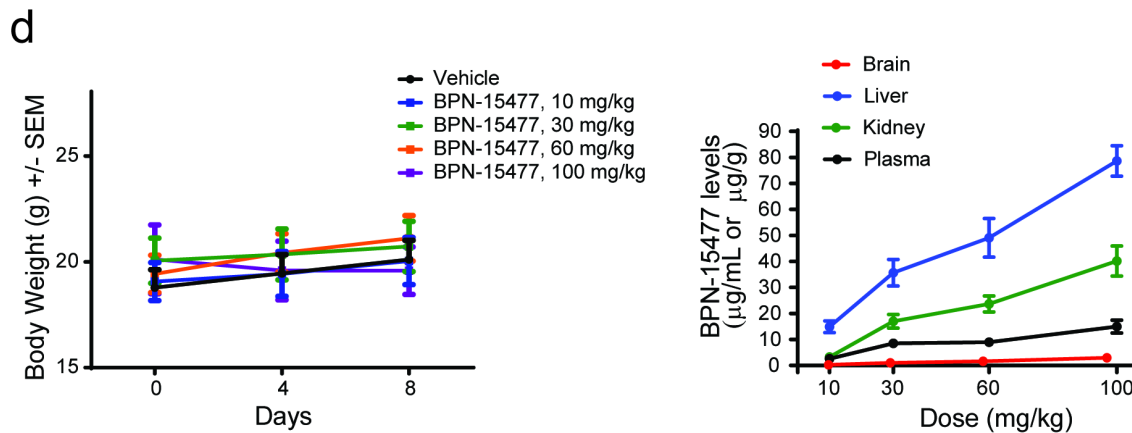
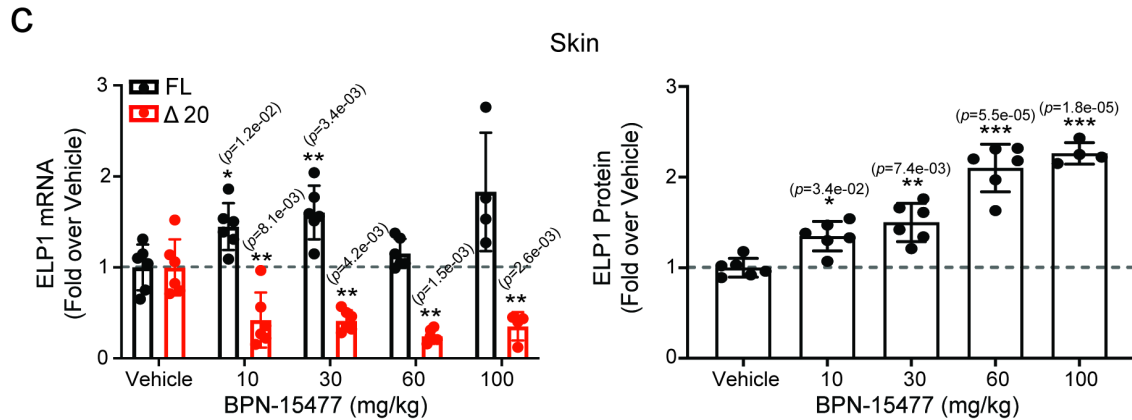
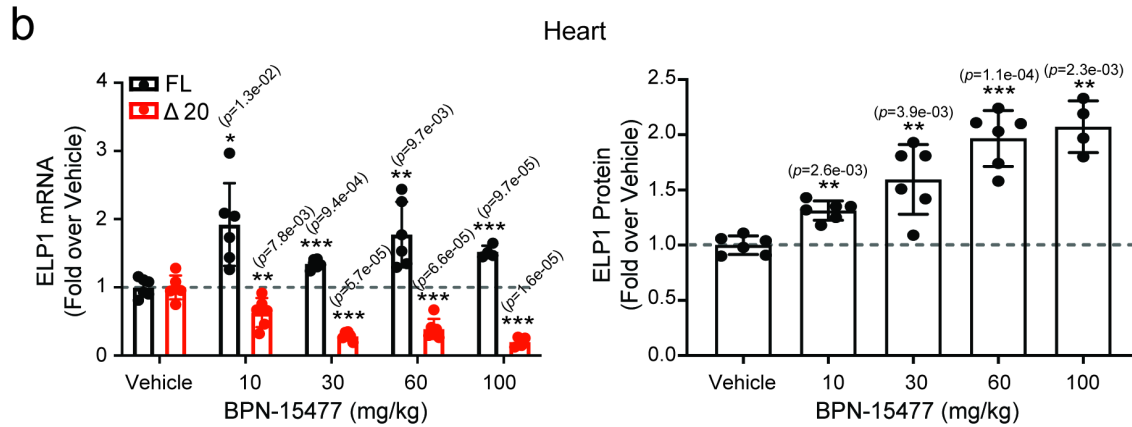
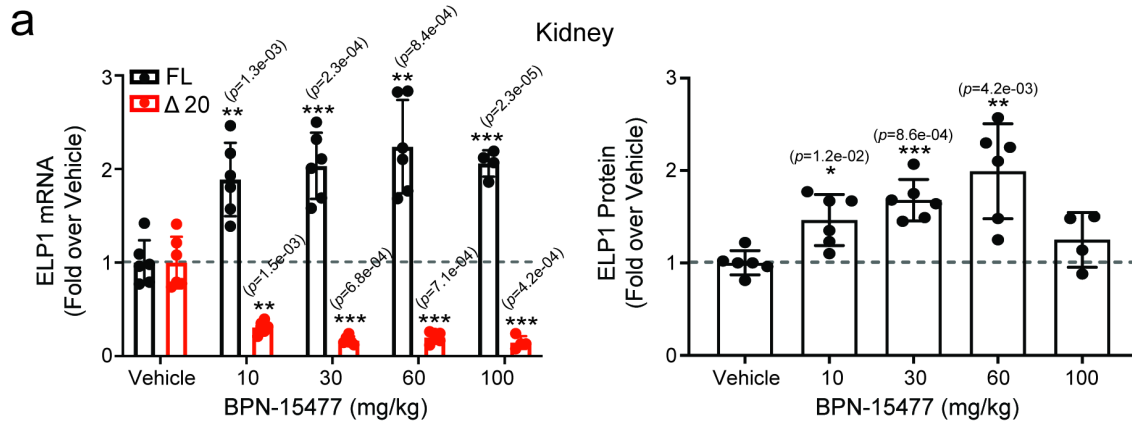


