

Chapter II. Functional selection of intronic splicing elements provides insight into their regulatory mechanism

Abstract

Despite the critical role of alternative splicing in generating proteomic diversity and regulating gene expression, the sequence composition and function of intronic splicing regulatory elements (ISREs) have not been well elucidated. Here, we employed a high-throughput *in vivo* Screening PLatform for Intronic Control Elements (SPLICE) to identify 125 unique ISRE sequences from a random nucleotide library. Bioinformatic analyses reveal consensus motifs that resemble splicing regulatory elements and binding sites for characterized splicing factors and that are enriched in the introns of naturally-occurring spliced genes, supporting their biological relevance. *In vivo* studies, including an RNAi silencing study, demonstrate that ISRE sequences can exhibit combinatorial regulatory activity and that complex splicing regulatory networks (SRNs) are involved in the regulatory effect of a single ISRE. From our results, we propose three mechanisms through which ISREs interact with splicing factors to achieve regulatory function: direct binding / competition, recruitment, and agonist interaction.

2.1. Introduction

Post-transcriptional gene regulatory mechanisms play central roles in programming the complexity of biological systems. One such process is alternative splicing, a dynamic mechanism that produces multiple protein isoforms from a single gene by altering the ways in which exons are joined from a single pre-mRNA¹. Splicing patterns are regulated by the interplay between auxiliary *cis*-acting elements that include exonic and intronic splicing enhancers (ESEs and ISEs, respectively) and exonic and intronic splicing silencers (ESSs and ISSs, respectively) and the trans-acting factors that modulate them, leading to a ‘splicing code’². The lack of high-throughput *in vivo* methods for analyzing the function of spliced variants and the *cis*-acting elements involved in the regulation of these transcripts has hindered the functional validation of spliced transcripts discovered through recent genome-wide mRNA sequencing studies³⁻⁵. Bioinformatic and experimental analyses have identified several RNA motifs that regulate splicing, where much of this effort has been directed toward the functional characterization of *cis*-acting exonic regulatory sequences⁶⁻⁹. Despite the widespread importance of intronic splicing regulatory elements (ISREs), knowledge regarding their sequence composition, the mechanisms through which they regulate splicing and the regulatory networks of trans-acting splicing factors by which they are bound, or splicing regulatory networks (SRNs), is limited. The development of a functional definition of ISREs and the elucidation of corresponding SRNs is of great interest given that > 90% of human genes are alternatively spliced¹⁰ and that up to 50% of disease-causing mutations affect splicing¹¹.

Several properties of ISREs have complicated their functional characterization. ISSs and ISEs have been identified near alternatively spliced exons; however, their actions appear to be antagonistic¹² suggesting that they behave in a combinatorial manner¹³. In addition, the activities of some sequences are context dependent^{10,14}. ISSs may inhibit exon inclusion by recruiting splicing repressors that directly antagonize splicing factor binding or by recruiting repressors to multiple binding sites resulting in a ‘zone of silencing’¹⁵. While several ISEs have been characterized¹⁶, the trans-acting factors that bind these sequences remain unknown¹⁷.

To begin to generate a functional definition of ISREs, we have developed a generalizable *in vivo* screening strategy for ISREs, which we call SPLICE (Screening Platform for Intronic Control Elements). SPLICE was used to identify intronic sequences that regulate the inclusion of an alternatively spliced exon that triggers rapid transcript decay through nonsense-mediated decay (NMD). Our high-throughput approach combines a systematic screening strategy, extensive genome-wide bioinformatic analyses and experimental characterization, including an RNAi silencing study, to identify ISRE consensus motifs, characterize the SRNs associated with these global regulatory elements and generate a model for ISRE regulatory function. Our results indicate that *cis*-acting intronic regulatory sequences function through combinatorial effects from multiple elements and trans-acting factors, and that the immediate transcript context has a dominant effect on ISRE function. In addition, our results support three mechanisms for ISRE regulatory function: direct binding / competition, recruitment, and agonist interaction.

2.2. Results

2.2.1. *SPLICE: a Screening PPlatform for Intronic Control Elements*

SPLICE is a high-throughput *in vivo* screen for ISRE function based on a reporter construct encoding the green fluorescent protein (GFP) fused 5' of a three-exon, two-intron mini-gene. The alternatively-spliced middle exon harbors a premature termination codon (PTC) that triggers mRNA degradation through the NMD pathway¹⁸. Auxiliary elements that regulate alternative splicing are normally positioned in proximity to splice sites^{16,19,20} and have been shown to vary in length between 10- to 30-nt¹⁹. We implemented SPLICE with the SMN1 mini-gene containing a random 15-nucleotide (nt) library positioned 45-nt upstream of the 3' ss in the first intron (Figure 2.1a). Therefore, cells with a high level of exon 7 inclusion display lower GFP fluorescence than cells in which this exon is excluded. By coupling NMD to splicing efficiency, ISREs with a range of activities can be selected using fluorescence activated cell sorting (FACS).

To test the utility of NMD as the basis of SPLICE we examined the difference in fluorescence between a NMD-based reporter construct (NMD control), containing a 15-nt control insert and a PTC in exon 7, relative to a construct lacking a PTC (GFP-SMN1 control). All constructs were stably transfected into HEK-293 FLP-In cells to generate isogenic cell lines. Flow cytometry (Figure 2.1a and Figure S2.1a) and fluorescence microscopy analyses (Figure S2.1b) reveal that the fluorescence difference between the GFP-SMN1 and NMD controls is ~22-fold. Transcript isoform analysis through quantitative real time-PCR (qRT-PCR) indicates that the level of exon 7 inclusion in the NMD control is ~60-fold less than the GFP-SMN1 control (Figure S2.1c,d), supporting that differences in fluorescence are due to exon 7 inclusion.

A library of synthetic DNA oligonucleotides containing a random 15-nt region ($\sim 1 \times 10^9$ sequences) was ligated into the NMD control construct and transformed into *Escherichia coli*. Library constructs were purified from $\sim 1 \times 10^6$ pooled transformants, representing $\sim 0.1\%$ of possible sequences. The pooled library was stably transfected into HEK-293 FLP-In cells and $\sim 450,000$ stable transformants were generated (Methods). Sequencing of the library before and after transfection demonstrated minimal sequence bias at each position (Figure S2.1e). FACS analysis indicated that $\sim 0.05\% - 0.1\%$ of the cell population exhibits fluorescence levels greater than the NMD control, corresponding to putative ISSs. Positive cells were bulk sorted, grown 2-3 weeks and re-analyzed by flow cytometry. The round-one pool exhibits an approximate six fold increase in mean fluorescence compared to the NMD control (Figure 2.1a) and was re-sorted into different groups based on fluorescence ranges (A, B, and C) to further enrich the population and select for sequences varying in splicing regulatory activity (Figure 2.1a and Figure S2.2). The enriched populations were analyzed by flow cytometry, and the mean fluorescence levels correlated well with their sorted sections (Figure 2.1b).

