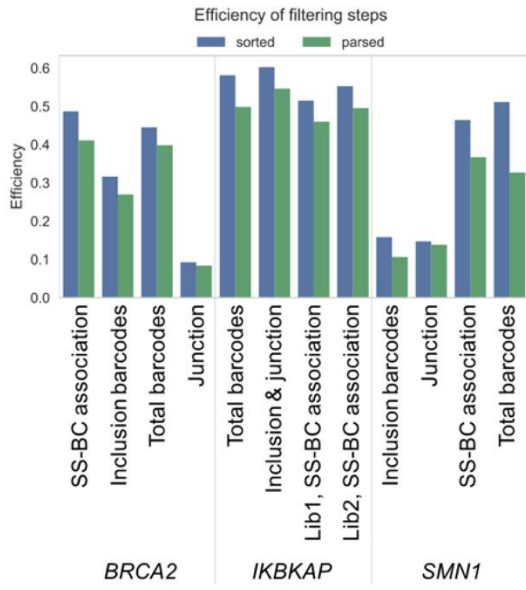
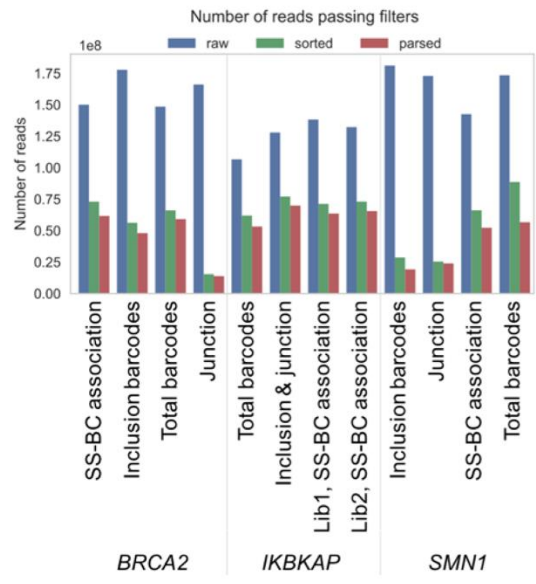


Figure S1. Detailed diagram of high-throughput method used to assess all 5'ss sequences. Related to Figure 2A. RE1 and RE2 represent restriction enzyme sites used to insert the minigenes into the pcDNA5 expression vector, which has a cytomegalovirus (CMV) promoter and a bGH polyadenylation site (pA).

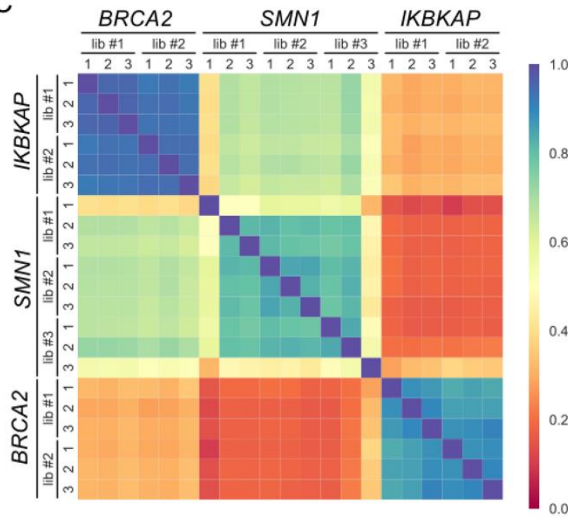
**A**



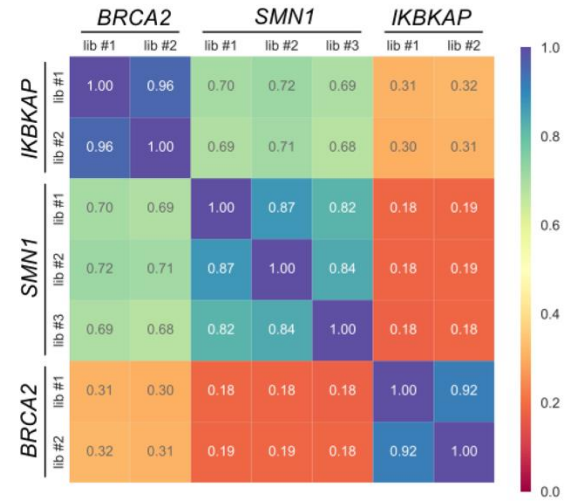
**B**



**C**



**D**



**E**

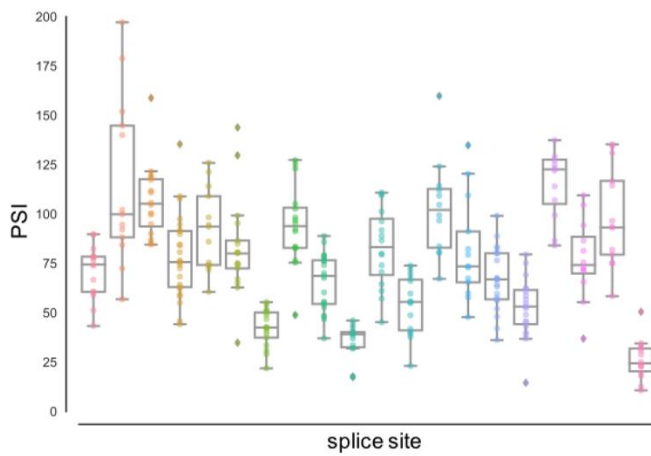
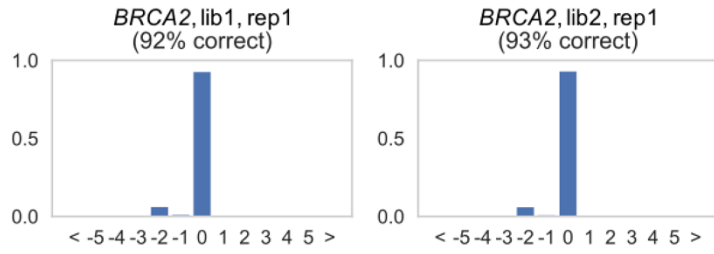