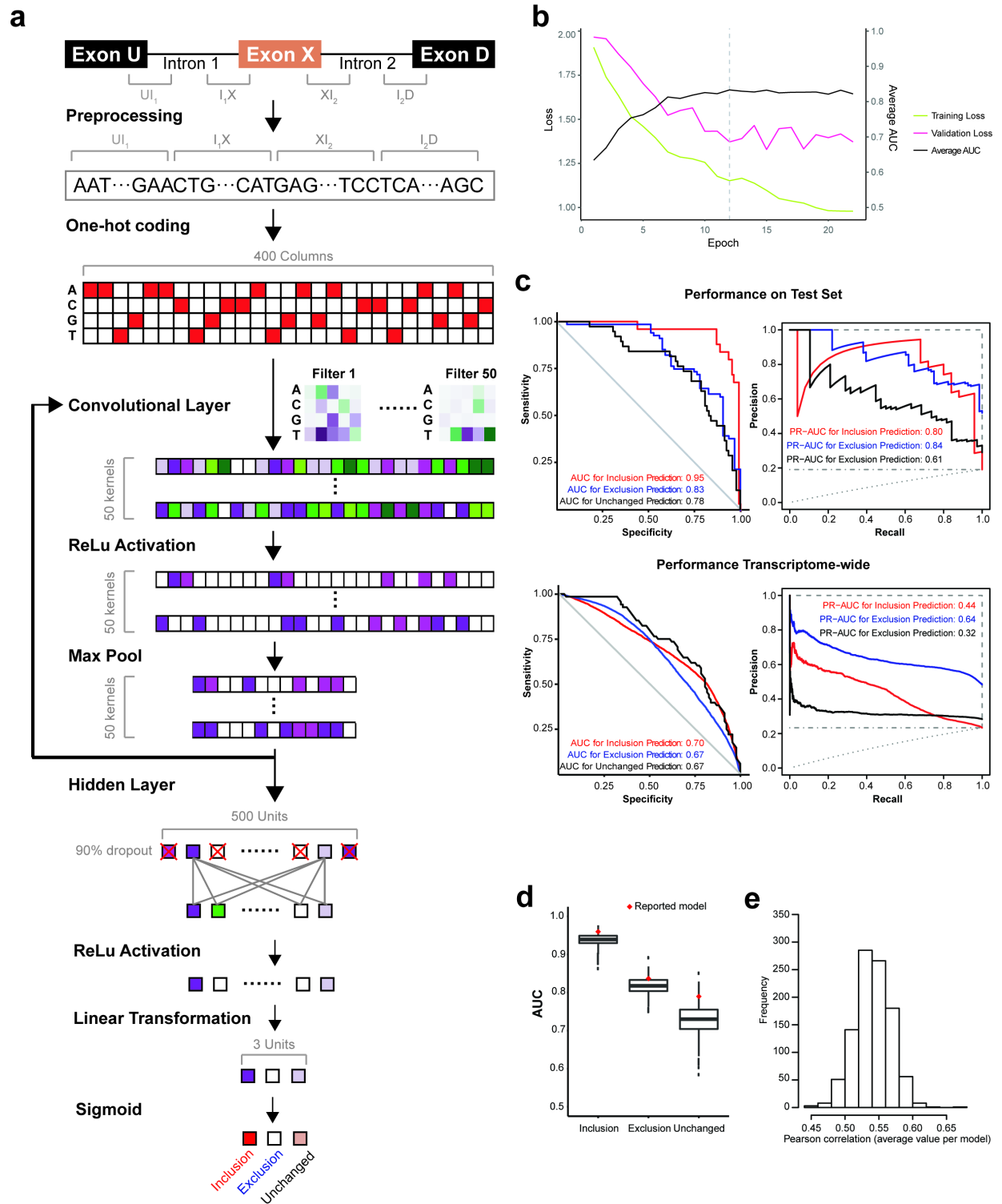
Supplementary Information

Deep learning approach to identify targets of a novel therapeutic for human splicing disorders.

Dadi Gao^{1,2,3,8}, Elisabetta Morini^{1,2,8}, Monica Salani^{2,8}, Aram J. Krauson², Anil Chekuri^{1,2}, Neeraj Sharma⁴, Ashok Ragavendran^{1,3}, Serkan Erdin^{1,3}, Emily M. Logan², Wencheng Li⁵, Amal Dakka⁵, Jana Narasimhan⁵, Xin Zhao⁵, Nikolai Naryshkin⁵, Christopher R. Trotta⁵, Kerstin A. Effenberger⁵, Matthew G. Woll⁵, Vijayalakshmi Gabbeta⁵, Gary Karp⁵, Yong Yu⁵, Graham Johnson⁶, William D. Paquette⁷, Garry R. Cutting⁴, Michael E. Talkowski^{1,2,3*}, Susan A. Slaughaupt^{1,2*}



Supplementary Fig. 1 Bpn-15477 treatment increases full length *ELP1* transcript and protein amount in several tissues of the *TgFD9* mouse. (a-c) Relative expression of full-length (FL) and $\Delta 20$ *ELP1* mRNA (left graphs), and ELP1 protein quantification (right graphs) in kidney (a), heart (b) and skin (c) after oral doses of BPN-15477 ranging from 10 to 100 mg/kg in adult transgenic *TgFD9* mouse (n=4-6 mice in each treatment group). Comparisons are done within the same color-coded group, against the vehicle-treated mice under two-tailed Welch's *t*-test. Data are presented as mean values +/- SD. The unadjusted *p* values are displayed. (d) Weight assessment of *TgFD9* mice in different treatment groups. (e) BPN-15477 distribution in in brain, liver, kidney and plasma. The levels of compound were measured using mass spectrometry. In the figure, **p* < 0.05; ***p* < 0.01; ****p* < 0.001.

