log10 scale. All the points were fitted by a linear regression red line.
- (F) Distribution of 129,406 exitrons in TCGA based on their locations in exon of different sizes. Medium-sized exons contain the largest number of exitrons. All the points were fitted by a generalized additive model.
- (G) The proportion of genes with and without exitrons in TCGA tumor and GTEx healthy samples.
- (H) Saturation analysis of exitron splicing events in TCGA by adding more samples. Each point is a random subset of samples. All the points were fitted by a smoothed red line.

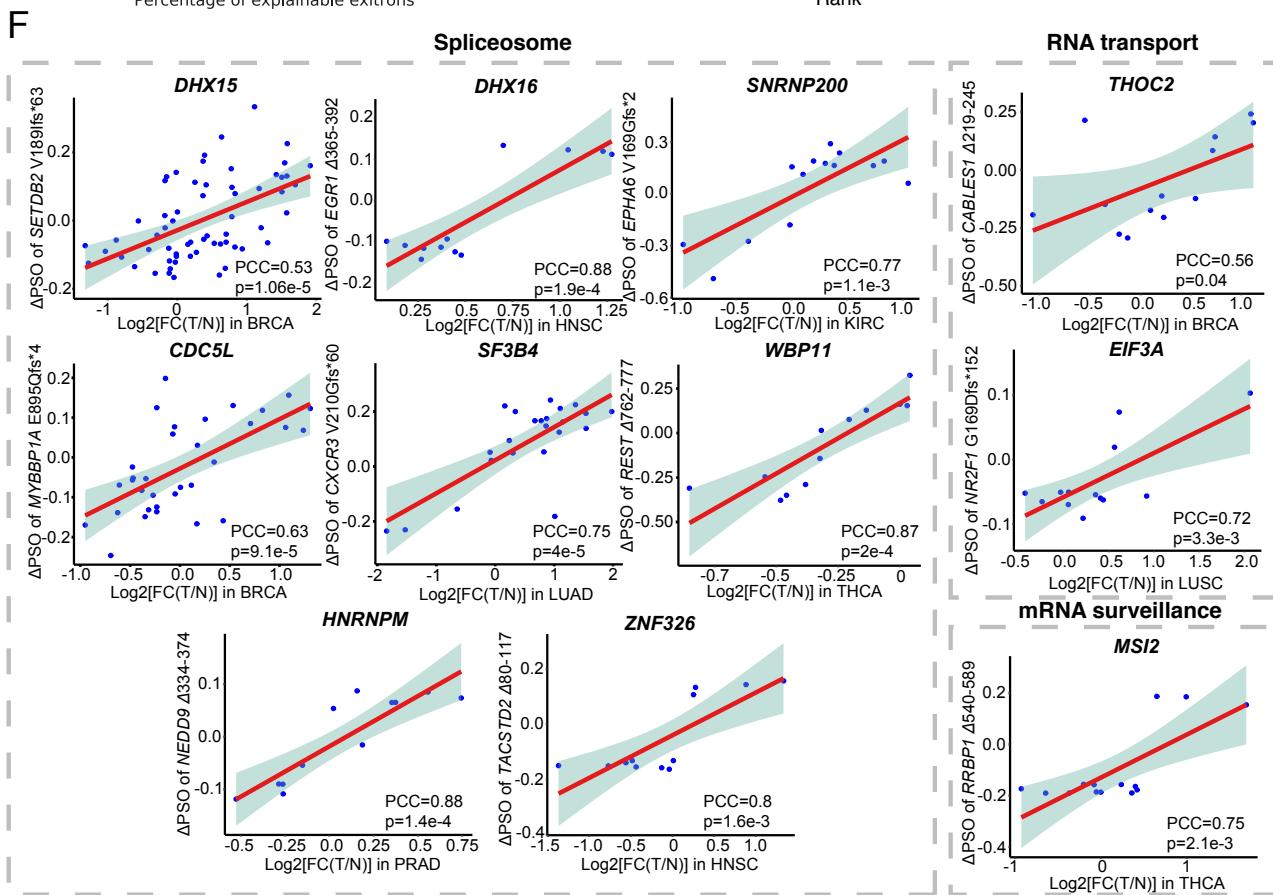
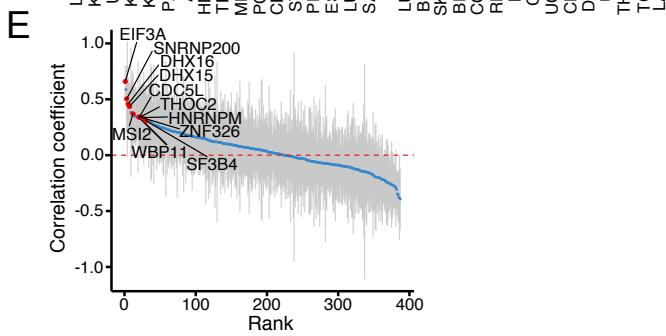
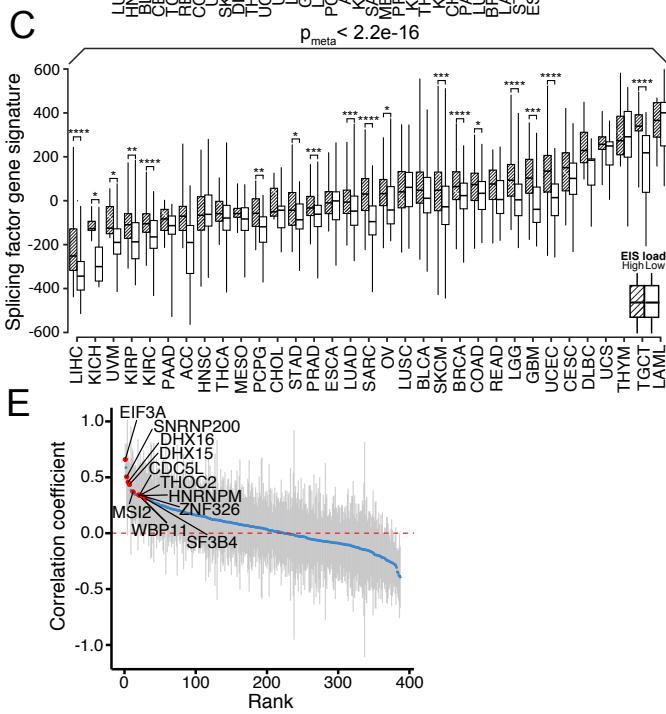
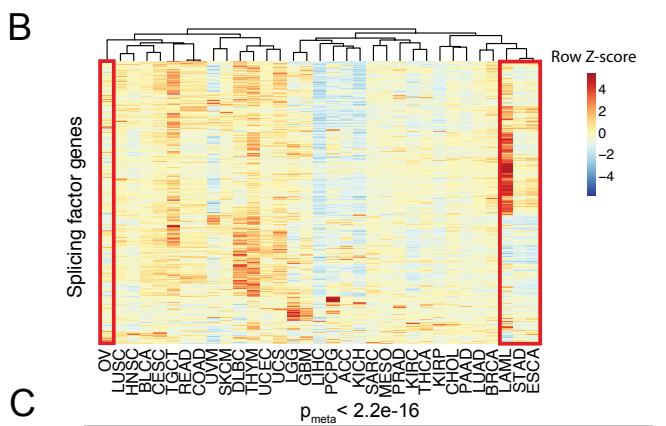
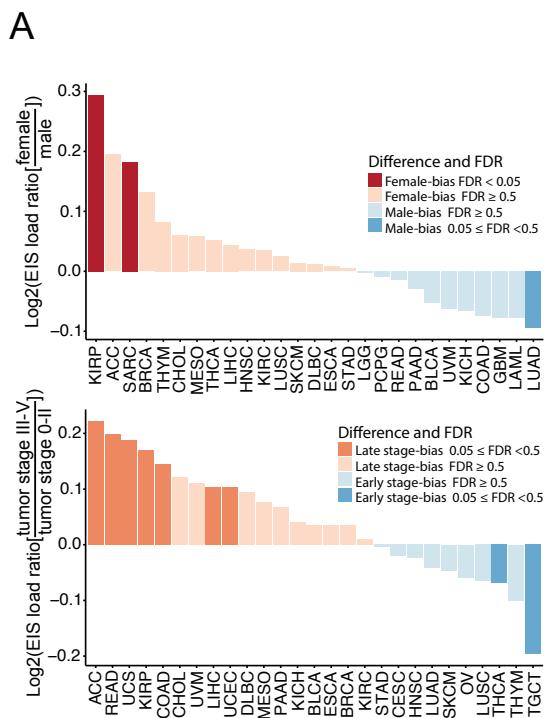