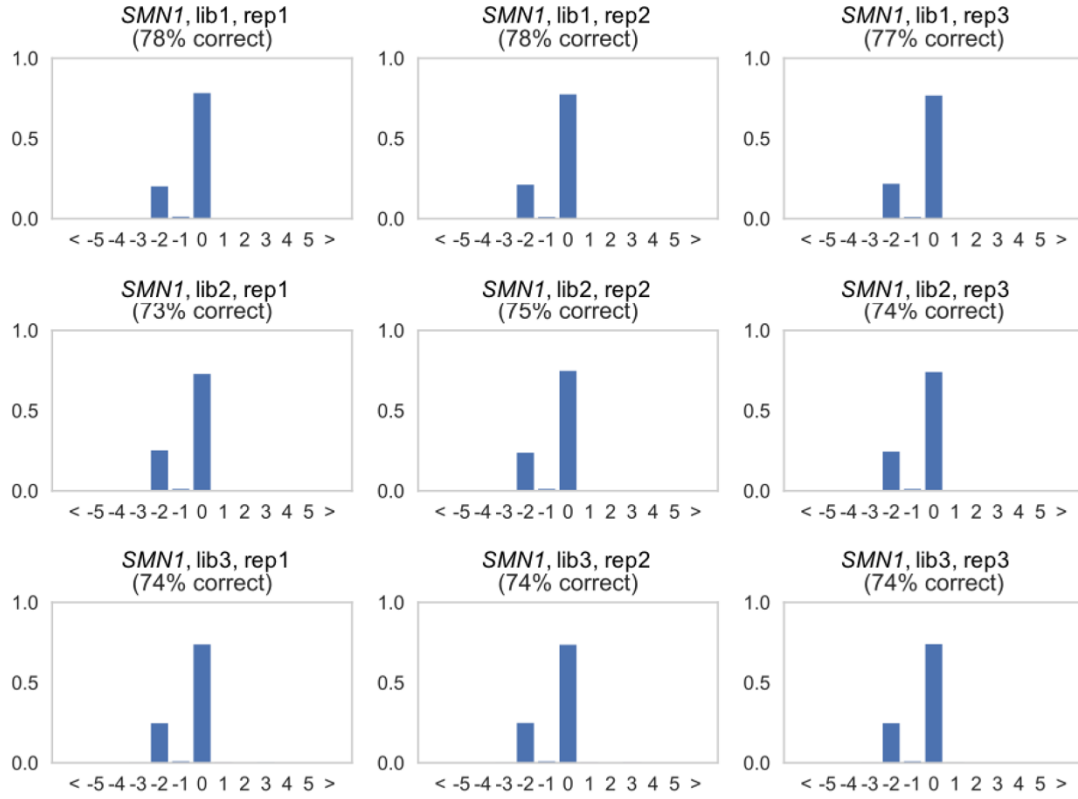
Figure S2. Parsing of raw high-throughput results and association matrix of the coefficient of determination comparing libraries. Related to Table and Figure 2D.

- (A) The efficiency of filtering of raw sequencing data through the bioinformatics pipeline.
- (B) Number of reads for analysis after filtering.
- (C) Two low-quality datasets (*SMN1* library 1, replicate 1, and *SMN1* library 3, replicate 3) were removed from subsequent analyses.
- (D) Coefficient of determination values show that the independently derived libraries highly correlate with each other within a context.
- (E) Box plot of the PSI of individual barcodes for 20 randomly selected 5'ss. For this analysis, each barcode is required to have at least 10 counts in the total RNA sample in order to accurately calculate a PSI. The median PSI of the 5'ss is required to be  $\geq 20$ . For each 5'ss, a minimum of 10 associated barcodes is necessary for the 5'ss to be included in the analysis. The central rectangle spans the first to the third quartile, with the median line segment. The vertical line presents the maximum and minimum.