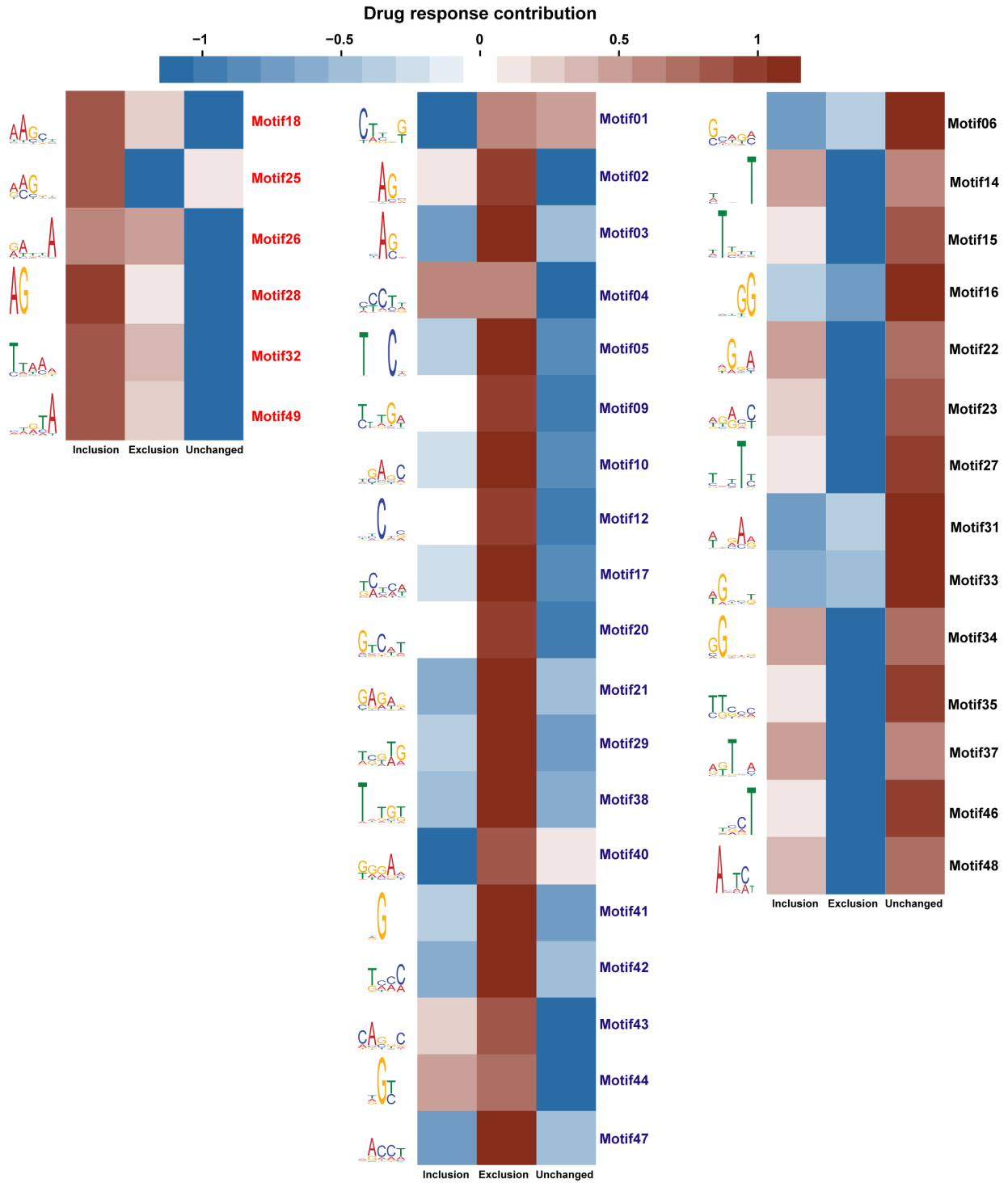


Supplementary Fig. 2 Training process of CNN

(a) CNN model workflow. For each exon-triplet, the sequences from UI_1 , I_1X , XI_2 and I_2D are concatenated and then one-hot coded. In the model, two rounds of convolution were applied before

the hidden layer. Each round of convolution consists of a convolution layer of fifty filters, a ReLU activation layer and a max pooling layer of size 2. After two rounds of convolution, the output is converted and connected to a hidden layer with 90% dropout rate. The output from the hidden layer is ReLU transformed again and is then linearly transformed into a vector of three, representing three different treatment responses. The final nonlinear sigmoid maps each element in the vector to a value between 0 and 1, and it is considered as the probability of drug responsiveness. **(b)** Training progress of the CNN model. The x axis represents the number of epochs iterated during training. The left y axis shows loss score measured by binary entropy while the right y axis shows the average AUC of prediction from three classes. The red line and blue line show training and validation loss respectively along the growth of epochs. The black line shows the improvement of AUC along the growth of epochs. The vertical dashed line indicates the stop of training to avoid overfitting. **(c)** *Left panels*: ROC-AUC curves of prediction on each class using the test set and all transcriptome-wide events (see Methods), respectively. The x axis represents specificity while the y axis represents sensitivity. The solid grey line indicates the boundary beneath which the prediction is not better than a random guess. *Right panels*: PR-AUC curves prediction on each class using the test set and all transcriptome-wide events (see Methods), respectively. The x axis represents recall while the y axis represents precision. The three dashed grey lines, from top to bottom, indicate the best performance, the performance from a random classifier and the worst performance, respectively. In both panels, the performance curve for inclusion, exclusion and unchanged response are marked in red, blue and black, respectively. **(d)** AUCs for the different treatment responses obtained from 1,000 random-initiated CNN models. In the boxplot, the middle lines inside boxes indicate the medians. The lower and upper hinges correspond to the first and third quartiles. Each box extends to 1.5 times inter-quartile range (IQR)