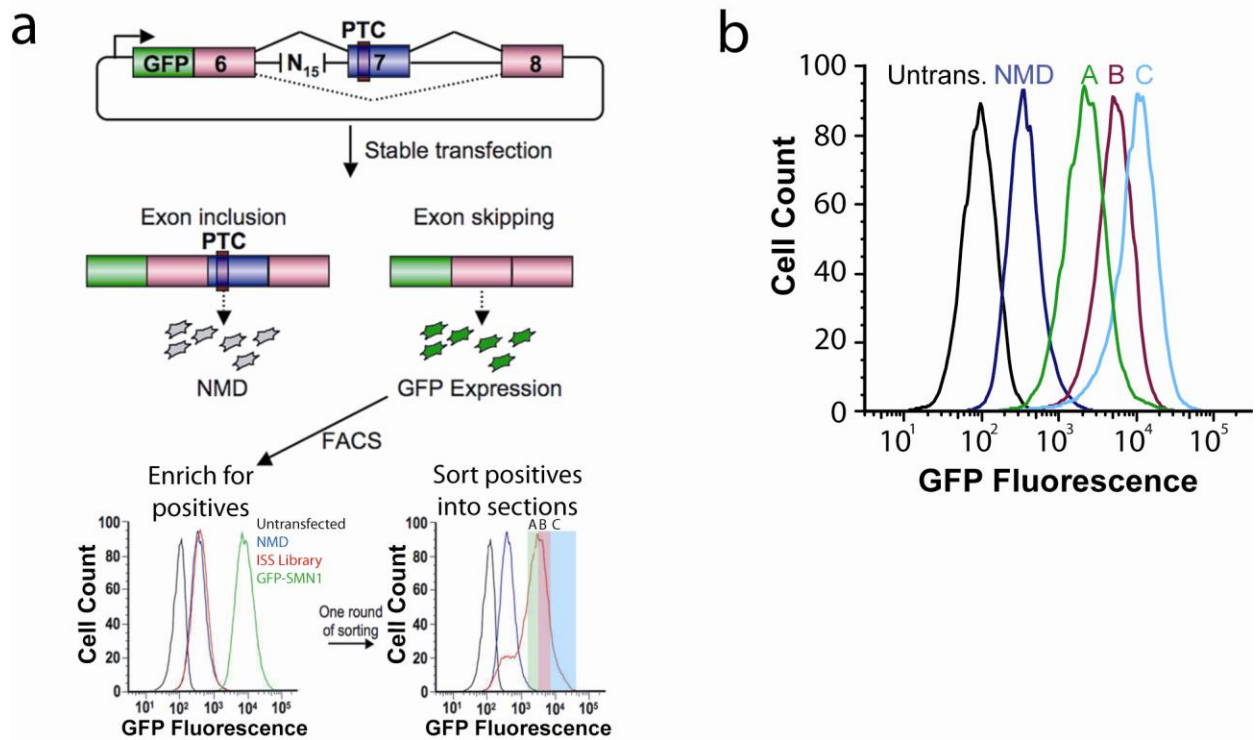


Figure 2.1. A Screening Platform for Intronic Control Elements (SPLICE) provides a generalizable *in vivo* screening strategy for ISREs. **(a)** The application of SPLICE to the screening of ISRE sequences. SPLICE couples an exon inclusion event in a mini-gene (SMN1) to the expression level of a fluorescent reporter protein (GFP) through a NMD-based reporter system. A random nucleotide library cloned into unique restriction sites in intron 6 is screened for ISSs by sorting cells exhibiting fluorescence levels higher than the negative control (NMD). The enriched cells are expanded and later sorted into sections (A, B, C) based on user-designated fluorescence levels in a second screening round. **(b)** The enriched cell populations maintain the fluorescence levels of the sorted sections (A, B, C). Following the second round of sorting, the fluorescence levels of expanded populations were re-analyzed through flow cytometry to confirm maintenance of expression levels.

2.2.2. *Recovered ISRE sequence composition correlates with sorted sections*

We identified 125 unique sequences with enhanced fluorescence from 480 sequenced isolates (Table S2.1). Three of the recovered ISRE sequences exhibit significant (12 of 15-nt) similarity to portions of the SMN1 mini-gene, suggesting that these sequences may be involved in the cooperative assembly of repressor elements on the SMN1 transcript (Figure S2.3). The sequences have a higher level of G (35.8%) and reduced levels of T (22%), C (18.6%), and A (23.6%) (Figure S2.4a). The dinucleotide CC is overrepresented in the ISRE dataset, while others, such as AC, AG, CA, GT, TA, TC, and TG, are only slightly enriched (Figure S2.4b).

Recovered 15-mers were subjected to hierarchical clustering to determine the overall sequence similarity between elements (Figure 2.2 and Figure S2.5)^{6,9}. SPLICE-generated sequences are generally diverse (>95% of sequences differ by more than 1-nt), indicating that the majority of the recovered sequences arose from independent selection. We evaluated the association between clusters of sequences (using a dissimilarity score cutoff of 1.1) and the fluorescent section from which they were sorted. In particular, clusters 11 and 13 show a significant association with the sorted sections, while clusters 10 and 18 do not. The resulting clusters generally contain sequences from identically sorted sections suggesting that sequence composition correlates with cellular fluorescence (Figure 2.2).

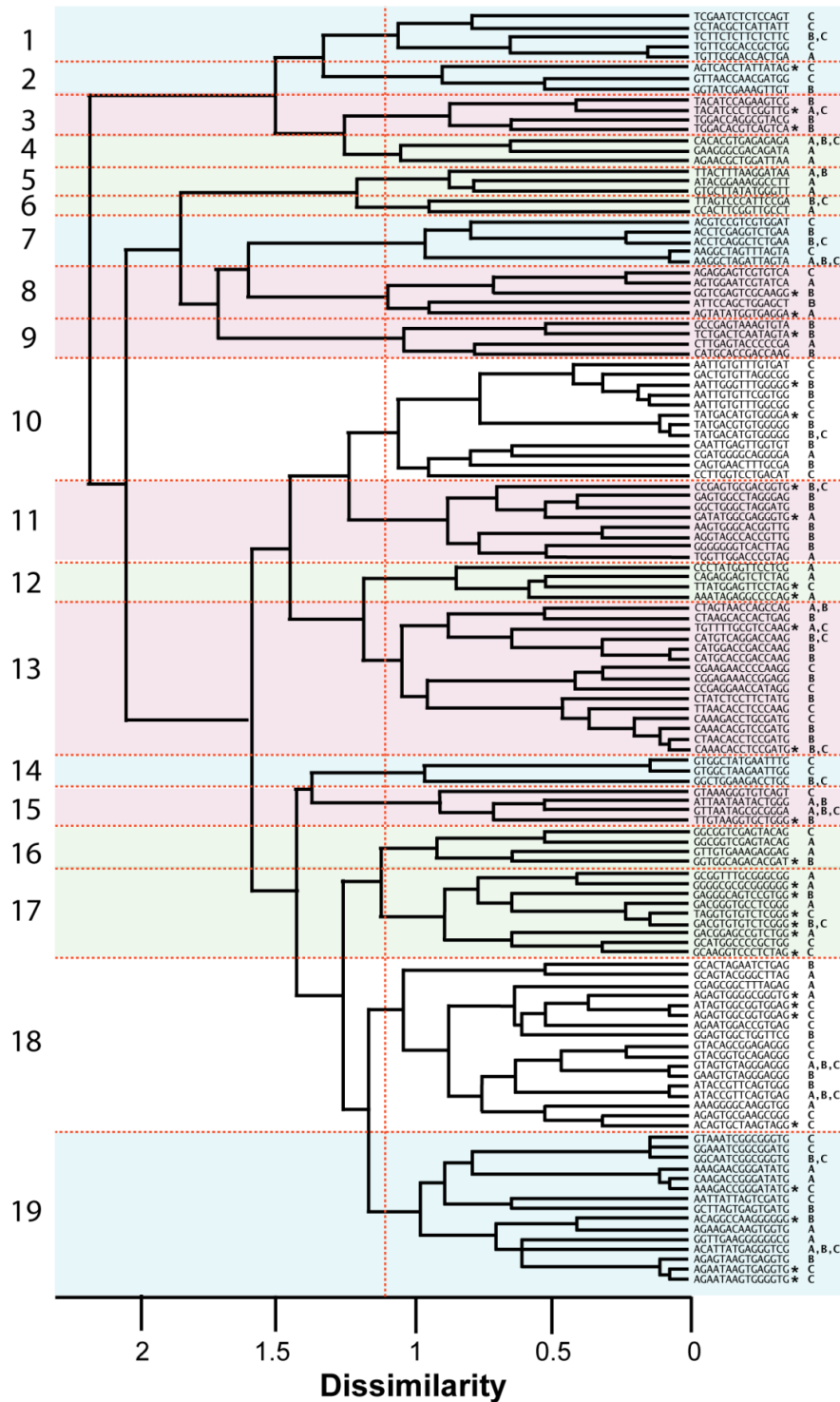


Figure 2.2. Hierarchical clustering of recovered ISREs indicates sequence composition correlates with sorted sections. Hierarchical clustering applied to 125 recovered ISRE sequences identified 19 clusters using a dissimilarity cutoff of 1.1. Clusters that have over

40% sequence representation from one sorted section are indicated (A, green; B, red; C, blue). Sorted sections for each sequence are denoted. Starred sequences were subjected to additional studies to examine regulatory activity.