A



B



C

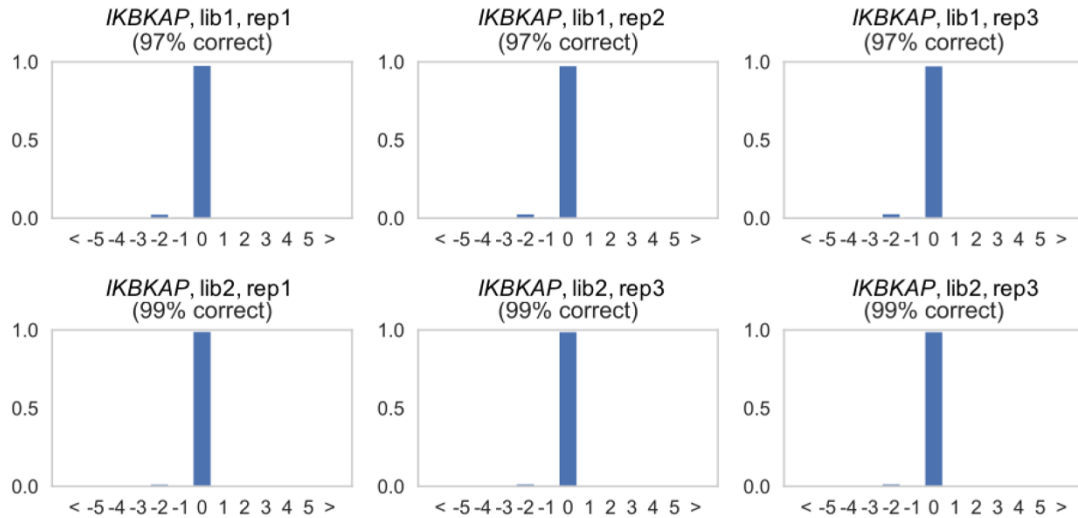
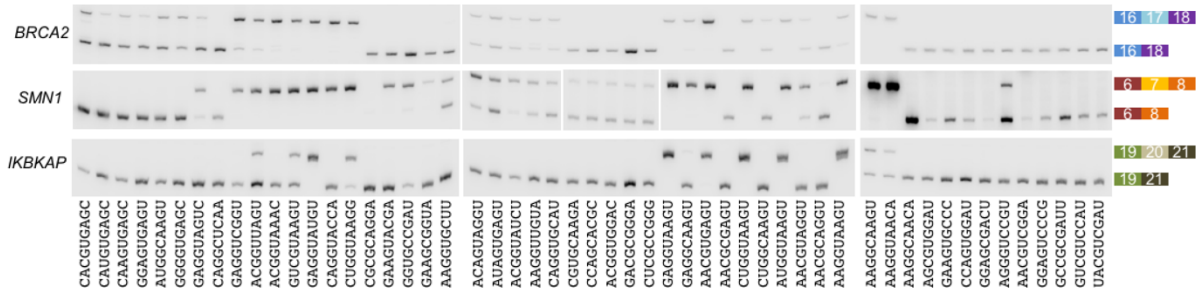


Figure S3. Exon-exon junction analysis of libraries. Related to Figure 2 and “high-throughput analysis of the activity of all 5’ss sequences”.

- (A) Sequencing results for the exon-exon junction reveal that a secondary GU at the -2 and -1 positions is preferentially used when the GU or GC at the +1 and +2 positions escapes recognition. 5’ss sequences with the secondary -2G-1U were removed from further analysis.
- (B) Same as A, but in *SMN1*.
- (C) Same as A, but in *IKBKAP*.

A



B

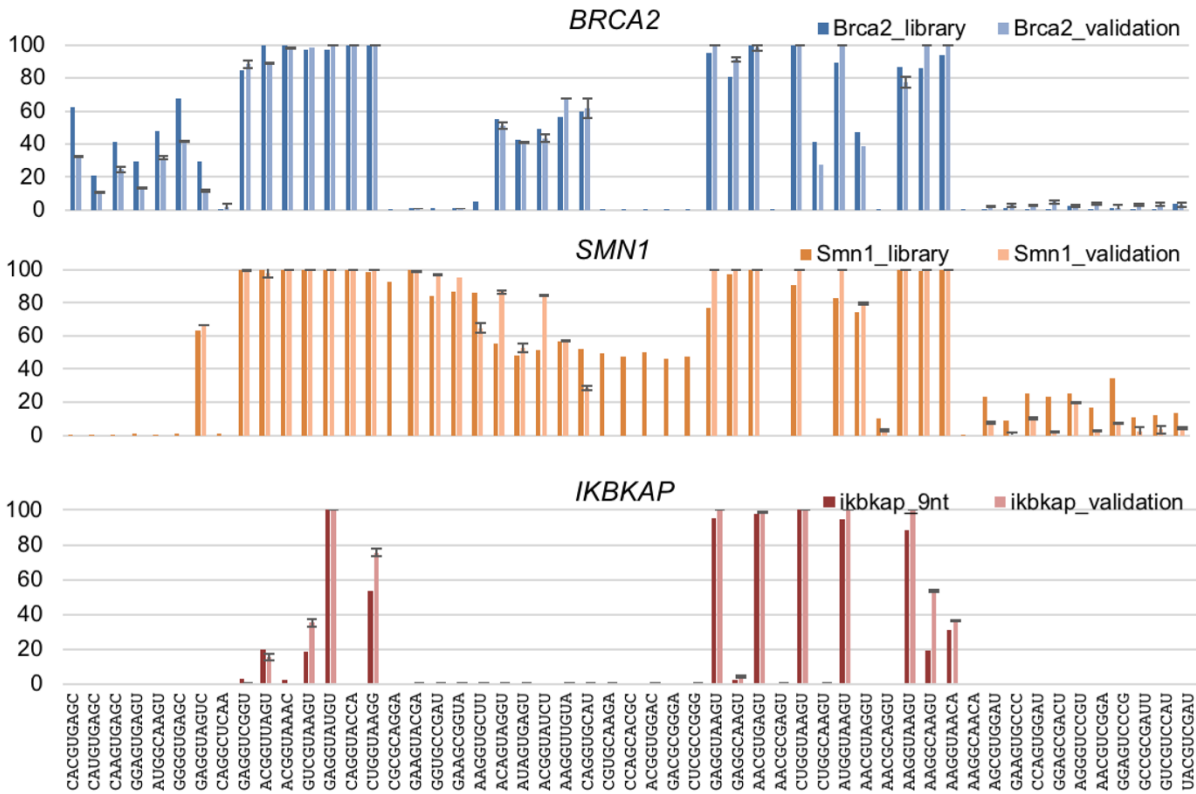


Figure S4. Manual validations of randomly-selected 5'ss in the three contexts. Related to Figure 2F.