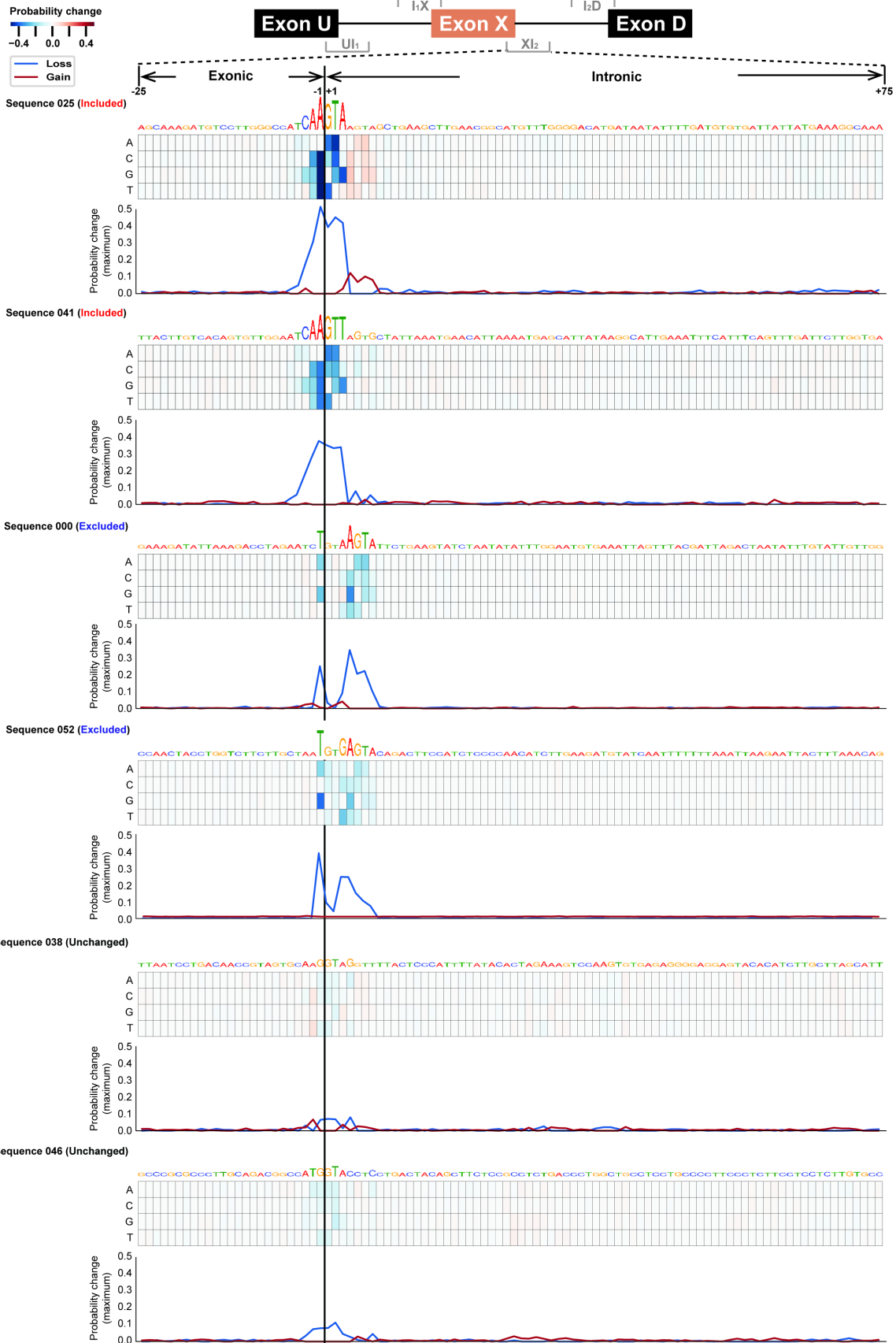
from upper and lower hinges respectively. For each box, $n=1000$. Outliers were not shown. The red diamond indicates the performance of our original CNN model. **(e)** The distribution of Pearson correlation between the top thirteen most important motifs in the reported CNN model and the top thirteen most important motifs in each of the 1,000 random-initiated models.



Supplementary Fig. 3 Motifs identified by the CNN model

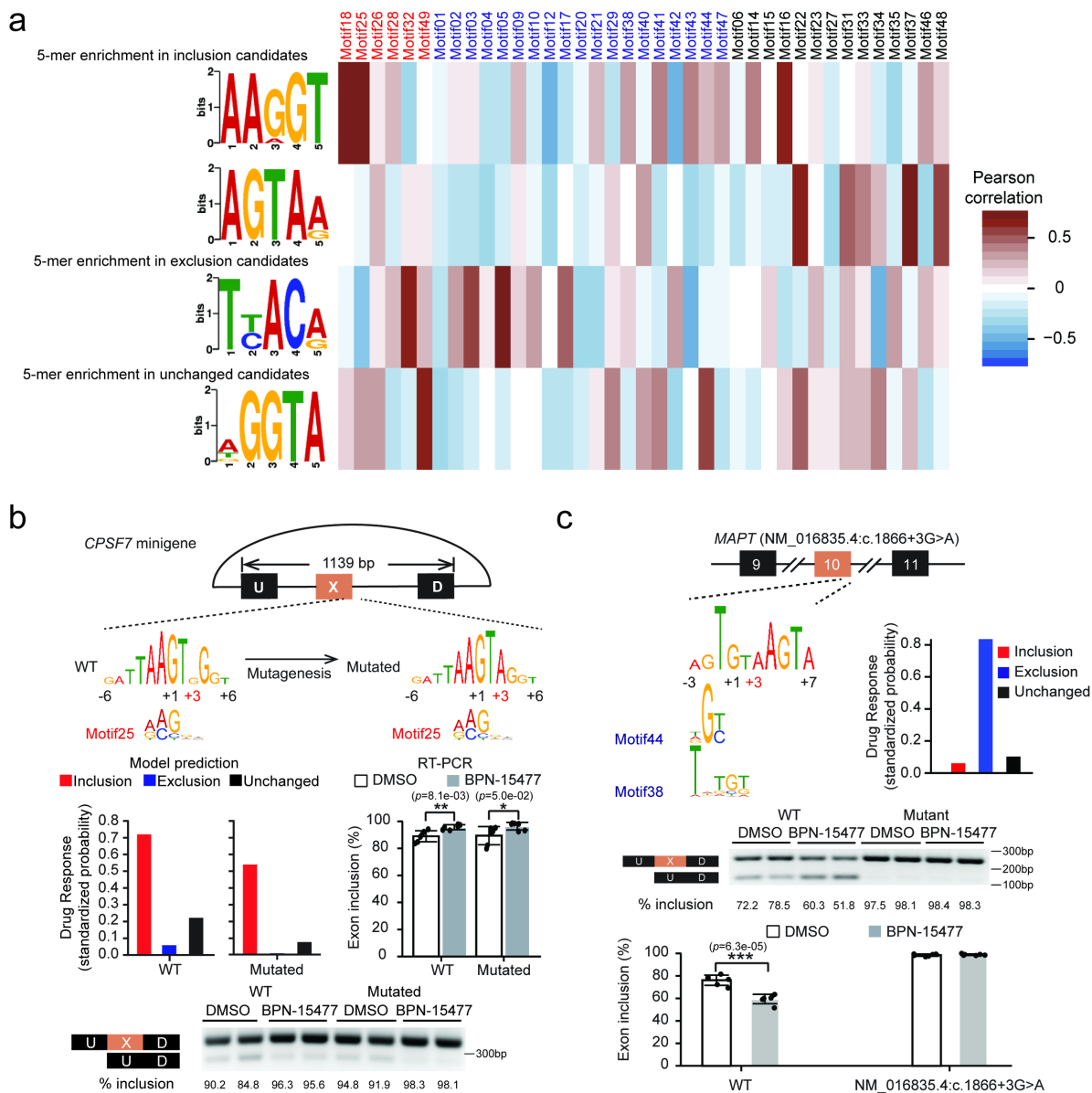
(a) Heatmap shows all motifs identified by the model. The color indicates the directional contribution of each motif. The brown domain indicates positive contribution while the blue