2.2.3. GCCS clustering of recovered ISREs identifies motifs similar to known splicing factor binding sites

The sequence alignment results indicate that multiple motifs of varying lengths occur within the 15-nt dataset (Figure S2.5). For analyzing datasets of this nature, Graph Clustering by Common Substrings (GCCS)²¹ is better suited than hierarchical clustering. We analyzed a 19-nt region including the 15-mer ISRE sequence and 2-nt of the flanking regions for sequence enrichment. Since RNA binding proteins typically recognize short sequence motifs, we restricted our analysis to n-mers ranging from 4–6-nt. We determined the enrichment of n-mers in a sample of 125 sequences using a confidence interval for the binomial distribution based on probabilities expected for 19-nt oligonucleotides containing 15-nt of uniformly random bases flanked by the 2 constant bases present in the experimental system. In the ISRE dataset, 241 n-mers consisting of 39 4-mers, 93 5-mers, and 109 6-mers were significantly enriched ($\alpha_{1\text{-tailed}} = 0.1$; Figure 2.3a and Table S2.2). The GCCS analysis grouped 80.1% of the statistically enriched 4-6-nt n-mers into 30 consensus motif clusters (Figure 2.3b and Tables S2.3 and S2.4, Methods).

Many of the consensus motifs identified by the GCCS analysis resemble known binding sites for trans-acting splicing factors (Figure 2.3b and Table S2.5). The largest number of motifs resembles binding sites of the hnRNP family of proteins (class 1). In

particular, class 1 motifs resemble binding sites for several known repressors of splicing: hnRNP A1 (TAGGG)²², hnRNP F/H (GGGGG)²³, the polypyrimidine tract binding protein PTB (hnRNP I, CT-rich)²⁴ and hnRNP L (CA-rich)²⁵. The significant similarity between binding sites for the hnRNP family of proteins and the enriched motifs supports the possible functional role of selected ISREs.

Several identified motifs resemble known binding sites for the SR protein family (class 2) whose members act as general splicing factors¹. Enriched ISREs within class 2 resemble binding sites for SF2/ASF (GAAGAA)²⁶, SRp40 (ACAAG)²⁷, SRp30c (CTGGATT)¹⁴, SC35 (AGGAGAT)²⁸, 9G8 (GACC)²⁸, and Tra2 β (GAA)_n²⁹. While the examples of SR proteins involvement in splicing repression are limited, the enrichment of motifs similar to binding sites for members of this family suggests that their role in intronic regulation may be more widespread than previously thought.

Several of the enriched motifs identified in our dataset resemble the major 5' splice site (ss) consensus sequence GT[A/G]AGT (class 3)³⁰. All four motifs in class 3 contain an AGT core element, and the enriched motif TAAGTG is almost identical to the canonical 5' ss sequence and the hnRNP G binding motif AAGT³¹. The occurrence of 5' ss motifs within intronic regulatory elements has been noted³² and computational analyses have identified conserved elements that are similar to the consensus 5' ss within mammalian intronic regions^{21,33}. In addition, the enrichment of 5' ss motifs was previously observed in an *in vivo* screen for ESSs⁹. Taken together, these results add support to the role of cryptic 5' ss in regulating alternative splicing.

Other enriched motifs in our dataset (GTGT, GGTGG, TTGTGT, and GGTT) resemble known binding sites for the CELF/Bruno-like family (class 4). This family of

proteins regulates alternative splicing patterns by binding sequences that contain CTG repeats and exhibit a higher affinity for GT repeats³⁴. The motifs GTGT and TGTG resemble binding sites to a well-characterized member of this family, CUG-BP1, which has been shown to bind TGT-containing sequences³⁴. The GTGT motif may also serve as a binding site for hnRNP M³⁵.

GCCS identified 5 motifs that represent either novel regulatory elements or weak binding sites for characterized splicing factors (class 5). The [A/G]TGGC motif is similar to a degenerate CELF protein binding site and the motif TCGG[G/C] shares up to 80% sequence identity to a hnRNP A1 binding site. Strikingly, the GCTGG, CGA[T/G] and TATG motifs have not been previously identified. Therefore, in addition to identifying elements resembling binding sites for characterized trans-acting splicing factors, SPLICE generated novel regulatory elements.

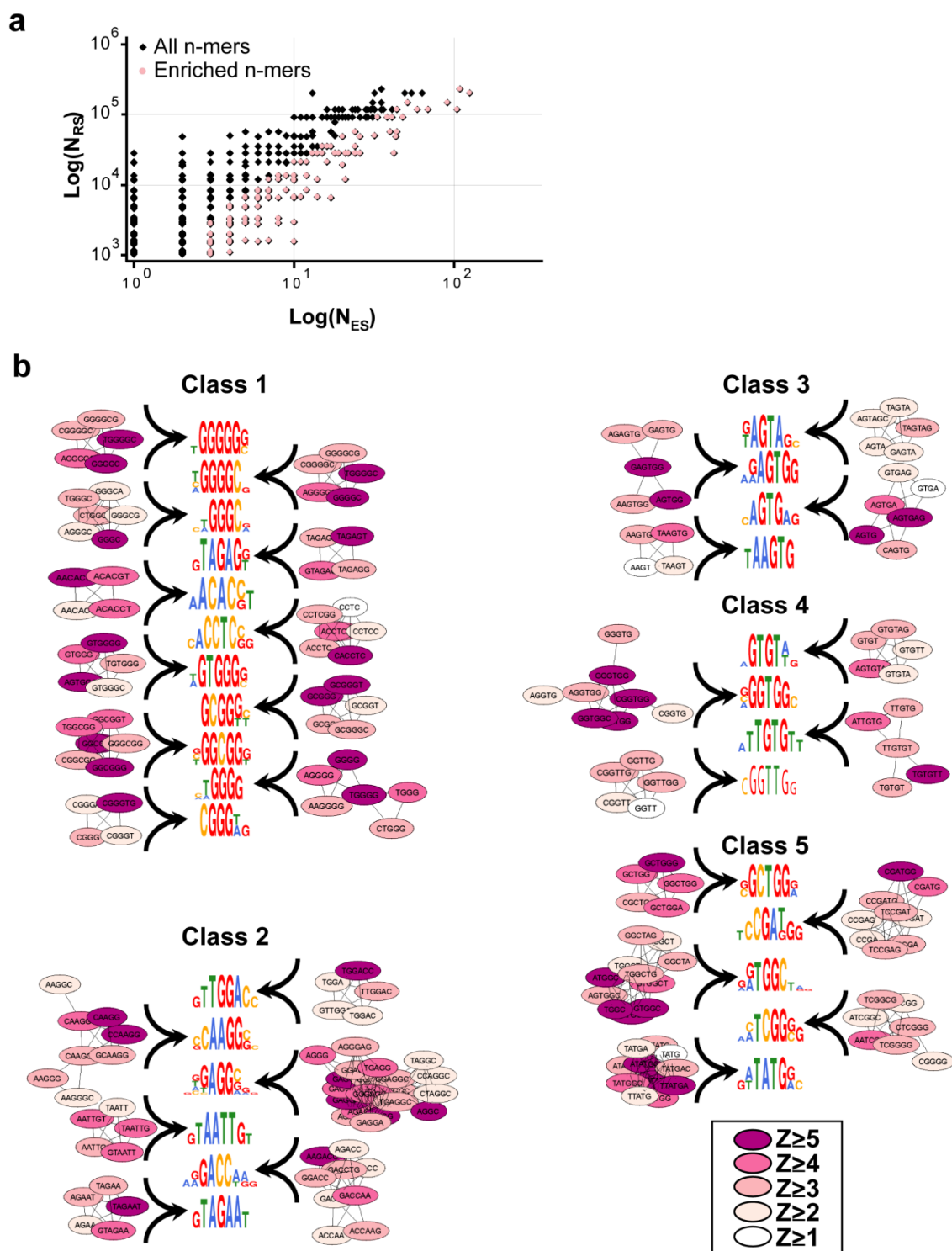


Figure 2.3. Enriched motifs and GCCS clusters derived from recovered ISRE sequences

map to known and unknown splicing factors. **(a)** Scatter-plot for the occurrence

frequency of all 4–6-nt n-mers in the enriched sample set (N_{ES}) vs. a corresponding random sample set (N_{RS}) (black). A similar scatter-plot based on n-mers determined to be significantly enriched in the recovered ISREs is overlaid (pink). **(b)** Consensus motif groupings according to resemblance to binding sites for trans-acting splicing factors. Motif classes include enriched ISREs that are similar to the binding sites of the hnRNP, SR and CELF families of proteins and the 5' ss (classes 1–4, respectively). Class 5 consists of previously unidentified elements and may represent novel regulatory sequences. The graph clusters representing the enriched n-mers used to construct each consensus motif are shown. Vertices are colored according to the enrichment Z-scores.