Figure S2. EIS dysregulation related to gender, tumor stage and splicing factors, Related to Figure 2.

(A) (top panel) Association between exitron splicing load and gender difference. (bottom panel) Association between exitron splicing load and tumor stage difference.

(B) Cancer types were clustered by the median expression of 404 splicing factor genes.

(C) The comparison of gene expression signature (measured by Z-score) of splicing factor genes on the top (high) vs. bottom (low) quartile patients by exitron splicing (EIS) load across 33 cancer types. p-values were calculated by Mann-Whitney U test with Bonferroni correction. (* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, **** $p < 0.0001$).

(D) The percentage of exitron splicing events in each cancer type that can be explained by a Generalized Additive Model (GAM) of differential expression of splicing factor between tumor and normal samples.

(E) Splicing factor genes ranked by Pearson correlation coefficient of GAM-derived exitron/splicing-factor pairs. Grey bar indicates the 95% confidence interval. Genes of interest are highlighted.

(F) Linear regression indicates a positive correlation of exitron splicing changes (Δ PSO) and differential mRNA expression ($\text{Log}_2(\text{FC})$) of corresponding splicing factors involving spliceosome, RNA transport and mRNA surveillance according to KEGG annotation. FC, Fold change; T, tumor; N, normal.

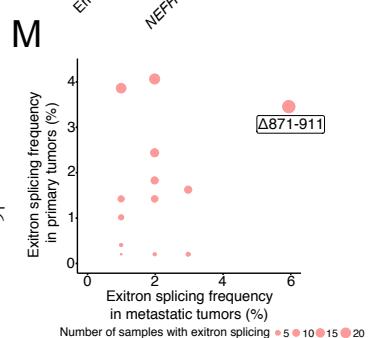
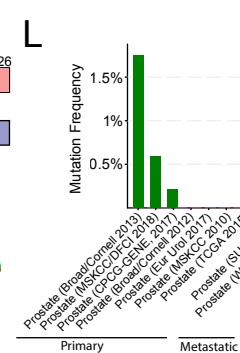
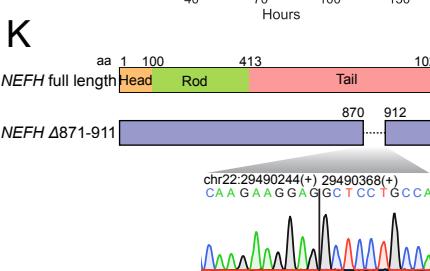
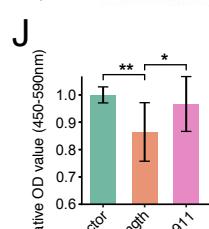
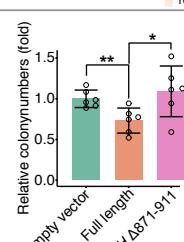
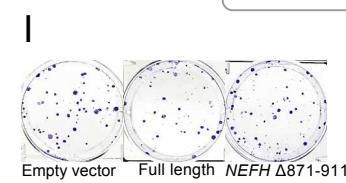
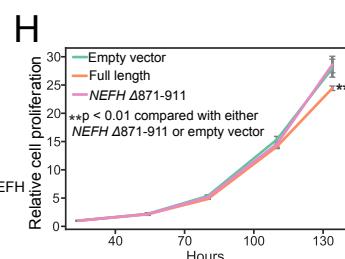
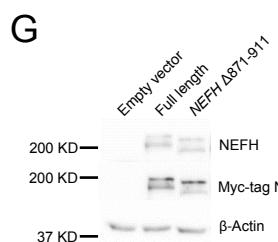
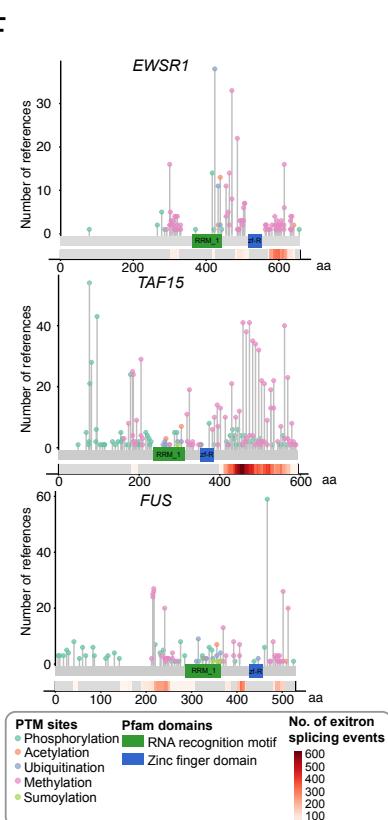
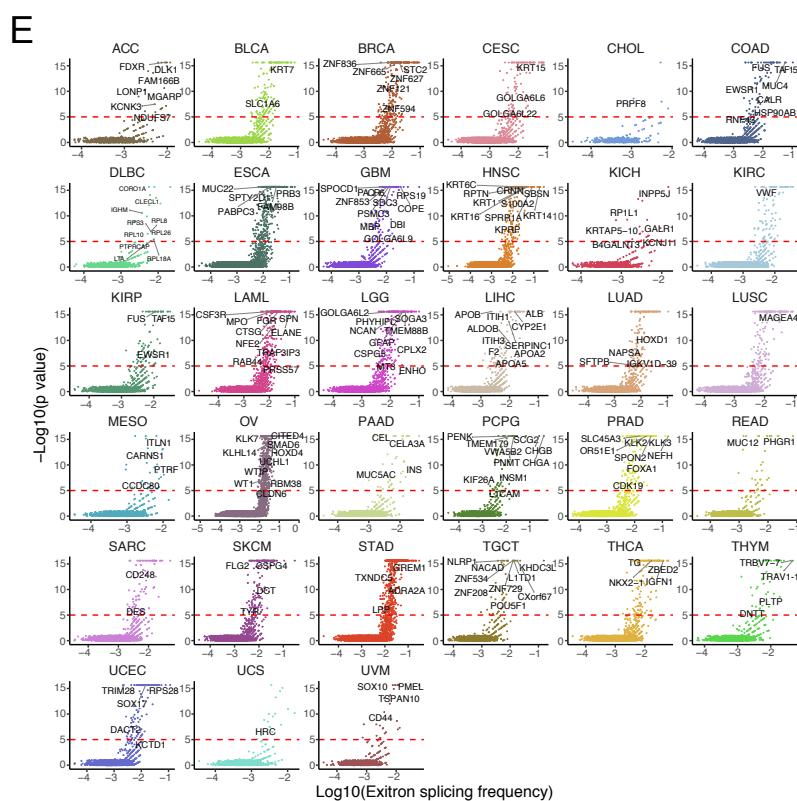
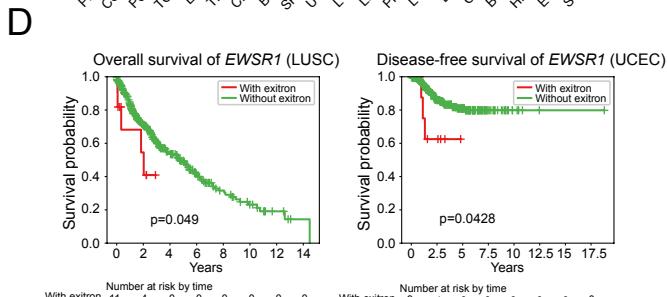
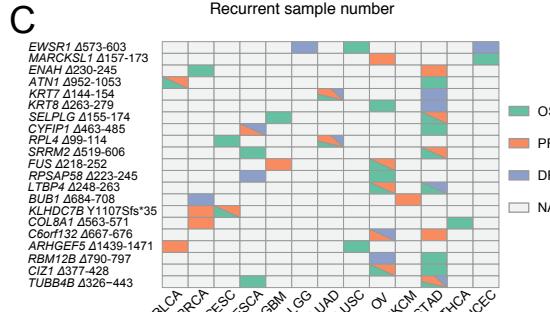
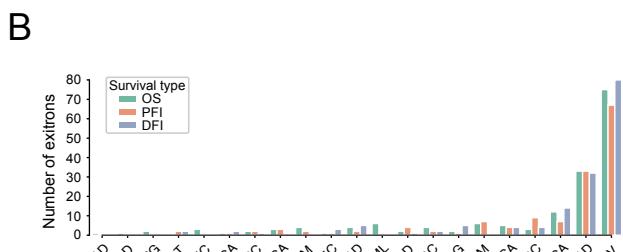
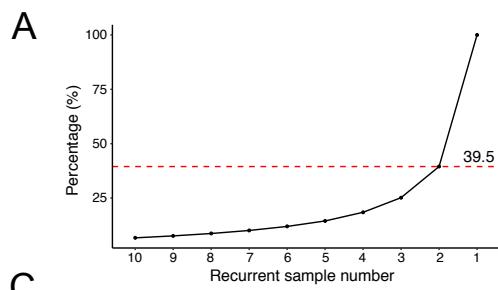


Figure S3. Clinical and functional impacts of tumor-specific exitrons, Related to Figure 3.