domain indicates negative contribution. The LOGO plot of each motif is shown on the left side of the heatmap. The motif names are displayed on the right side of the heatmap. Colors assigned to the motif names, suggest their contribution to inclusion (red), exclusion (blue) and unchanged (black) response respectively.



Supplementary Fig. 4 *In silico* saturation mutagenesis analysis on selected sequences reveals specific patterns at the 5' splice sites of the middle exons.

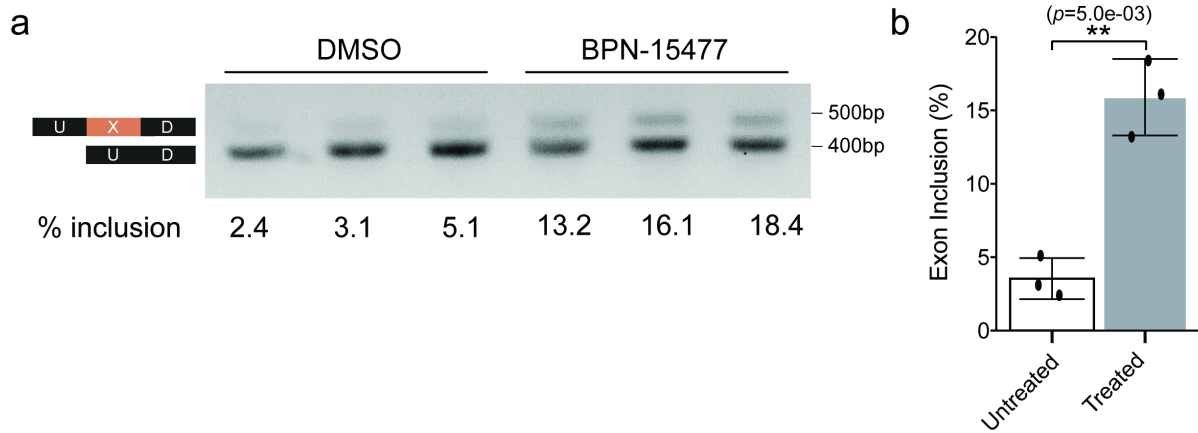
The indices of the selected sequences are provided and inclusion or exclusion responses are indicated in the brackets. For each sequence, three rows are used to demonstrate the results from *in silico* saturated mutagenesis. *Upper row*: The LOGO plot represents the original input sequence before *in silico* mutagenesis while the nucleotide height represents the maximal prediction change during *in silico* mutagenesis. *Middle row*: The heatmap represents the prediction changes according to four different nucleotides. The blue color domain indicates that the new prediction has a lower score to support the class the sequence is supposed to be while the red color domain indicates the new prediction is more in favor of the class the sequence is supposed to be. The darker the color, the stronger is the prediction. The white color indicates no change in the prediction. *Bottom row*: The curve plot represents the maximal prediction change among the four nucleotides at each position.



Supplementary Fig. 5 Validation of the CNN model by *k*-mer analysis, minigene with random mutagenesis, and minigene mimicking pathogenic mutation.

(a) Heatmap of Pearson correlation between the motifs identified by 5-mer enrichment analysis (each row) and the motifs identified by CNN model (each column). The brown domain indicates a positive correlation while the blue domain indicates a negative correlation. **(b)** *Upper row:* Schematic representations of *CPSF7* triplet minigene. The length of the exon triplets cloned into the minigenes are shown. The sequences adjacent to the 5' splice site of the middle exon are shown