2.2.4. Enriched ISRE n-mers resemble known splicing regulatory elements

To investigate the potential functional role and general significance of the identified ISRE motifs, we examined the number of pentamer motifs identified in our enriched n-mer dataset (Table S2.2) that are identical to pentamers in published sets of splicing regulatory elements (SREs). We analyzed data corresponding to four SRE classes: ESEs^{6,7}, ESSs^{7,9}, ISEs¹⁶, and computationally identified conserved intronic elements. The latter class includes conserved intronic sequences (CISs)²¹, ISREs³³, pentamers enriched in intronic regions of excluded exons in neural progenitor (NP) cells³⁶ and motifs enriched upstream of weak polypyrimidine (PY) tracts in AT- and GC-rich introns³⁷. Significant overlap exists between the enriched pentamers and ESSs, ISEs, donor intronic (DI) elements in NP cells and motifs enriched upstream of weak PY tracts ($P < 0.05$ for ESSs, $P < 0.0001$ for ISEs, NP DI elements and weak PY tracts) (Figure 2.4a). The dominant motifs that overlap between SPLICE-generated pentamers and ISEs

and weak PY elements are G-rich elements, similar to hnRNP A/B and hnRNP F/H binding sites and the canonical 5' ss. The results suggest that the selected elements may function as general splicing silencers and intronic modulators of splicing (as both silencers and enhancers) depending on their context across various cell types and are likely regulated by general splicing factors. The comparison between enriched pentamers and conserved acceptor intronic elements (AI) for CIS and ISRE datasets demonstrate some overlap ($P < 0.05$). However, the observed overlap is far less than expected, suggesting that SPLICE selected against these elements.

2.2.5. Genome-wide analysis demonstrates that enriched ISREs associate with spliced exons

The biological relevance of selected motifs was examined by assessing the association of enriched motifs with naturally occurring alternative and constitutive splicing events. The occurrence of enriched motifs in the region 80-nt upstream of the AI regions flanking skipped exons was determined using a database of alternatively spliced junctions throughout the human genome²¹. A portion of SPLICE-generated ISREs significantly associate with alternative splicing (2 of 30; $P_{t\text{-test}} < 0.01$; class 4 only) and constitutive splicing (10 of 30; $P_{t\text{-test}} < 0.05$; all classes except 3) (Figure 2.4b). The entire population of consensus n-mers significantly associates with constitutive splicing ($P_{t\text{-test}} = 1.8e^{-8}$). This association is unexpected since selected ISREs are located within an alternatively spliced gene. However, the alternative exon 7 of the SMN1 mini-gene strongly favors inclusion, such that it may display regulatory signals similar to those involved in constitutive splicing, potentially biasing the sequence composition of selected

ISREs towards the association with constitutive splicing. In addition, our ISREs are likely to be more enriched in ISSs, which have been shown to be enriched in the intronic flanks of constitutively spliced exons². Results from our genome-wide association analysis suggest that selected ISREs serve an important role in defining constitutive and alternative splice sites.

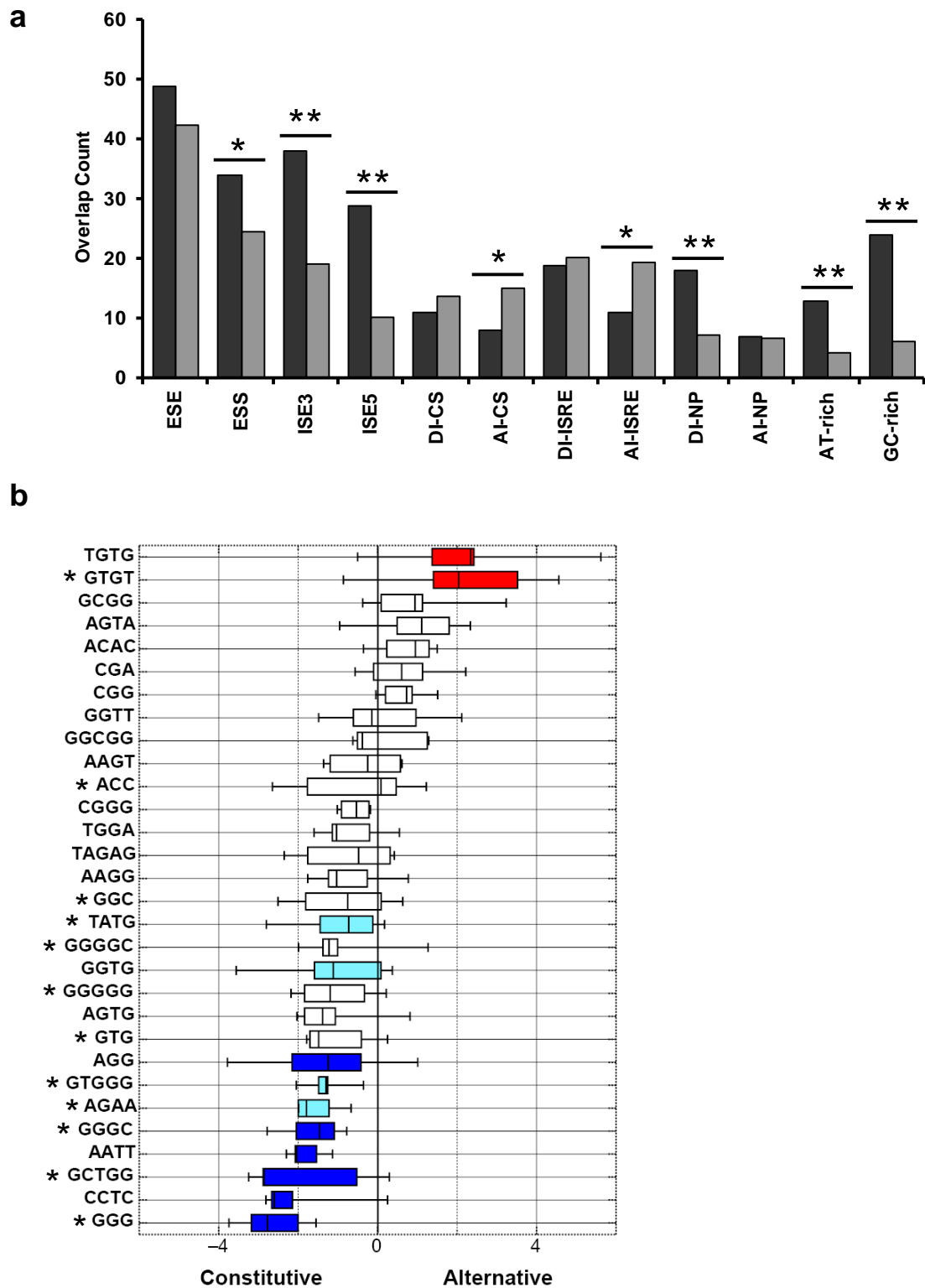


Figure 2.4. Enriched n-mers overlap with both experimentally and computationally derived SREs and associate with constitutive and alternative splicing. **(a)** Overlap of

enriched n-mers from recovered ISRE sequences with known classes of SREs. Observed (black) and expected overlap (gray) between datasets is shown. P -values derived from the chi-squared test of association are as follows: * $P < 0.05$ and ** $P < 0.0001$. **(b)** Box-plots revealing the distribution of TA-scores for GCCS derived ISREs. The GCCS consensus motifs that are significantly associated with alternative splicing are shown in red ($P_{\text{t-test}} < 0.01$) and those that are significantly associated with constitutive splicing are shown in shades of blue (dark blue, $P_{\text{t-test}} < 0.01$; light blue, $P_{\text{t-test}} < 0.05$). In total, 9 consensus motifs are biased toward alternative splicing and 21 consensus motifs display a bias towards constitutive splicing. Elements exhibiting no significant association with either category are not shaded. Starred motifs are present in hexamers subjected to RNAi silencing studies to examine regulated splicing.