- (A) The percentage of recurrently spliced TSEs at varied sample numbers.
- (B) The distributions of clinically informative TSE splicing events in different cancer types. OS stands for overall survival, PFI stands for progression-free interval, and DFI stands for disease-free interval.
- (C) The clinical relevance of fourteen clinically informative TSE splicing events were identified in multiple cancer types. The colored box indicates the association with the different survival types.
- (D) Kaplan-Meier curves demonstrate clinical relevance of clinically informative TSE splicing event at *EWSR1* Δ573-603 in LUSC and UCEC. The p value was calculated by a log-rank test.
- (E) SEGs identified across TCGA cohorts. For COAD and KIRP, SEGs overlapping with the COSMIC cancer census set are highlighted. For the remaining cancer cohorts, tissue-specific SEGs are highlighted. Coloring is according to cancer type.
- (F) Exitron splicing hotspots in FET family proteins are enriched in post-translational modification (PTM) sites. Any position with a PTM site obtains a circle. Circles are colored with respect to the corresponding PTM types, and the height of the line depends on the number of references supporting the PTM site at the specified position. The grey bar represents the entire protein with the different amino acid positions (aa). The colored boxes are specific Pfam domains.
- (G) Western blotting detected overexpression of Myc-tagged wild-type and exitron-spliced NEFH in PC-3 cells.
- (H) Cell proliferation is inhibited by overexpressing wild-type NEFH in PC-3 cells.
- (I) Colony formation assays in PC-3 cells with overexpression of wild-type and exitron-spliced NEFH. The number of colonies were quantified and shown at the right panel.
- (J) BrdU ELISA assay of PC-3 cells overexpressing wild-type and exitron-spliced NEFH with 24 hrs of BrdU label (n = 10). Y axis, absorbance of 450-590 nm relative to empty vector. There was less incorporation of BrdU in cells expressing wild-type NEFH. All p-values are calculated using unpaired, two-tailed Student's t-test. Error bars indicate ± s.d. *p < 0.05, **p < 0.01.
- (K) Graphical depiction of NEFH protein domains annotated by Uniprot and the location of NEFH Δ871-911 exitron. Electropherogram illustrates the splicing junction of NEFH Δ871-911. aa, amino acid.
- (L) Mutation frequency of NEFH among prostate cancer studies from cBioPortal (<https://www.cbioportal.org/>) and the literature (Quigley et al., 2018 for WC-SU2C).
- (M) *NEFH* Δ871-911 is the most frequent TSE splicing commonly present in primary (TCGA-PRAD, n=492 samples) and metastatic prostate tumors (WC-SU2C, n=101 samples, dbGap: phs001648.v1.p1). Dot color indicates the protein domain (Tail) where TSE splicing events occur.