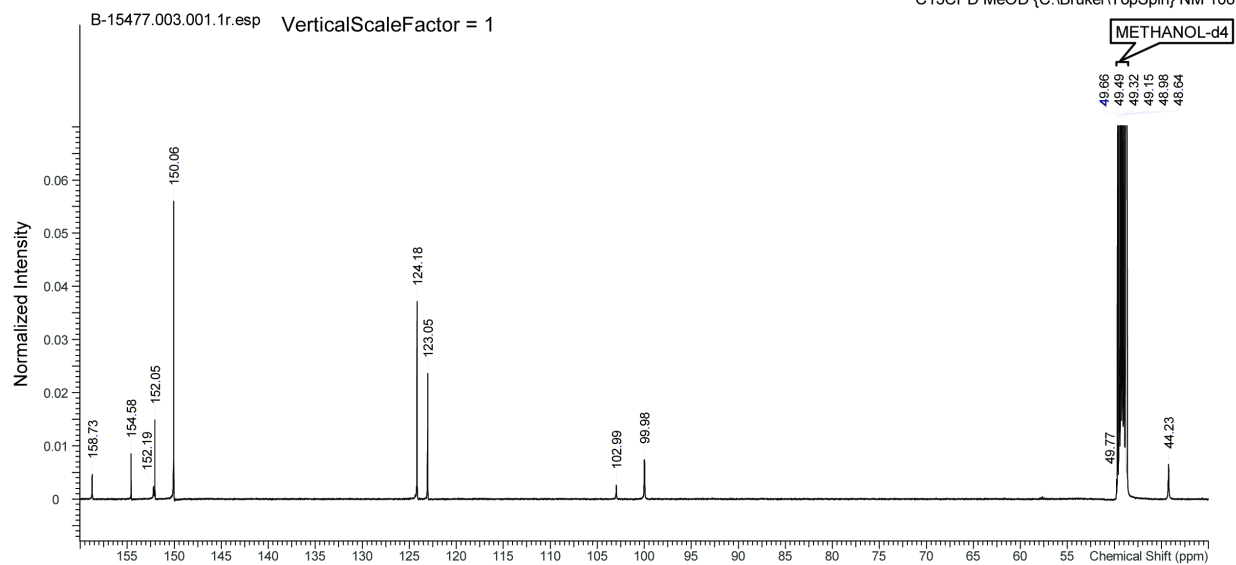
in LOGO plots. The height of each nucleotide was estimated using *in silico* saturated mutagenesis (See Methods). The red coordinate numbers indicate the positions of mutations relative to the 5' splice sites. Their closely matched CNN motifs are indicated beneath. *Middle row*: Splicing changes of the middle exons in both wildtype and mutated exon triplets, predicted by the CNN model (*left*) and measured by RT-PCR of the minigene (*right*). To make the bar plots each experiment was repeated six times (n=6, two-tailed Welch's *t*-test). Data are presented as mean values +/- SD. The unadjusted *p* values are displayed. *Bottom row*: Example of splicing changes induced by the treatment using a minigene splicing assay. The percentage of middle exon inclusion is indicated beneath each lane. **(c)** *Upper row*: Schematic of the *MAPT* minigene construct. The sequence around the 5' splice site of the middle exon is shown in LOGO plots, with closely matched CNN motifs indicated below. The red coordinate number indicates the position of the mutation relative to the 5' splice site. The bar plots show that the CNN model predicts increased exon exclusion. *Middle row*: RT-PCR validation of treatment responses in cell lines expressing the mutated minigene. To generate the bar plots, each experiment was repeated six times. *Bottom row*: The bar plots demonstrate the splicing change promoted by BPN-15477 treatment (n=6). Data are presented as mean values +/- SD. The statistical significance is determined via two-tailed Welch's *t*-test. The unadjusted *p* values are displayed. In the figure, **p* < 0.05; ***p* < 0.01; ****p* < 0.001.



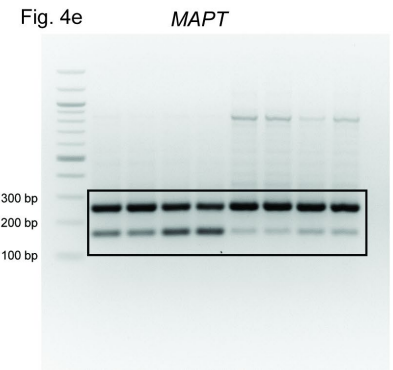
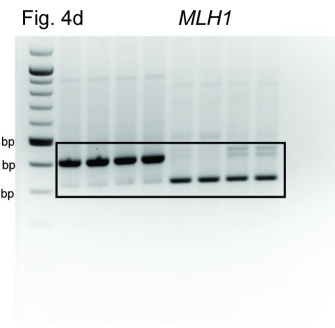
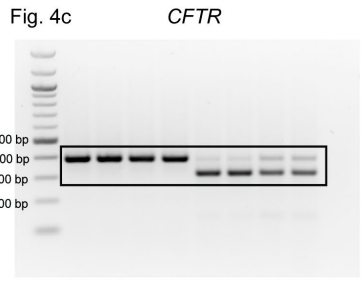
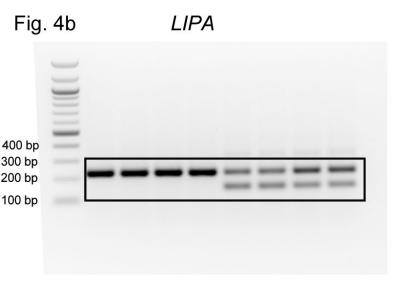
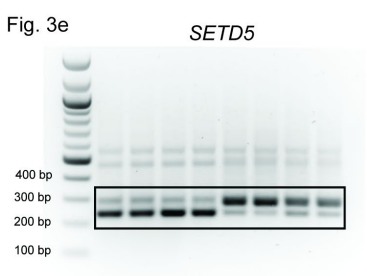
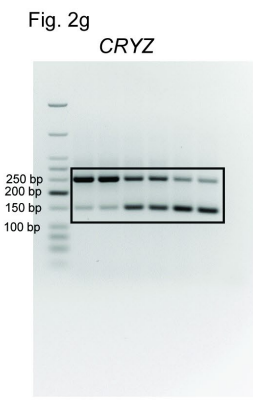
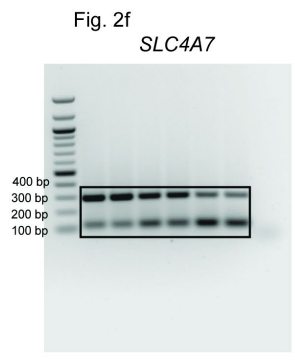
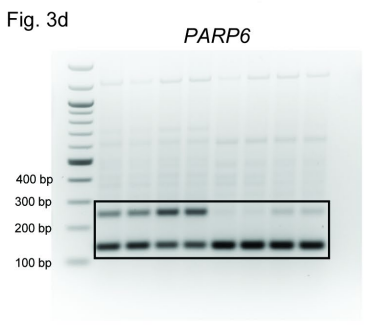
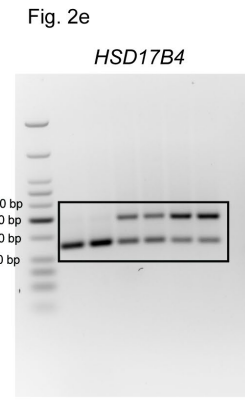
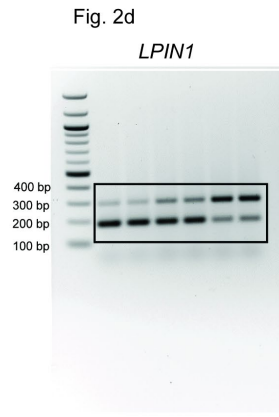
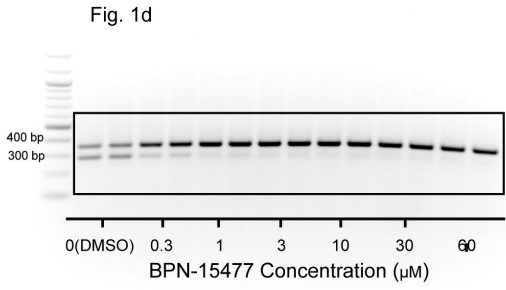
Supplementary Fig. 6 Correction of CFTR splicing in in treated CFBE cells. (a) RT-PCR validation of treatment responses in CFBE-Flpin cells stably expressing c.2988G>A-EMG-i14-i18. Cells were treated with 3 μ M of BPN-15477 for three days. (b) The bar plots demonstrate the splicing correction promoted by BPN-15477 treatment. To generate the bar plots, the experiment was performed in triplicate (n=3). Data are presented as mean values \pm SD. The statistical significance is determined via two-tailed Welch's *t*-test. The unadjusted *p* values are displayed. In the figure, ***p* < 0.01.