2.2.6. Recovered ISRE sequences enable tuning of alternative splicing

The silencer activity of recovered 15-mers was validated by characterizing individual sequences that were selected randomly from 13 of the 19 hierarchical clusters (Figure 2.2). We analyzed an initial set of 18 15-mers (Figure 2.5a) and 4 known ISS sequences: an hnRNP H³⁸, 2 PTB^{39,40} and a U2AF65 binding sites³⁹. Individual sequences were cloned into our NMD-based reporter, stably transfected into HEK-293 FLP-In cells and analyzed by flow cytometry. Of the known ISSs tested, only the U2AF65 element demonstrates significant silencing activity relative to the NMD control, exhibiting an ~1.5-fold higher fluorescence level (Figure 2.5b). This result is in line with studies demonstrating that the silencing mechanisms of several characterized ISREs are context dependent⁴¹. In contrast, 16 of the selected sequences display significant silencer activity

($P \ll 0.001$) and 2 exhibit enhancer activity relative to the NMD control ($P < 0.05$) (Figure 2.5b), and over half exhibit silencing activities equal to or greater than the U2AF65 element. Similar trends were seen upon examination of an additional 12 recovered sequences (Figure S2.6b). In addition, we arranged the sequences into groups representing low (ISS1-5), medium (ISS6-10), and high (ISS11-16) silencing activities and determined the section from which each sequence was recovered (Figure 2.2). The activities of the majority of tested sequences correlated with sectioned populations, where a subset of enriched n-mers GGGGC, GGGC, and GGG correlated significantly with their sorted section ($P \ll 0.01$) and those sequences that did not correlate were shown to cluster with the appropriate group by sequence. These results support that functional regulatory activity is related to sequence.

To directly examine changes in splicing patterns, we analyzed the transcript isoforms of 12 of the recovered sequences and the ISS controls by qRT-PCR. The total transcript levels and the levels of intron retention for the examined ISS and control sequences were similar to the NMD control, while these levels for the selected ISEs differed from the NMD control ($P < .05$) (Figure S2.7a–d). The GFP-SMN1 control exhibits a low level of the skipped exon isoform compared to the NMD control ($P < 0.05$) (Figure 2.5c). As expected, for most of the recovered and control ISS sequences the levels of the skipped exon isoform are significantly higher than the NMD control ($P < 0.05$), with the exception of ISS15 ($P = 0.51$) and ISS8 ($P = 0.40$). In addition, the ISE sequences exhibited lower levels of the skipped exon isoform relative to the NMD control ($P < 0.05$). Therefore, expression levels of the skipped exon isoform generally confirm the activity of the sequences observed by fluorescence measurements.

All constructs except for the GFP-SMN1 control are expected to exhibit low levels of the exon 7 included isoform, as this isoform should be rapidly degraded through NMD. The GFP-SMN1 control exhibits a high level of exon 7 inclusion (99.7%), ~60-fold more than the NMD control. Exon inclusion levels for the ISS controls do not differ from the NMD control ($P > 0.35$), with the exception of PTB(2), which had a higher level of exon inclusion ($P < 0.005$). Exon inclusion levels for 8 of the 10 recovered ISS sequences (ISS5, ISS8-13) and the ISE sequences range from 2 to 20-fold less than the NMD control ($P < 0.05$), whereas ISS15 and ISS16 exhibited increased levels of the exon included isoform relative to the NMD control ($P < 0.05$, Figure S2.7e). The elevated levels of exon inclusion observed from several ISSs is not a result of cryptic splice sites as determined by analyzing the sizes of the RT-PCR amplification products (data not shown). Overall, the majority of sequences that display increased fluorescence have decreased levels of exon 7 inclusion compared to the NMD control, supporting their silencer function.

2.2.7. ISRE sequences function in a different cell type

The relative levels⁴² and activities⁴³ of trans-acting splicing factors vary widely across different cell types, which may result in *cis*-acting sequences exhibiting different regulatory activities. To determine whether the selected ISREs are cell type specific, we examined their regulatory function in a second cell line. We first examined the fluorescence of ISS1-16, ISE1, and the NMD and GFP-SMN1 controls in a transient transfection assay in the HEK-293 cell line to verify that regulatory activity was observed under these conditions. Flow cytometry analysis reveals that transiently transfected cells

display increased expression levels and population distributions relative to stable cell line assays (Figure S2.8). As such, the relative expression of the GFP-SMN1 control is only ~4.1-fold that of the NMD control (Figure 2.5d). Despite the decreased sensitivity of the transient transfection assay, the qualitative activity of 15 of the recovered ISREs was maintained and 11 sequences exhibited significantly increased expression ($P < 0.05$).

We next investigated whether the recovered sequences function in HeLa cells. The GFP-SMN1 control displays a approximate six fold higher level of expression than the NMD construct in HeLa cells in the transient transfection assay (Figure 2.5e). The ISRE sequences display a range of expression levels, but all are significantly different than the NMD control ($P < 0.05$). The majority of examined sequences (12 of 16) maintain the same trend in activity in HeLa cells as was observed in HEK-293 cells and ANOVA analysis of the activities in both cell lines shows a strong correlation ($P < 0.0005$). In contrast, four of the tested sequences (ISS3, 5, 6, and 13) exhibit enhancer activity relative to the NMD control in HeLa cells, which may be due to differences in levels of trans-acting factors between the cell lines. The results support that most sequences recovered from SPLICE retain function in a cell line different from which they were selected and may represent global splicing regulators.

2.2.8. Analysis of recovered ISRE sequences in a different transcript supports context dependent function

The context dependence of *cis*-regulatory elements on splice site choice has been shown⁴⁴ and we have observed little activity from known silencers in the context of the SMN1-NMD reporter system (Figure 2.5b). A subset of the selected ISREs was tested for

context dependent function by examining their activity in a second NMD-based reporter, based on the BRCA1 gene consisting of exons 17, 18, and 19⁴⁵, via transient transfection in HEK-293 cells. Selected ISRE sequences were inserted 50-nt upstream from the 3' ss of exon 18. Analysis of the reporter constructs by flow cytometry reveals a approximate two fold difference between the positive and negative controls ($P < 0.05$, Figure 2.5f). Only 3 of the tested sequences (ISS14, 17 and 18) exhibit significant silencer activity ($P < 0.05$) in the context of the BRCA1 mini-gene, while 1 sequence (ISS15) exhibits enhancer activity ($P < 0.05$). Transcript isoform analysis indicates that the level of exon 18 inclusion in the NMD control is ~12.5 fold less than the GFP-BRCA1 control (Figure S2.7f) and that splicing patterns for a subset of tested ISRE sequences were similar to the NMD control (Figure S2.7g), validating the lack of ISRE activity observed by fluorescence measurements. A predicted secondary structure analysis of the intronic regions shows that individual ISREs change the overall structure of each intron very little (Figure S2.9). However, the predicted secondary structure of the SMN1 intron is significantly different than that for the BRCA1 intron. The results suggest that the regulatory activity of SPLICE-generated ISREs is likely dependent on specific properties of the mini-gene in which they are selected.