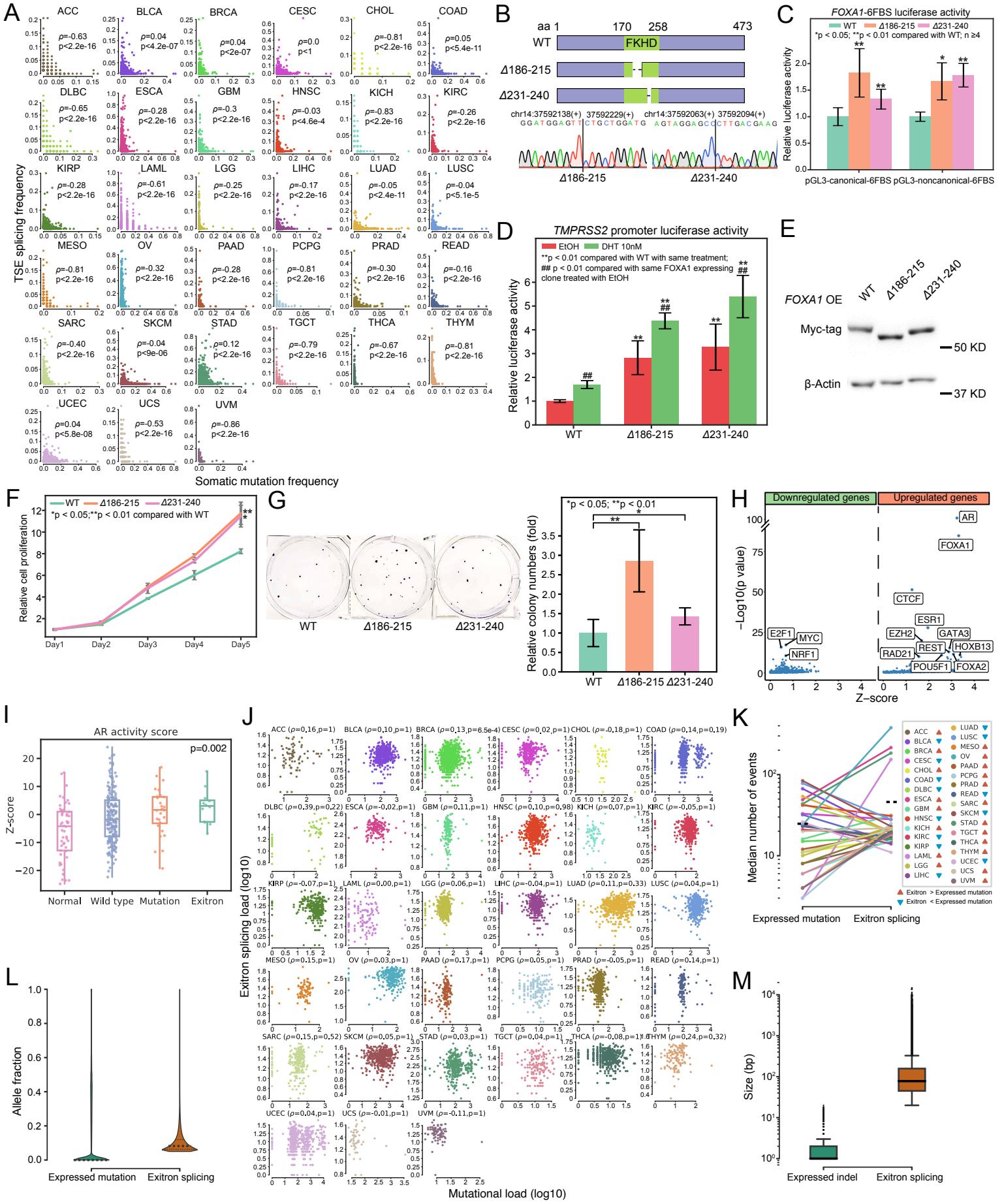


Figure S4. Mutual exclusivity between tumor-specific exitrons and somatic mutations, Related to Figure 4.

- (A) The frequency of somatic mutation and exitrons splicing occurred in genes are inversely correlated across TCGA cohorts. Each dot corresponds to a gene. Coloring is according to cancer type. The title of each sub-panel stats the Spearman correlation coefficient. All p values have been corrected for multiple-testing using Bonferroni's method.
- (B) Graphical depiction of FOXA1 protein Forkhead DNA-binding domain (FKHD) and locations of two FKHD exitrons. Electropherogram illustrates splicing junctions of FOXA1 exitrons shown in the upper diagram.
- (C) FOXA1-luciferase reporter assay with results normalized to level of wild-type (WT) FOXA1 activity.
- (D) TMPRSS2 promoter luciferase reporter activity with overexpression of WT or various FOXA1 mutants and DHT stimulation.
- (E) Representative western blot against C4-2 stable cell lines expressing Myc-tagged WT and exitron-spliced FOXA1.
- (F) CellTiter-Glo growth assays indicate that overexpression of exitron-spliced FOXA1 significantly promote growth in prostate cancer cells in standard media conditions.
- (G) Colony-formation ability of prostate cancer cell lines overexpressing exitron-spliced FOXA1 is increased relative to cells overexpressing WT FOXA1. Cells were fixed and stained with crystal violet. Assays were done at least thrice and in triplicate wells. The figure is a representative of three experiments with similar results. Quantification was performed by manual counting. All p values are relative to WT, calculated using unpaired, two-tailed Student's t-test. Error bars indicate \pm s.d. * $p < 0.05$, ** $p < 0.01$.
- (H) BART prediction of specific transcription factors mediating observed transcriptional changes. The significant and strong (z-score) mediators of transcriptional responses in patients with FOXA1 exitron splicing events are labelled (BART, Wilcoxon rank-sum test).
- (I) Extent of AR pathway activation in FOXA1 exitron spliced, mutated, wild-type tumors and adjacent normal prostate tissues in TCGA study. AR score was calculated using established gene signatures. P value was calculated by Kruskal-Wallis test.
- (J) Correlation between the number of tumor-specific exitron splicing events per sample and its corresponding somatic mutational load across TCGA cohorts. Each dot corresponds to a sample. Coloring is according to cancer type. Both exitron splicing load and mutational load are shown on log10 scale. The title of each sub-panel stats the Spearman correlation coefficient. All p values have been corrected for multiple-testing using Bonferroni's method.
- (K) The comparison on median number of expressed mutations and tumor-specific exitron splicing events per sample in each TCGA cohort.
- (L) The comparison on allele fraction of expressed mutations and exitron splicing events. Allele fractions of expressed mutation and exitron splicing were measured by VAF and PSO respectively.
- (M) The comparison on size between expressed indels (inframe and frameshift) and exitron splicing events.