- (A) Gel showing the splicing results of the same 53 randomly-selected 5'ss in the *BRCA2*, *SMN1*, and *IKBKAP* minigenes. Gel images are representatives of triplicates.
- (B) Graphs comparing the percent spliced in (PSI) of 5'ss derived from the library results and manual validations. Standard deviation is represented by error bars.

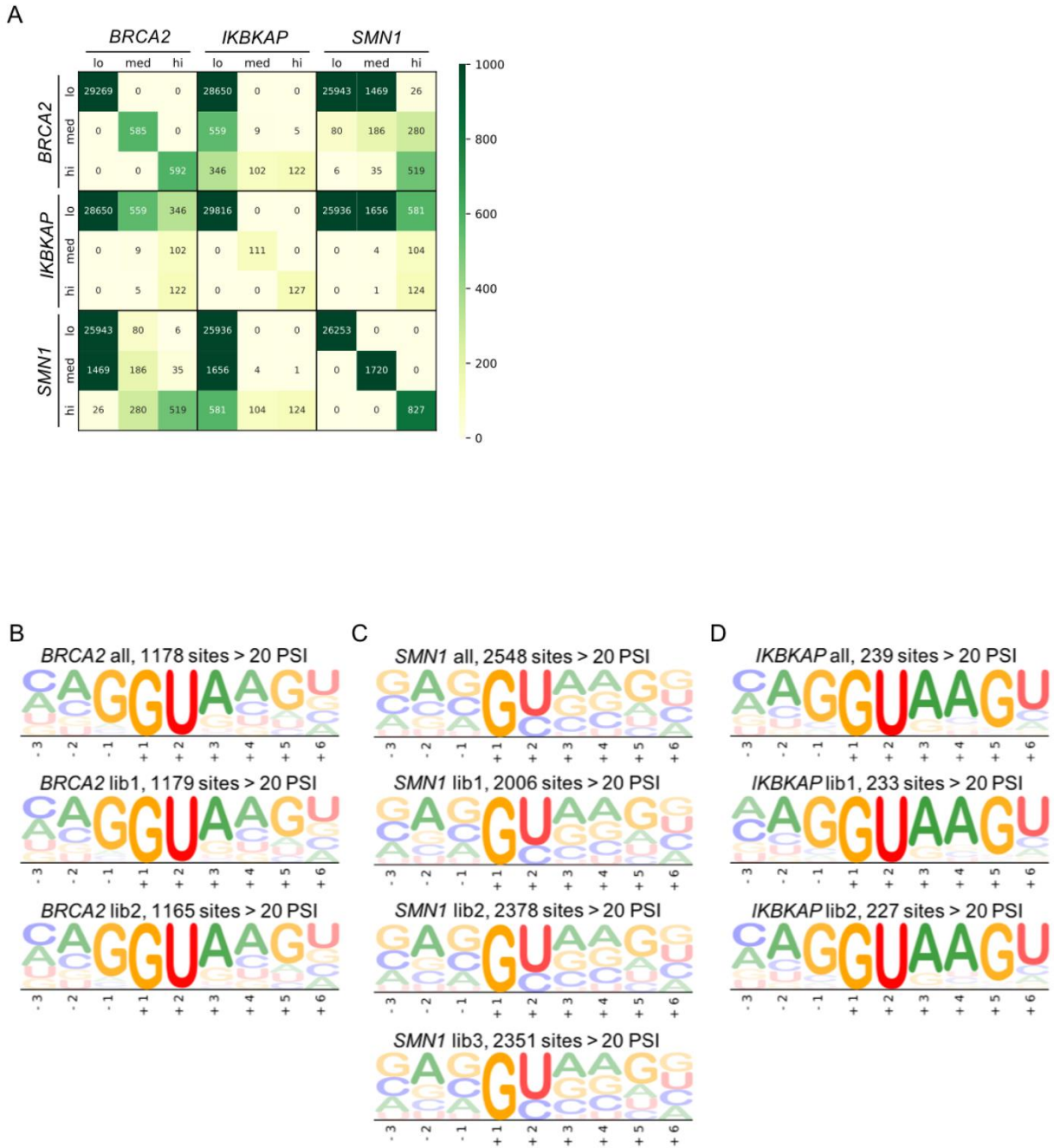


Figure S5. Sequence logo of 5'ss in each context. Related to Figure 3B and 3C.

- (A) Heat map illustrating the overlaps between 5'ss with lo (PSI < 20), med ( $20 \leq \text{PSI} < 80$ ) or hi ( $80 \leq \text{PSI}$ ) activity levels among the three gene contexts.
- (B) Sequence logo for all 5'ss compiled or separated by each independently-derived library in *BRCA2*.
- (C) Same as A, but in *SMN1*.
- (D) Same as A, but in *IKBKAP*.