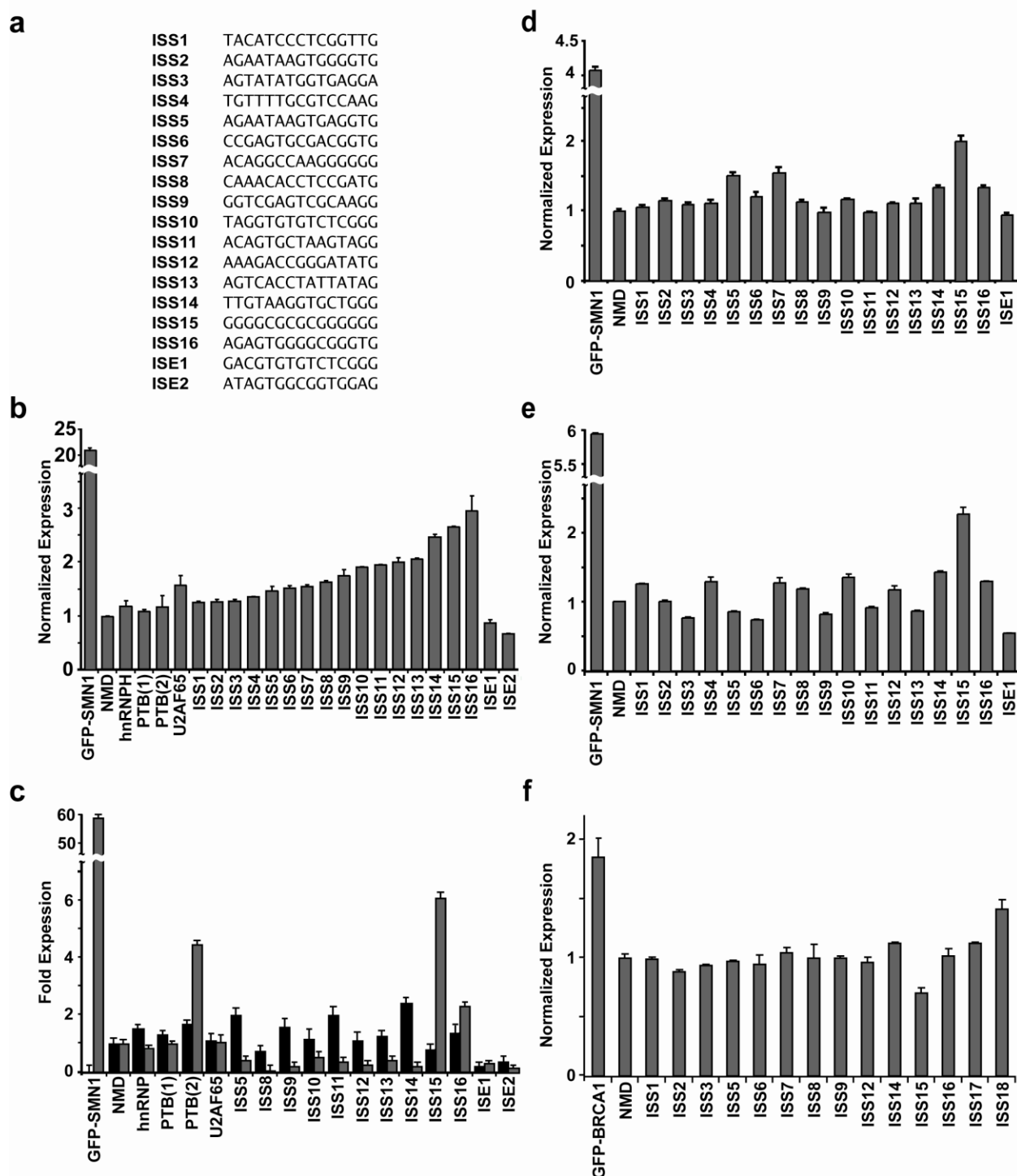


Figure 2.5. Functional analysis of recovered ISRE sequences. **(a)** Recovered ISRE sequences examined for regulatory activity. **(b)** Flow cytometry analysis of HEK-293

FLP-In stable cell lines generated for recovered ISRE sequences and control constructs. For all reported activities, the mean GFP levels from two independent experiments were determined and normalized to the NMD control. Normalized expression and average error are reported. ISRE sequences are labeled according to function. **(c)** qRT-PCR analysis of the ISS control sequences and 12 selected sequences with primer sets specific for exon 7 included (black bars) and excluded (gray bars) products. Expression levels of duplicate PCR samples were normalized to the levels of *HPRT*. Fold expression data is reported as the mean expression for each sample divided by the mean NMD expression value \pm the average error. **(d)** Flow cytometry analysis of recovered ISRE sequences and control constructs transiently transfected in HEK-293 cells. **(e)** Flow cytometry analysis of recovered ISRE sequences and control constructs transiently transfected in HeLa cells. **(f)** Flow cytometry analysis of recovered ISRE sequences and control constructs in the BRCA1 mini-gene transiently transfected in HEK-293 cells.

2.2.9. Analysis of enriched hexamers confirms independent and combinatorial function

We examined the silencer activities of representative hexamers from consensus motifs within the GCCS clusters through transient assays in HEK-293 cells to confirm the activity of individual motifs. Hexamers resembling the PTB, hnRNP H, SF2/ASF, Sam and the CELF protein binding sites, the 5' ss and an unknown motif were examined (classes 1–5; Figure 2.3b). Silencing activity was investigated by comparing expression levels of the hexamer alone to the hexamer with double point mutations (loss of function) and to the hexamer in duplicate (Figure 2.6a). A majority of the mutated hexamers (PTB, hnRNP H, Sam and unknown motifs) exhibits significant loss of function ($P < 0.05$;

classes 1, 2, and 5). However, only one of the hexamers (class 4) displays an increase in silencer activity when present in duplicate. The results indicate that while individual hexamers exhibit silencing activity and likely represent core ISREs, they do not necessarily behave in an additive manner likely due to context and spacing requirements. For example, the duplicate hnRNP H hexamer does not exhibit increased silencing, whereas ISS15, which differs from the duplicate hexamer by 3 cytosine residues, exhibits strong silencer activity. The additional residues may provide spacing between the G-rich hexamers important for functional activity.

To test the possibility of combinatorial control within the context of a selected 15-mer, we examined two ISS sequences that contained multiple enriched hexamers. Most of the recovered ISRE sequences contain several enriched n-mers, where 88% of all extended 15-mers (plus 2-nt flanking region) contain at least one enriched hexamer (Table S2.6). ISS5 contains 7 enriched hexamers resembling the 5' ss and binding sites for the CELF and SF2/ASF proteins, which cluster into 3 main regions within the sequence (Figure 2.6b). We introduced 2 point mutations within each region and in combination and assessed their activity through transient transfection assays in HEK-293 cells. Mutations within each region of ISS5 decrease expression to levels comparable to the NMD control, indicating that each silencing zone has regulatory activity. Region 1 contains two overlapping hexamers resembling the SF2/ASF binding motif, where one of these is the SF2/ASF representative hexamer that did not demonstrate silencing activity in the hexamer analysis studies, suggesting that the regulatory function of this hexamer is context dependent (Figure 2.6a). Simultaneous mutations to regions 2 and 3 resulted in expression levels comparable to or slightly higher than the individual mutations ($P <$

0.05), indicating that the regulatory function of ISS5 is likely not due to combinatorial recognition of motifs.

In contrast, analysis of a two-zone ISS sequence, ISS8, demonstrated that enriched hexamers can exhibit combinatorial control over ISS activity (Figure 2.6c). The extended ISS8 sequence contains 8 enriched hexamers resembling preferred binding sites for the PTB and hnRNP L proteins and a novel element that overlaps regions 1 and 2. A similar analysis of ISS8 shows that the individual mutations within each zone disrupt silencer activity to levels below the NMD control ($P < 0.05$), resulting in an ~18% decrease in activity (Figure 2.6c). Mutations to both regions in combination result in an ~25% decrease in silencer activity, suggesting that the hexamer regions work together to effect silencer activity ($P < 0.005$). Therefore, the ‘zones of silencing’ in our recovered ISRE sequences consisting of enriched hexamers exhibit regulatory function independently and in combination with other zones, but the effects are context dependent and may depend on the specific trans-acting factors involved.