BPN-15477 C13 500MHz

4/7/2021 9:46:13 AM
C13CPD MeOD (C:\Bruker\TopSpin) NM 108



Supplementary Fig. 7 The identity and purity of the BPN-15477 established through ^{13}C -NMR.



Supplementary Fig. 8 Full scan of RT-PCR gels in Figures 1-4.

Fig. 5a

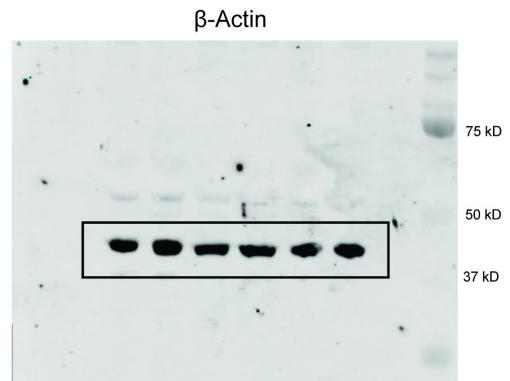
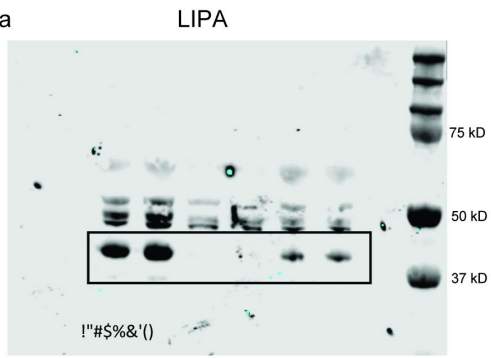
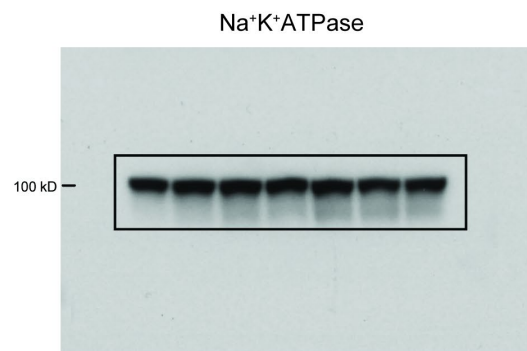
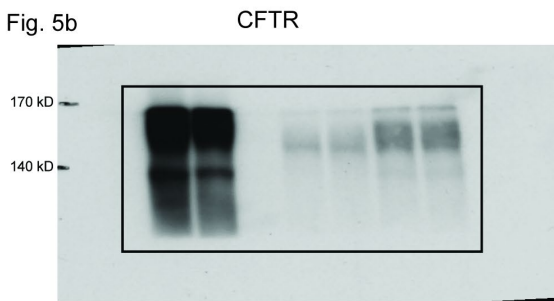


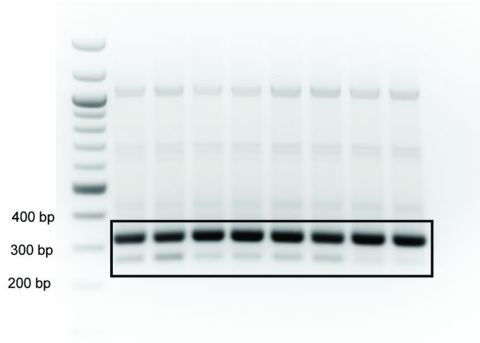
Fig. 5b



Supplementary Fig. 9 Full scan of Western blot gels in Figure 5.

Supplementary Fig. 5b

CPSF7



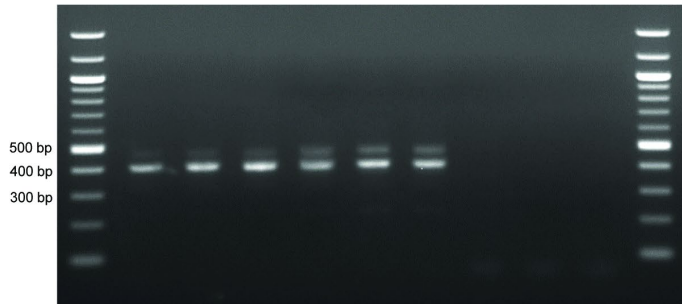
Supplementary Fig. 5c

MAPT



Supplementary Fig. 6

CFTR



Supplementary Fig. 10 Full scan of RT-PCR gels referring to Supplementary Figures 5-6.