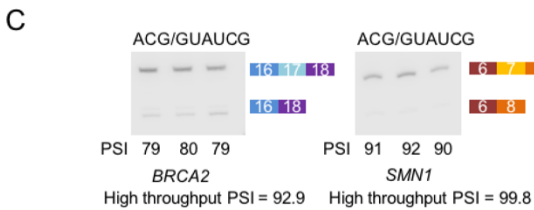
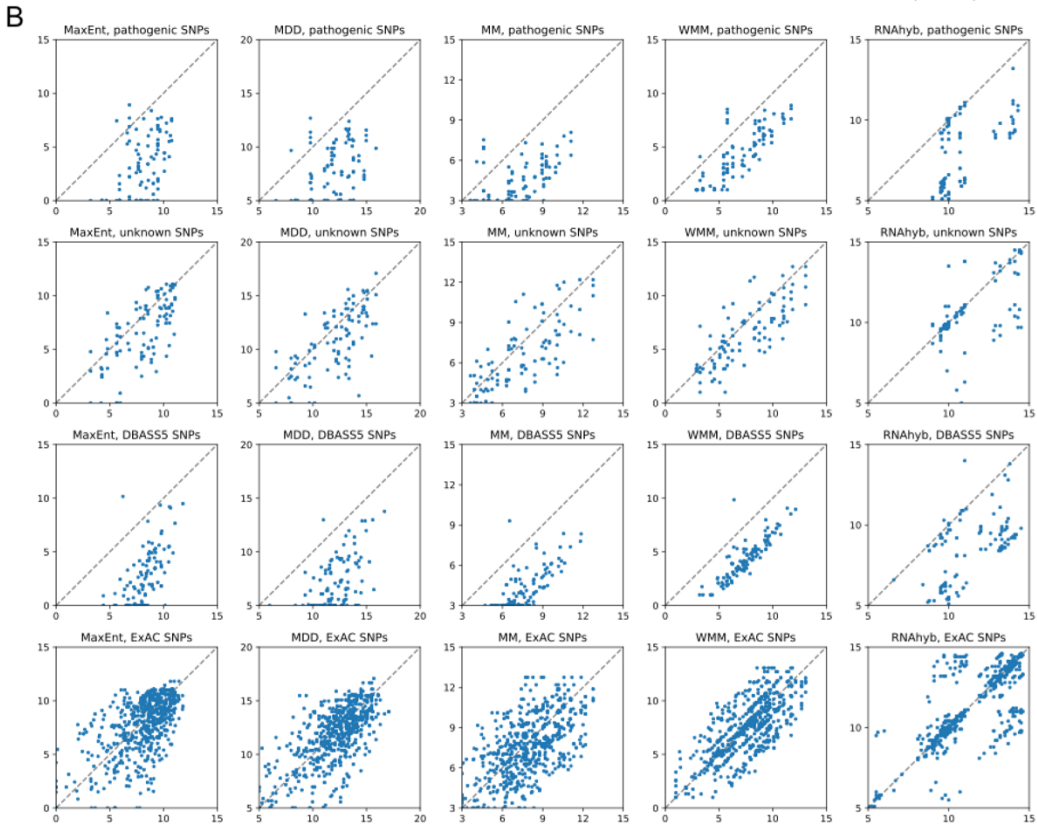
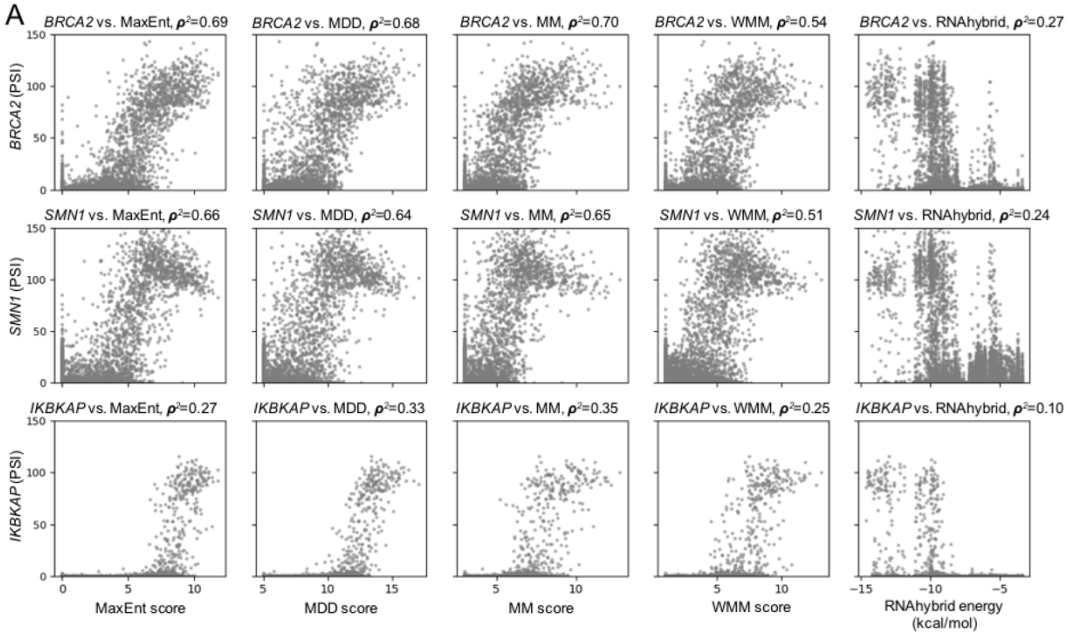
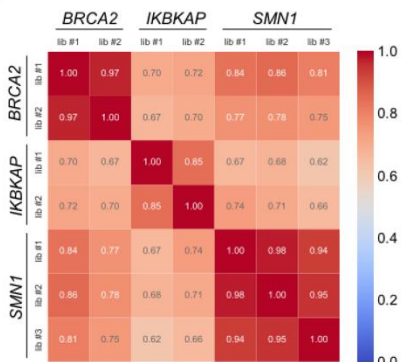




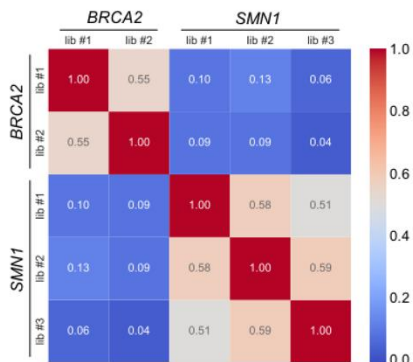
Figure S6. Comparison between library results and several conventional prediction algorithms. Related to Figure 3D, 3E, and 6.