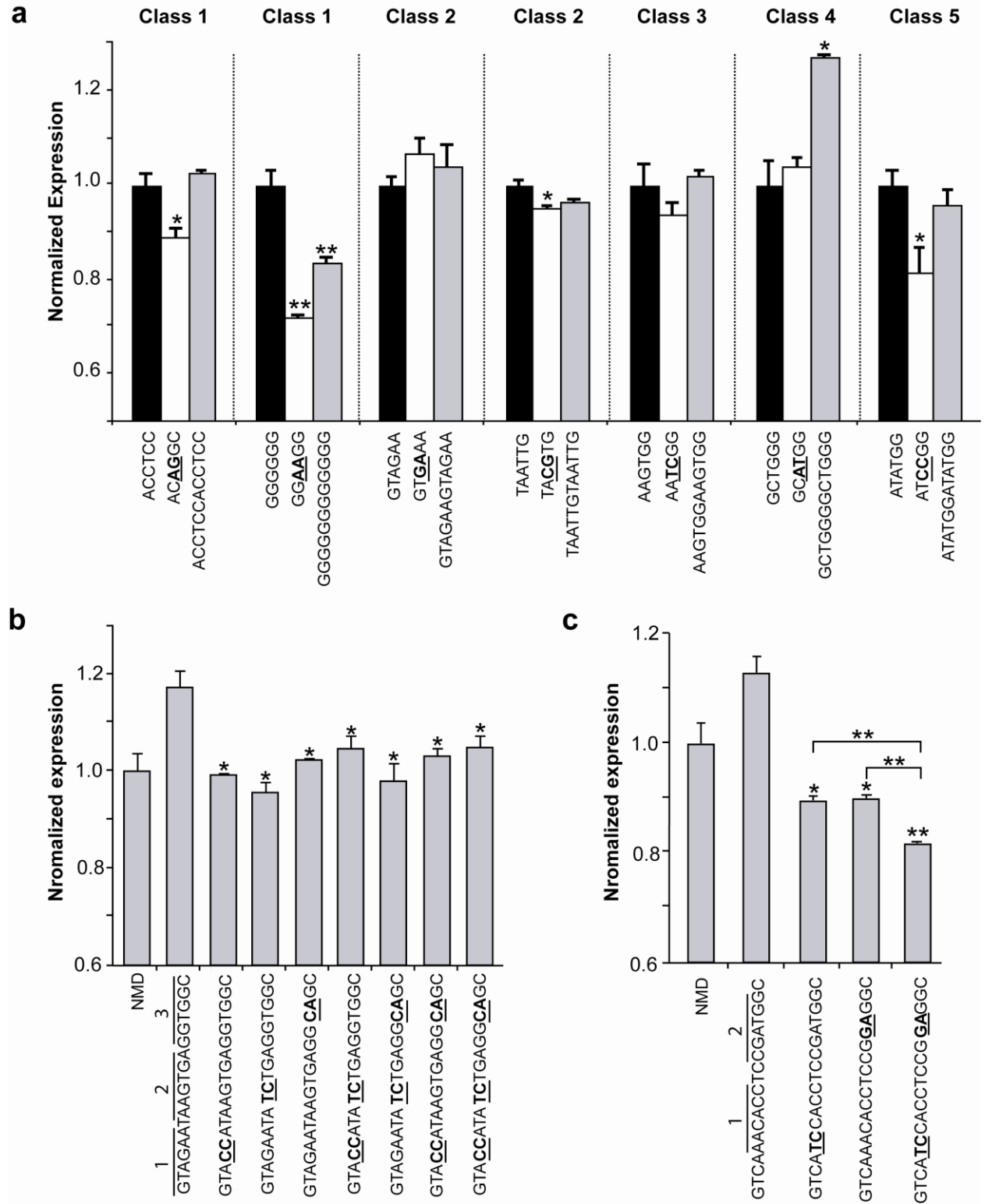


Figure 2.6. Enriched ISRE hexamers demonstrate silencer activity. **(a)** Individual hexamer analysis supports the silencing activity of enriched hexamer sequences.

Representative hexamers from each class of GCCS clusters and corresponding mutant and duplicate sequences were characterized in transient transfection assays in HEK-293 cells. For all reported data, silencing activity was assessed by flow cytometry analysis, where the mean GFP levels from two independent experiments were normalized to the wild-type hexamer construct. Normalized expression and average error are reported. P - values derived from the Student's t-test are as follows: * $P < 0.05$ and ** $P < 0.01$. **(b)** Mutational analysis of an ISS sequence supports the silencing activity of individual hexamer regions. The combined and individual activity of hexamer regions within the context of an ISS sequence was examined by introducing 2 point mutations into 3 regions, in combination and separately into all 3 hexamer regions of ISS5. **(c)** Mutational analysis of an ISS sequence supports the silencing activity of combined hexamer regions. The combined and individual activity of hexamer regions within the context of an ISS sequence was examined by introducing 2 point mutations separately and in combination into 2 hexamer regions of ISS8.

2.2.10. Splicing factor depletion influences ISRE regulated splicing in vivo

Many of the sequence classes identified by GCCS analysis resemble known or predicted binding sites for trans-acting splicing regulators. To validate the functional significance of the GCCS-identified sequence classes and uncover the associated trans-acting factors or SRNs, we screened a panel of siRNAs targeting known splicing regulators (hnRNP H, hnRNP A1, PTB, CUG-BP1, and SF2/ASF) for resulting effects on the splicing patterns of selected hexamers in stable cell lines. RNAi-mediated silencing of each gene resulted in a substantial reduction ($\geq 70\%$) of the targeted protein (Figure 2.7a)

and displayed minimal effects on the other splicing factors examined (Figure S2.10a). Hexamers from classes 1 (ACCTCC, GGGGGG), 2 (GTAGAA), 4 (GCTGGG) and 5 (ATATGG), which harbor potential binding sites for the selected trans-acting factors, and a random insert control were subjected to the RNAi-based screen using a mini-gene lacking a PTC to avoid any siRNA-mediated effects on the NMD pathway.

We analyzed the splicing patterns of the hexamer and control constructs through qRT-PCR analysis. In the presence of the mock siRNA, four of the hexamers exhibit silencing activity, a higher ratio of exon exclusion to inclusion, relative to the GFP-SMN1 control (Figure 2.7b and Figure S2.10b). In contrast, one hexamer (GGGGGG) exhibits enhancer activity in the presence of the mock siRNA. Splicing of constructs containing the GCCS hexamers were significantly affected by the depletion of at least one, and in some cases multiple, trans-acting factors (Figure 2.7c and Figure S2.10c). In contrast, siRNA-mediated depletion of the selected trans-acting factors had statistically insignificant effects ($P \gg 0.05$) on the splicing pattern of the GFP-SMN1 control.

The splicing pattern of three hexamer constructs exhibited significant changes in response to the depletion of one of the trans-acting factors. The GGGGGG enhancer motif matches the hnRNP F/H binding site²³. Depletion of hnRNP H leads to a 2.3-fold increase in exon exclusion, demonstrating that this factor enhances the recognition of the 3' ss, most likely through direct binding to the hexamer. The ACCTCC motif is similar to the CT-rich PTB binding site²⁴, but depletion of PTB leads to only a marginal decrease in exon exclusion of the construct. However, depletion of CUG-BP1 leads to a significant two fold decrease in exon exclusion. Although PTB has been shown to act antagonistically to CELF proteins¹², it is unlikely that CUG-BP1, which binds CTG and

GT-rich motifs, directly binds to the ACCTCC hexamer, suggesting it may be recruited through interactions with other regulatory proteins. Analysis of the splicing of a novel motif, ATATGG, reveals that depletion of hnRNP A1 leads to an increase in exon exclusion levels (~2.3-fold). The hexamer and flanking regions contain a GGG motif that may be a weak binding site for hnRNP A1. However, any direct binding of hnRNP A1 likely competes with other regulatory factors since its depletion leads to an increase in exon exclusion. Alternatively, modulation of hnRNP A1 levels may affect the levels of other trans-acting factors that play a role in the splicing regulatory effect of the hexamer.

Two hexamers constructs exhibit significant changes in their splicing pattern in response to depletion of multiple factors. The GTAGAA motif closely resembles the SF2/ASF SELEX-derived binding site (GAAGAA)²⁶, although the hexamer and flanking regions contain two GT repeats, which may serve as binding sites for CUG-BP1, and a TAGA motif, which may be a weak binding site for hnRNP A1. Depletion of hnRNP H, hnRNP A1, CUG-BP1 and SF2/ASF led to a 2.5-fold or greater reduction in exon exclusion levels for the construct. One possible mechanism is that SF2/ASF, CUG-BP1 and hnRNP A1 directly compete for binding to the GTAGAA hexamer and that hnRNP H acts positively in the recruitment of these factors. Both hnRNP H and CUG-BP1 have been shown to form an RNA-dependent suppressor splicing complex⁴⁶, suggesting that many of these factors may be involved in an inhibitory splicing complex that aids in the recruitment of a factor that directly binds to the transcript. The GCTGGG motif and flanking regions contain GT and TG dinucleotides and a CTG element that resemble CUG-BP1 binding sites. Depletion of CUG-BP1 results in a four fold decrease in exon exclusion of the construct. Depletion of PTB and SF2/ASF also cause significant

decreases in exon exclusion, although the preferred binding sites of these factors don't resemble any motifs within the GCTGGG hexamer and flanking regions. The results suggest that CUG-BP1 may be directly involved in binding to the GCTGGG hexamer, while PTB or SF2/ASF may be recruited by CUG-BP1 or other trans-acting factors.