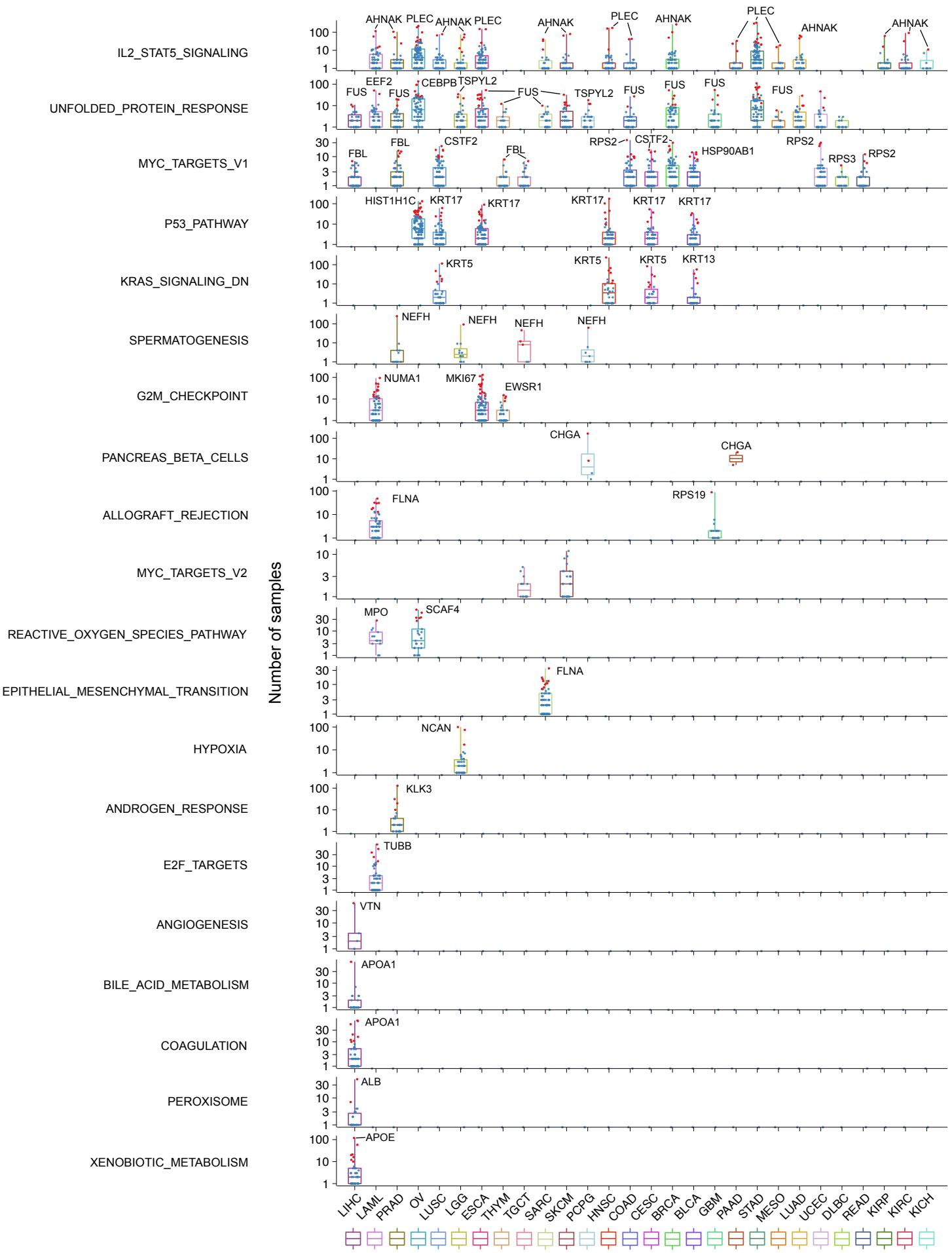


Figure S5. TSE splicing enriched hallmark gene sets identified across TCGA cohorts, Related to Figure 5.
 SEGs in each tumor cohort and gene set are colored in red. The SEG with the highest number of exitron splicing altered samples in each gene set is highlighted.

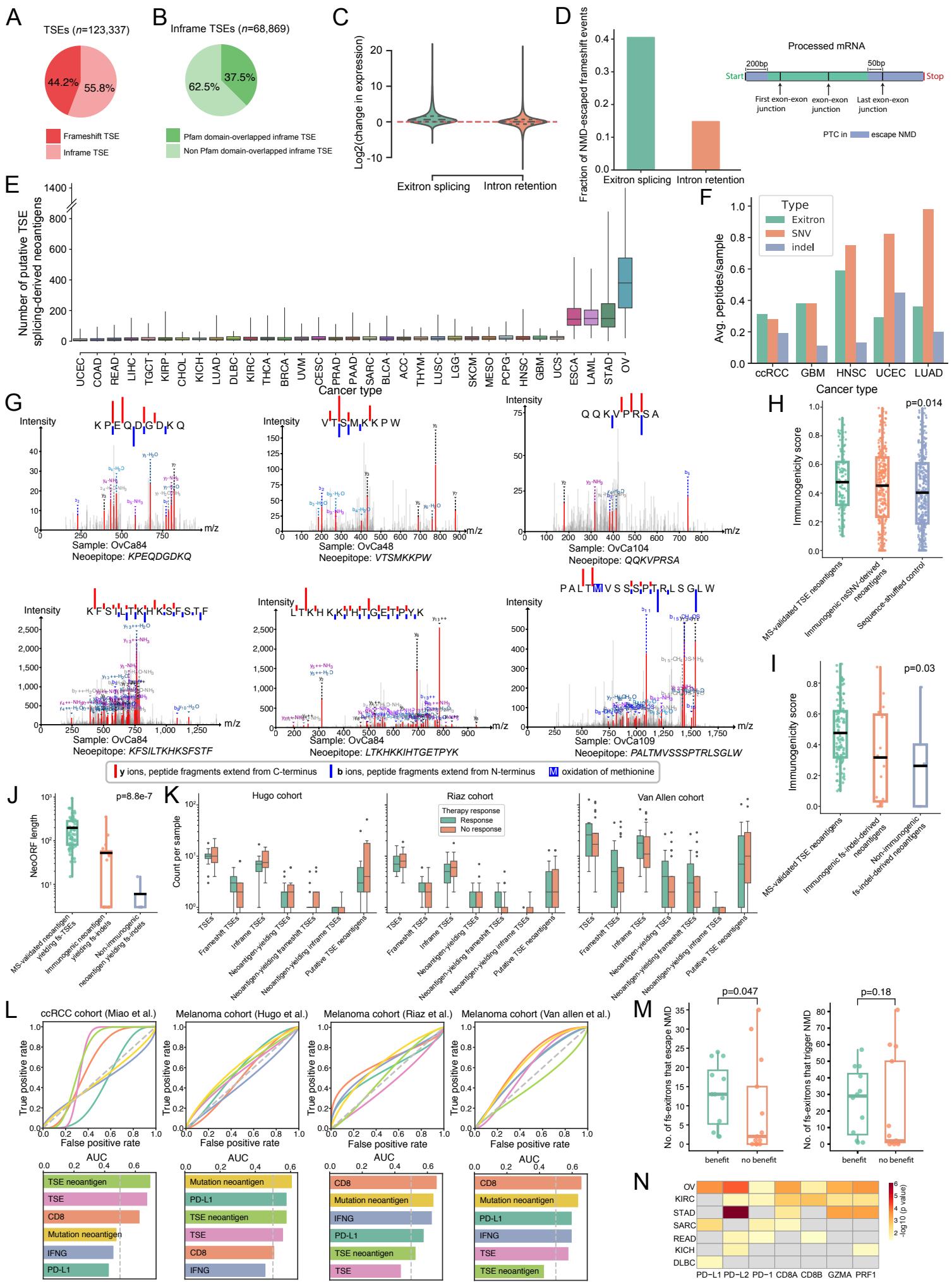


Figure S6. Tumor-specific exitrons (TSEs) represent a source of neoepitopes, Related to Figure 6.