- (A) Scatter plots comparing the predicted values for each 5'ss using maximum entropy (MaxENT; Yeo et al., 2004), maximum dependence decomposition (MDD; Burge et al., 1998), first-order Markov model (MM; Krogh et al., 1994), weight matrix model (WMM), and RNAhybrid (RNAhyb; Kruger et al., 2006) versus the experimentally derived library results.
- (B) Scatter plots, corresponding to the computational scoring matrices used above, comparing the predicted PSI of the WT 5'ss sequences to mutant 5'ss sequences known to be pathogenic in *BRCA1* and *BRCA2* (pathogenic SNPs), with unclassified or uncertain significance in *BRCA1* and *BRCA2* (unknown SNPs), 5'ss mutations across a broad range of genes and diseases (DBASS5), and for 5'ss SNPs with >10% frequency found in the human population (ExAC SNPs).
- (C) RT-PCR validation of the usage of the 5'ss with the sequence ACG/GUAUCG, which showed high inclusion ratio in *BRCA2* and *SMN1* library results but does not occur naturally as a 5'ss in the human transcriptome. Gel image is representative of triplicates. Percent spliced in (PSI) is indicated below each lane.

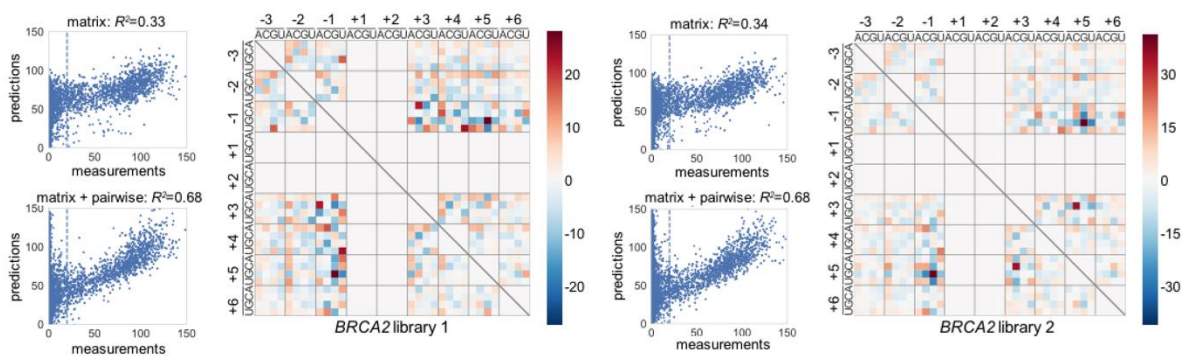
**A**



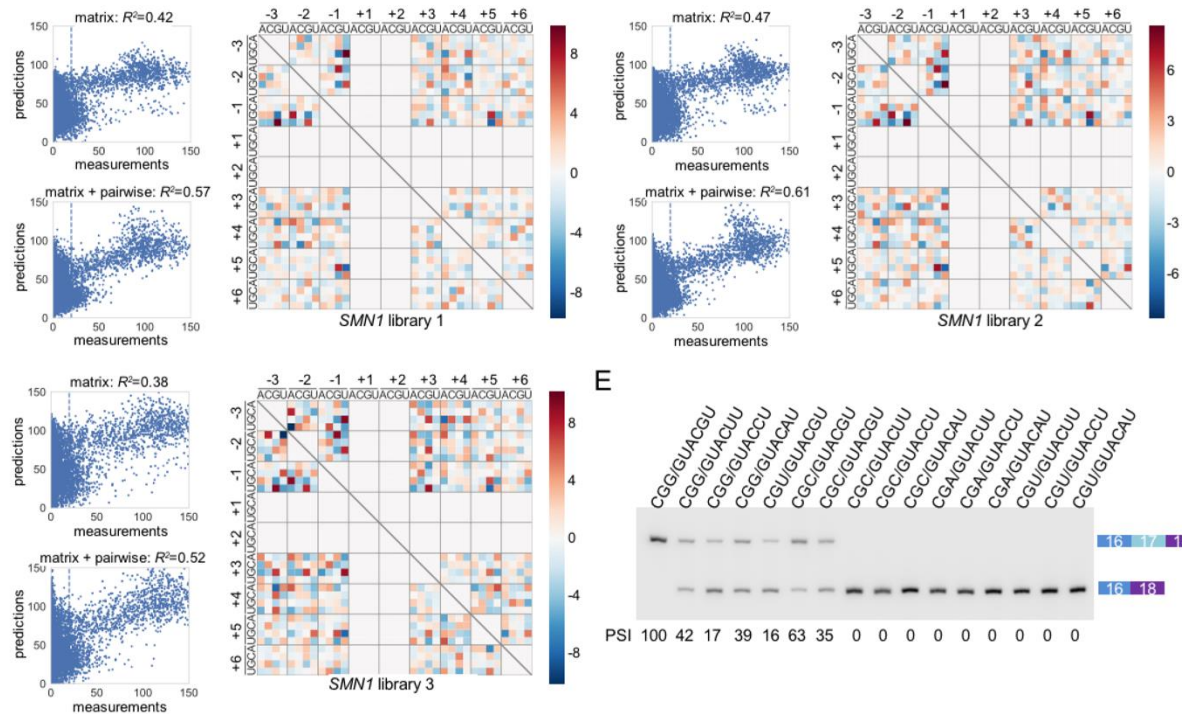
**B**



**C**



**D**



**E**

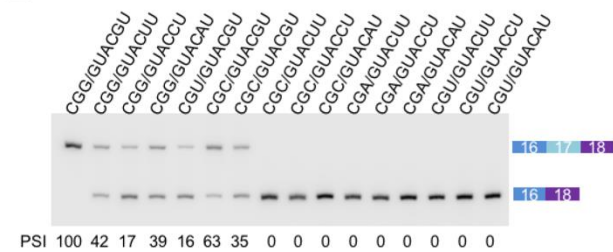


Figure S7. Accounting for pairwise associations between nucleotides improves the predictability of the results. Related to Figure 4.

- (A) Association matrix showing correlations of parameters in the linear matrix model.
- (B) Same as A, but in pairwise associations. The number of active 5'ss was insufficient to analyze pairwise associations for *IKBKAP*.
- (C) Scatter plot showing the predicted PSI versus the measured PSI when considering a linear matrix model only (top) and a matrix model together with pairwise interaction (bottom). The heat map shows pairwise associations between nucleotides at different positions in *BRCA2*, separated by libraries. Red indicates a positive interaction, and blue indicates a negative interaction.
- (D) Same as C, but for *SMN1*.
- (E) Validation of a comprehensive series of mutations at -1 and +5 positions of the 5'ss CGG/GUACGU, in the *BRCA2* context. Gel image is representative of triplicates. Percent spliced in (PSI) is indicated below each lane.

Table S1 - 5'ss sequences that have activity in library results, but do not occur naturally in the human transcriptome. Related to Figure 3E.