Supplementary Table 1 Cell lines information

Coriell #	Genotype	Sex	Age	Race
AG16409	WT	Male	12 years	Caucasian
GM03348	WT	Male	10 years	Caucasian
GM08402	WT	Male	32 years	Caucasian
GM01652	WT	Female	11 years	Caucasian
GM02036	WT	Female	11 years	Caucasian
GM00041	WT	Female	3 months	Caucasian

Supplementary Table 2 RT-PCR validation of PSI changes

Gene Symbol	PSI change in RT-PCR	PSI change in RNASeq
<i>U2AF1L4</i>	-0.257166667	-0.471499182
<i>LPINI</i>	0.4005	0.414356642
<i>CYLD</i>	0.532666667	0.41517918
<i>SLC4A7</i>	-0.349416667	-0.426139237
<i>CRYZ</i>	-0.57025	-0.51621752
<i>KLC1</i>	0.343	0.425290445
<i>HSD17B4</i>	0.602666667	0.597164946
<i>AFMID</i>	0.453166667	0.519212881
<i>EPB41L2</i>	0.4412	0.410758255
<i>CD99P1</i>	-0.522633333	-0.457793242
<i>LRRC28</i>	-0.397	-0.417275893
<i>EVC</i>	-0.549583333	-0.496281655
<i>MEGF6</i>	-0.344	-0.357743325
<i>MYEF2</i>	0.4263	0.3531151
<i>ASXL1</i>	-0.0715	-0.154897036
<i>SPTANI</i>	-0.1619167	-0.249074873
<i>COPS8</i>	0.27758333	0.230359308
<i>PDZD11</i>	0.3545	0.257397272
<i>KTN1</i>	0.11175	0.107924419
<i>SPPL2A</i>	0.371	0.378300362

Supplementary Table 3 Primers and melting temperature (T_m) used for RT-PCR analysis

Gene	Forward 5'-3'	Reverse 5'-3'	T_m
<i>ELP1</i>	CCTGAGCAGCAATCATGTG	TACATGGTCTTCGTGACATC	58 C°
<i>KLC1</i>	CGC AGT GGT TCC TTT AGC	CAC TGC TGC TGC TGT CG	60 C°
<i>LPIN1</i>	GCT GTG ATT TAC CCT CAG TCA GC	CTT AGC AGC CTG CGG CAG C	64 C°
<i>HSD17B4</i>	GCA GAA AGA GGA GCG TTA	GTT GGC CAC TGC TTT TCC	56 C°
<i>SLC4A7</i>	GCT ACA GAG GAC TGG ACG	CTA GAA CTG GAC CTG TGC TCC	60 C°
<i>LRRC28</i>	GAT ATA GTG CTG CAG CGT GC	CAA CTA TGT TAT TTG AGT GCA GG	60 C°
<i>AFMID</i>	GCC TTT CTT CCT GTT CTT TCA CG	GGT GAG CAC GTT GTC CTT CT	60 C°
<i>CD99P1</i>	CGA CCC AGC ACC TCT TAA TTC	CGG TGG AAT CAG GCT GCT TG	62 C°
<i>CRYZ</i>	GCA CTG CTG GTA CTG AGG AAG	CTT TGC CAT GGT GTC TCG TGG	64 C°
<i>IP6K2</i>	AAC AAG CCA AGG AGC CAA GA	ATT CAG GCC ACA CTT CCC TG	62 C°
<i>EVC</i>	TGC CCT GAA GCT GAT GAA GG	GGT GCC AGC GTC TGC TTC	62 C°
<i>EPB41L2</i>	GGA GAA GTA CCT GAT GCC GAC	CTC ACT CTC ACT GCT GCT G	60 C°
<i>CYLD</i>	GAT GGT TCT ACA CAG CCA CC	CTT CCC AGT AGG GTG AAG TGA C	62 C°
<i>MEGF6</i>	CTGGTTTGGAGAGGCCTGTG	GGGACGGACTGCAACCTCA	60 C°
<i>MYEF2</i>	GTACCGTGGTGCGATGACTA	TCTGACAAATATCTGGTTGCCT	60 C°
<i>ASXL1</i>	GCCTCGAGTTGTCCTGACTC	TTCAGGCAGGAGGAAGAGGA	60 C°
<i>SPTANI</i>	CGATCGTCAGGGTTTTGTGC	TTCTGGAGCACCTCAACCTG	60 C°
<i>COPS8</i>	CTGAGGGACAGTCTGGGGTT	ATGGAGCAAATATAAAGCTAGAAGC	60 C°
<i>PDZD11</i>	GTCAGTGAGCGGAGTCTGAG	TGAGGAGGAATCCATGCTGG	60 C°
<i>KTNI</i>	AGGCAGAGATGGAACGATCT	TGCAAATCACCAGCTACCTTCT	60 C°
<i>SPPL2A</i>	TCATGGTTGAACTCGCAGCT	TGAGGCACACACTCATTACTGA	60 C°
<i>LIPA</i>	CCCAGAGTGCCTTTTTGAA	CCCAGTCAAAGGCTTGAAAC	60 C°
<i>CFTR</i> HEK	GCAGTGATTATCACCAGC	GGAGGAAATATGCTCTCAAC	58 C°
<i>PTEN</i>	TTGCACAATATCCTTTTGAAGACCA	TTAGCATCTTGTTCTGTTTGTGGA	60 C°
<i>CPSF7</i>	ATTGCCCTTGACCCAGAGTT	GAGCGAATAACCTGGATCAGC	60 C°
<i>SETD5</i>	TGTGGTGGAAATTGCCCTTAC	GCAGAATTGCCCTTCTGATA	60 C°
<i>MLH1</i>	GTGGAATTGCCCTTGAGC	CACTGTGCTGGATATCTGCTG	60 C°
<i>PARP6</i>	GGAATTGCCCTTGATCATCT	AATTGCCCTTCTGCAGTTTG	60 C°
<i>MAPT</i>	ACCCAAGTCGCCGTCTTCCGCC	CACCTTGCTCAGGTCAACTGGT	60 C°
<i>CFTR</i> CFBE	GTGGCTGCTTCTTTGGTTGTG	GGAGGAAATATGCTCTCAAC	55 C°