- (A) The proportion of frameshift and inframe TSEs in TCGA tumors.
- (B) The proportion of inframe TSEs overlapping with Pfam domains.
- (C) RNA-Seq expression changes in alternatively-spliced versus wild-type samples. The extent of NMD by exitron splicing and intron retention is estimated by comparing the mRNA expression level in samples with the splicing event to the median mRNA expression level of the same gene across all other wild-type tumor samples.
- (D) Proportion of NMD-escape events between frameshift exitron splicing events and intron retention. Rules of NMD evasion is determined based on the premature termination codon (PTC) position according to the literature (Lindeboom et al., 2016).
- (E) TSE neoantigen counts by cancer type.
- (F) Comparison of the contribution of TSE splicing, somatic SNVs and indels to proteomic-confirmed putative neoantigens across five cancer types offered by the CPTAC3 study.
- (G) Mass spectra show TSE neoepitopes bound to (left panel) MHC class I molecules and (right panel) MHC class II molecules in Schuster et al. ovarian cancer cohort.
- (H) Comparison of the immunogenicity scores of mass spectrometry (MS)-validated TSE neoantigens from CPTAC with nsSNV-derived neoepitopes inducing T-cell responses collected from dbPepNeo database and corresponding sequence-shuffled control.
- (I) Comparison of the immunogenicity scores of MS-validated TSE neoantigens from CPTAC with functionally validated immunogenic frameshift (fs) indels-derived neoantigens and non-immunogenic indel-derived neoantigens collected in the literature (Litchfield et al., 2020).
- (J) Comparison of neoORF length between MS-validated frameshift TSEs and frameshift indels. P values were calculated by Kruskal-Wallis test. Mean values are highlighted in black.
- (K) Association of TSEs, frameshift TSEs, inframe TSEs, neoantigen-yielding TSEs, neoantigen-yielding frameshift TSEs, neoantigen-yielding inframe TSEs and putative TSE neoantigens with response to checkpoint inhibitor therapy in Hugo ($n=14$ response, $n=13$ no response), Riaz ($n=10$ response, $n=23$ no response) and Van Allen ($n=13$ response, $n=21$ no response) melanoma patient cohorts. All comparisons show a non-significant trend ($p > 0.05$, Mann-Whitney rank test). Boxplots show the median, first and third quartiles, whiskers extend to $1.5 \times$ the interquartile range, and outliers are plotted individually.
- (L) ROC curves for the performance and the area under the ROC curve (AUC) of the TSE neoantigen load, TSE splicing load, mutation neoantigen load, and immune signatures of CD8, PD-L1 and interferon gamma (IFNG) in predicting checkpoint inhibitor response in ccRCC and melanoma cohorts. Grey dash line indicates random predictions.
- (M) Efficacy of immunotherapy is predicted by the burden of NMD-escaped frameshift (fs) exitrons (left panel) but not NMD-triggered fs-exitrons (right panel) in clear cell renal cell carcinoma (ccRCC). P values were calculated using Mann-Whitney rank test.
- (N) Comparison of gene expression of PD-L1, PD-L2, PD-1, CD8A, CD8B, GZMA and PRF1 between the top and bottom quartile patients based on their TSE splicing-derived neoantigen load across TCGA cancer types. Cancer types with at least one significant p value for the seven gene markers are showed. TSE splicing events with PSO ≥ 0.1 are used for comparison. The p values were calculated by a Mann-Whitney rank test.

Table S1. Cancer abbreviation in TCGA cohort, Related to Figures 1, 2, 3 and 6.

Abbreviation	Study name
ACC	Adrenocortical carcinoma
BLCA	Bladder Urothelial Carcinoma
BRCA	Breast invasive carcinoma
CESC	Cervical squamous cell carcinoma and endocervical adenocarcinoma
CHOL	Cholangiocarcinoma
COAD	Colon adenocarcinoma
DLBC	Lymphoid Neoplasm Diffuse Large B-cell Lymphoma
ESCA	Esophageal carcinoma
GBM	Glioblastoma multiforme
HNSC	Head and Neck squamous cell carcinoma
KICH	Kidney Chromophobe
KIRC	Kidney renal clear cell carcinoma
KIRP	Kidney renal papillary cell carcinoma
LAML	Acute Myeloid Leukemia
LGG	Brain Lower Grade Glioma
LIHC	Liver hepatocellular carcinoma
LUAD	Lung adenocarcinoma
LUSC	Lung squamous cell carcinoma
MESO	Mesothelioma
OV	Ovarian serous cystadenocarcinoma
PAAD	Pancreatic adenocarcinoma
PCPG	Pheochromocytoma and Paraganglioma
PRAD	Prostate adenocarcinoma
READ	Rectum adenocarcinoma
SARC	Sarcoma
SKCM	Skin Cutaneous Melanoma
STAD	Stomach adenocarcinoma
TGCT	Testicular Germ Cell Tumors
THCA	Thyroid carcinoma
THYM	Thymoma
UCEC	Uterine Corpus Endometrial Carcinoma
UCS	Uterine Carcinosarcoma
UVM	Uveal Melanoma

Table S4. CPTAC-confirmed neoantigens derived from exitrons, Related to Figure 6.