<b><i>BRCA2</i></b>	<b><i>SMN1</i></b>
ACGGUAUCG	ACGGUAUCG
<u>CGCGUACGU</u>	<u>CGCGUACGU</u>
AAAG <u>CGCUG</u>	AACGUACGG
CCAGUACCG	CACGUAC <u>GC</u>
CCCGCC <u>GUG</u>	CCGGUAUAC
<u>GCGGUA</u> AAC	<u>CGCGUAC</u> GA
<u>GUGGCAUCG</u>	GAAGCG <u>GUA</u>
	GACGCG <u>GUA</u>
	GACGUACGA
	<u>GGCGCGUAU</u>

Table S2 – Predicted strength of the upstream and downstream 3'ss. Related to Figure 5.

	Gene (input sequence, 5'→3')	MAXENT	MM	WMM
Upstream 3'ss	<i>BRCA2</i> (TTCTACTTTTATTTGTTTCAGGGC)	8.33	9.46	10.78
	<i>SMN1</i> (TTCCTTTATTTTCCTTACAGGGT)	10.92	13.08	15.51
	<i>IKBKAP</i> (ACTGCTTTAATTTATTTAAGATG)	6.36	6.51	4.57
Downstream 3'ss	<i>BRCA2</i> (ATTTTTGTTTTCACTTTTAGATA)	11.50	12.16	12.62
	<i>SMN1</i> (TTCTAATTTCTCATTTCAGGAA)	10.77	10.86	10.52
	<i>IKBKAP</i> (CTTTCTCTGTCTTCTCACAGACT)	11.96	12.06	13.35

Table S3 – Primer sequences. Related to STAR methods.

Name	Sequence (5'→3')
FRT F	CTGGCTAACTAGAGAACCCACTGC
BRCA2 18R	GCTGTGTCATCCCTTTCCATTATC
BRCA2 7R	GAGCACAGTAGAACTAAGGGTGG
SMN1 R	TAGTGGTGTCAATTTAGTGCTGC
IKBKAP R	GATTGATTCTCAGCTTTCTCATGC
BRCA2 Bsu36I ss top	CATCATCCTAAGGAATTTGCTAATAGATGCCTAAGCCCAGAAAGGG TGCTTCTTCAACTAAAATANNGNNTTAAAGCAGCAGGTGGAT GCACATGATGACATAAT
BRCA2 NotI bc bot	TACTACCGCCGCGGNNTTCTAGAATGCA GGTGATTATGTCATCATGTGCATC
BRCA2 PE1 top	AATGATACGGCGACCACCGAGATCTACACTCTTCCCTACACGACG CTCTCCGATCTNNNNNNNNNT
BRCA2 PE1 bot	CATGTNNNNNNNTCTAGCCTTCTCGCAGCACATCCCTTTCTCA CATCTAGAGCCACCAGCGGCATAGTAA
BRCA2 PE2 top	CGTNNNNNNNAGATCGGAAGAGCGGTTTCAGCAGGAATGCCGA GACCGATCTCGTATGCCGTCTTCTGCTT
BRCA2 PE2 bot	TTCGTCTTCTGCCGTATGCTCTAGCCAGAGCCGTAAGGACGACTTG GCGAGAAGGCTAGANNNNNNNNT
BRCA2 insert F	CATCATCACCTGCAGAGTTAAAGCATTACATTACG
BRCA2 insert R	TACTACCACCTGCACACTCTAGAATTACTACTTTAAC
BRCA2 17F	AGATGCCTAAGCCCAGAAAG
BRCA2 18F	GGCTCTCCTGATGCCTGTAC
BRCA2 18R	GCTGTGTCATCCCTTTCCATTATC
BRCA2 BC R	GGCAACTAGAAGGCACAGTCG
BRCA2 BC-LID F	CCCTACACGACGCTCTTCCGATCTNNNNNNNNNGGCTCTCCTGAT GCCTGTAC
BRCA2 BC-LID R	CATTCCTGCTGAACCGCTCTTCCGATCTNNNNNNNNNGGCAACTA GAAGGCACAGTCG
BRCA2 JNCT F	CCCTACACGACGCTCTTCCGATCTNNNNNNNNNAGATGCCTAAGC CCAGAAAG
BRCA2 JNCT R	CATTCCTGCTGAACCGCTCTTCCGATCTNNNNNNNNNGCTGTGTC ATCCCTTTCCATTATC
SMN1 BseRI ss top	CATCATGAGGAGCTTAAATTAANNNGYNNNNCTGCCAGCATGCAGG TGGATGCACATGATGACATAA
SMN1 NotI bc bot	ATGATGGCGGCCGNNNNNNNNNNNTCTAGAATGCA GGTGATTATGTCATCATGTGCATC
SMN1 PE1 top	AATGATACGGCGACCACCGAGATCTACACTCTTCCCTACACGACG CTCTCCGATCTNNNNNNNNNGTGCT
SMN1 PE1 bot	CNNNNNNNNNAGATCGGAAGAGCGTGTAGGGAAAGAGTGTA GATCTCGGTGGTCGCCGTATCATT