Neoantigen	Gene	Chr	Start	End	Sample	Cohort
MTRAAPPLM	ARVCF	chr22	19981341	19981603	TCGA-61-1907-01A-01R-1567-13	OV
RAAWPLMTR	ARVCF	chr22	19981341	19981603	TCGA-61-1907-01A-01R-1567-13	OV
HGRAESAPR	DTX2	chr7	76482754	76482816	TCGA-61-1907-01A-01R-1567-13	OV
MAWRMTRA	ARVCF	chr22	19981341	19981603	TCGA-61-1907-01A-01R-1567-13	OV
RMTRAAPL	ARVCF	chr22	19981341	19981603	TCGA-61-1907-01A-01R-1567-13	OV
MAVPLLRRM	CTAG2	chrX	154653345	154653496	TCGA-61-1741-01A-02R-1567-13	OV
WMMRKTWMK	CCDC34	chr11	27362975	27363120	TCGA-29-1697-01A-01R-1567-13	OV
SLKVECMPK	CDH11	chr16	64982087	64982256	TCGA-29-1774-01A-01R-1567-13	OV
MTMDMGRFK	KHDRBS1	chr1	32037905	32037990	TCGA-29-1774-01A-01R-1567-13	OV
QSSWPCTFK	ZNF219	chr14	21090795	21090835	TCGA-13-1492-01A-01R-1565-13	OV
FQCAFLHRL	LRRK1	chr15	101065771	101065998	TCGA-24-1435-01A-01R-1566-13	OV
FYFGGPQYL	MYCN	chr2	15942179	15942814	TCGA-24-1435-01A-01R-1566-13	OV
RAAAAAGVR	AEBP2	chr12	19439858	19439919	TCGA-25-1316-01A-01R-1565-13	OV
AAAAGVRRR	AEBP2	chr12	19439858	19439919	TCGA-25-1316-01A-01R-1565-13	OV
AAAAAGVRR	AEBP2	chr12	19439858	19439919	TCGA-25-1316-01A-01R-1565-13	OV
KMRAAAAG	AEBP2	chr12	19439858	19439919	TCGA-25-1316-01A-01R-1565-13	OV
EMAPSWVLK	MYL6	chr12	56159988	56160073	TCGA-25-1319-01A-01R-1565-13	OV
KEMAPSWVL	MYL6	chr12	56159988	56160073	TCGA-25-1319-01A-01R-1565-13	OV
ERRRSLSNF	NCL	chr2	231455442	231455524	TCGA-25-2399-01A-01R-1569-13	OV
SAFPFPFER	NCL	chr2	231455442	231455524	TCGA-25-2399-01A-01R-1569-13	OV
SFCPSAFPF	NCL	chr2	231455442	231455524	TCGA-25-2399-01A-01R-1569-13	OV
SSRSWMLVR	EMILIN2	chr18	2892048	2892154	TCGA-25-2399-01A-01R-1569-13	OV
MRCWRTTRM	ZBTB46	chr20	63747061	63747212	TCGA-59-2351-01A-01R-1569-13	OV
ATPYSFSL	VWA5A	chr11	124123723	124123790	TCGA-61-1914-01A-01R-1567-13	OV
LMQATPYSF	VWA5A	chr11	124123723	124123790	TCGA-61-1914-01A-01R-1567-13	OV
MQATPYSFL	VWA5A	chr11	124123723	124123790	TCGA-61-1914-01A-01R-1567-13	OV
IQTTPPHDTL	HNRNPD	chr4	82373506	82373588	TCGA-AR-A0TT-01A-31R-A084-07	BRCA
TMWKNHYQK	WWC2	chr4	183261269	183261363	TCGA-A2-A0T7-01A-21R-A084-07	BRCA

Table S5. MS-based imunopeptidomics-confirmed neoantigens derived from exitrons, Related to Figure 6.

Neoantigen	Gene	Chr	Start	End	Sample	MHC type
WGRLLSEY	BTN3A1	chr6	26406024	26406142	OvCa48	classI
VTSMKKPW	BTN3A2	chr6	26368680	26368798	OvCa48	classI
MEMKMRKL	PTMA	chr2	231711925	231711947	OvCa48	classI
KPEQDGDKQ	HLCS	chr21	36936676	36937275	OvCa84	classI
QQKVPRSA	FGD5	chr3	14897954	14898025	OvCa104	classI
TMVSSSPTRLSGLWMT	PRPF8	chr17	1658603	1658682	OvCa58	classII
VARLPLCRREPEP	WDR97	chr8	144109435	144109926	OvCa58	classII
LVEDCLPNGGAQ	WDR97	chr8	144109435	144109926	OvCa58	classII
MKANPALTVMVSSPTRL	PRPF8	chr17	1658603	1658682	OvCa65	classII
KFSILTKHKAFNRS	ZNF208	chr19	21972981	21973953	OvCa70	classII
ILTKHKAYNTFSIL	ZNF208	chr19	21972477	21973953	OvCa70	classII
KFSILTKHKAYNT	ZNF208	chr19	21972477	21973953	OvCa70	classII
LTKHKKIHTGETPYK	ZNF208	chr19	21972057	21972477	OvCa84	classII
DTSSSSTGHATPL	MUC4	chr3	195781609	195783385	OvCa84	classII
KFSILTKHKSFSTF	ZNF208	chr19	21972057	21973953	OvCa84	classII
LHPHPQGAEAAAGLLQRAR	HSPA1B	chr6	31828262	31828846	OvCa99	classII
RGEGSARRQAGQGPDSRPG	HSPA1B	chr6	31828262	31828846	OvCa99	classII
GRGLHPHPQGAEAAAGLL	HSPA1B	chr6	31828262	31828846	OvCa99	classII
LTKHKAYNTFSIL	ZNF208	chr19	21972477	21972981	OvCa104	classII
NPITVRNVGKPFGTAQ	ZNF485	chr10	43616923	43617050	OvCa109	classII
NPALTMVSSSPTRLSGL	PRPF8	chr17	1658603	1658682	OvCa109	classII
ANPALTVMVSSSPTRLSGL	PRPF8	chr17	1658603	1658682	OvCa109	classII
PALTMVSSSPTRLSGLW	PRPF8	chr17	1658603	1658682	OvCa109	classII
KPFGTAQALLNIRD	ZNF485	chr10	43616923	43617050	OvCa109	classII
KFSILTKHKAYNTF	ZNF208	chr19	21972477	21973953	OvCa114	classII
KIMKANPALTIV	PRPF8	chr17	1658603	1658682	OvCa114	classII