SMN1 PE2 top	GGCCGCNNNNNNNNNNNAGATCGGAAGAGCGGTTTCAGCAGGAATG CCGAGACCGATCTCGTATGCCGTCTTCTGCTT
SMN1 PE2 bot	AAGCAGAAGACGGCATAACGAGATCGGTCTCGGCATTCTGCTGAA CCGCTCTTCCGATCTNNNNNNNNNNNGC
SMN1 insert F	CATCATCACCTGCTAGGGCCAGCATTATGAAACTGAATC
SMN1 insert R	ATGATGCACCTGCCCTATCTAGAATAACGCTTCACATTCCAGATC
SMN1 7F	GAAGGAAGGTGCTCACATTC
SMN1 8F	GACACCACTAAAGAAACGATCAG
SMN1 8R	CGCTTCACATTCCAGATCTG
SMN1 BC R	GGCAACTAGAAGGCACAGTCG
SMN1 BC-LID F	CCCTACACGACGCTCTTCCGATCTNNNNNNNNNNNGACACCACTAAA GAAACGATCAG
SMN1 BC-LID R	CATTCTGCTGAACCGCTCTTCCGATCTNNNNNNNNNNNGGCAACTA GAAGGCACAGTCG
SMN1 JNCT F	CCCTACACGACGCTCTTCCGATCTNNNNNNNNNNNGAAGGAAGGTG CTCACATTC
SMN1 JNCT R	CATTCTGCTGAACCGCTCTTCCGATCTNNNNNNNNNNNCGCTTCAC ATTCCAGATCTG
IKBKAP BseRl ss top	CATCATGAGGAGAGTGGTTGGANNNGYNNNNNGCCATTGTGCAGGT GGATGCACATGATGACATAAT
IKBKAP XhoI bc bot	ATGATGCTCGAGNNNNNNNNNNNNNNNNNNNTCTAGAATGCAGG TGATTATGTCATCATGTGCATC
IKBKAP PE1 top	AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACG CTCTTCCGATCTNNNNNNNNNNNTTAAT
IKBKAP PE1 bot	TAANNNNNNNNNNNAGATCGGAAGAGCGTCTGTAGGGAAAGAGTG TAGATCTCGGTGGTCGCCGTATCATT
IKBKAP PE2 top	TCGAGNNNNNNNNNNNAGATCGGAAGAGCGGTTTCAGCAGGAATG CCGAGACCGATCTCGTATGCCGTCTTCTGCTT
IKBKAP PE2 bot	AAGCAGAAGACGGCATAACGAGATCGGTCTCGGCATTCTGCTGAA CCGCTCTTCCGATCTNNNNNNNNNNNC
IKBKAP insert F	CATCGTCACCTGCAAGCGCCATTGTAAGTTTGGCGACTAGTTAGC
IKBKAP insert R	ATGATGCACCTGCCATGTCTAGAATACTTAGGGTTATGATCAT
IKBKAP 20F	GTTGTTTCATCATCGAGCCCTGG
IKBKAP 21F	GCATGAGAAAGCTGAGAATC
IKBKAP 21R	GATTCTCAGCTTTCTCATGC
IKBKAP BC R	GGCAACTAGAAGGCACAGTCG
IKBKAP BC-LID F	CCCTACACGACGCTCTTCCGATCTNNNNNNNNNGCATGAGAAAGCT GAGAATC
IKBKAP BC-LID R	CATTCTGCTGAACCGCTCTTCCGATCTNNNNNNNNNNNGGCAACTA GAAGGCACAGTCG
IKBKAP JNCT F	CCCTACACGACGCTCTTCCGATCTNNNNNNNNNNNGTTGTTTCATCAT CGAGCCCTGG
IKBKAP JNCT R	CATTCTGCTGAACCGCTCTTCCGATCTNNNNNNNNNNNGATTCTCA GCTTTCTCATGC
PE1_v4	AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACG CTCTTC

PE2_v4	AAGCAGAAGACGGCATACGAGATCGGTCTCGGCATTCCTGCTGAA CCGCT
ACTB-F	AGAGCTACGAGCTGCCTGAC
ACTB-R	AGCACTGTGTTGGCGTACAG