2.2.11. Splicing factor depletion alters splicing of endogenous genes containing ISREs

Our genome-wide analysis revealed that selected ISREs are enriched in the intronic regions flanking constitutively and alternatively spliced endogenous genes (Figure 2.4b). To determine the biological significance of these associations, we analyzed the splicing patterns of 10 alternatively spliced endogenous genes containing an intronic hexamer under depletion of trans-acting splicing factors (Figure 2.7d and Table S2.10). Each target gene was analyzed through qRT-PCR in the presence of a mock siRNA and a siRNA targeting the splicing factor that showed the most significant effect on the splicing pattern of each hexamer in the SMN1 mini-gene depletion studies (Figure 2.7c). We observed significant changes in the alternative splicing patterns of all targeted genes upon splicing factor depletion ($P < .05$), where 7 of the 10 genes showed increased exon inclusion supporting the ISS activity of the selected hexamers. In contrast, 3 genes (*RREB1*, *CAMK2G*, and *HNRNPA2B1*) displayed higher levels of exon exclusion, indicating that hexamers GTAGAA, GCTGGG, and ATATGG can function as ISEs within endogenous genes. The significant changes in splicing of the synthetic SMN1 mini-gene and endogenous genes containing selected hexamers upon splicing factor depletion support the functional role of SPLICE identified ISREs through known trans-acting factors. These studies further highlight the context dependent nature of ISRE

function, where a given sequence can display enhancer and silencer functions in different transcripts.

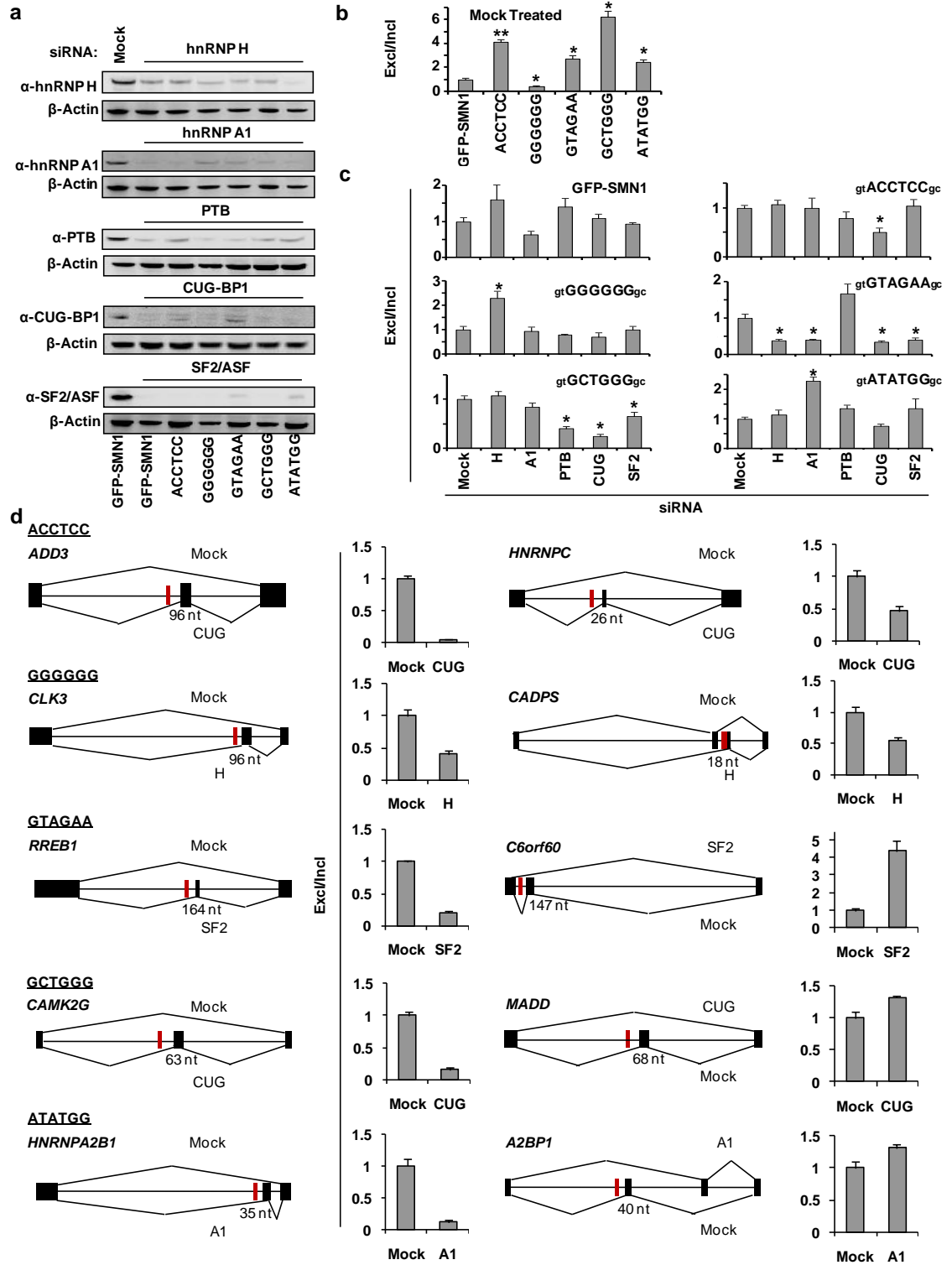


Figure 2.7. The effects of *in vivo* depletion of splicing factors on ISRE regulated splicing patterns of synthetic and endogenous genes. **(a)** Western blot analysis of total cell lysates prepared from the ISRE hexamer and GFP-SMN1 control cell lines treated with siRNAs targeted to trans-acting splicing factors and a mock siRNA negative control. β -Actin was used as a loading control for all blots. The results of the GFP-SMN1 mock treated lysate is representative of all mock treated cell lines. **(b)** qRT-PCR analysis of the mock treated ISRE hexamer and GFP-SMN1 control cell lines with primer sets specific for exon 7 included and excluded products. Expression levels of duplicate PCR samples were normalized to the levels of *HPRT*. Data is reported as the ratio of the mean expression of the exon excluded isoform to the exon included isoform normalized to the ratio for the GFP-SMN1 control \pm the average error. **(c)** qRT-PCR analysis of the siRNA treated ISRE hexamer and GFP-SMN1 control cell lines with primer sets specific for exon 7 included and excluded products. Data is reported as the ratio of the mean expression of the exon excluded isoform to the exon included isoform normalized to the ratio for the mock siRNA treated cell line control \pm the average error. *P* -values derived from the Student's t-test are as follows: * *P* < 0.05 and ** *P* < 0.01. **(d)** qRT-PCR analysis of the siRNA treated GFP-SMN1 control cell lines with primer sets specific for exon included and excluded products of 10 endogenous genes. The splicing patterns of each gene are diagrammed where black bars represent exons and red bars represent the location of conserved ISRE hexamer motifs. Data is reported as the ratio of the mean expression of the exon excluded isoform to the exon included isoform normalized to the ratio for the mock siRNA treated cell line control \pm the average error.

2.2.12. Models of ISRE mediated regulation of alternative splicing

Based on the location of our ISREs in the mini-gene construct, these elements likely exhibit their regulatory activity through interacting with trans-acting factors that enhance or inhibit the binding of general splicing factors, such as U2AF65 or the U2 snRNP complex, at the 3' ss. Based on our studies, we propose three models for ISRE regulation of alternative splicing based on direct or indirect interactions with trans-acting factors in the SRN. The first model suggests that direct binding of a specific factor or the competitive binding of multiple factors to the ISRE sequence plays a role in splicing regulation and is supported by results from the class 1 (GGGGGG), 2 (GTAGAA) and 5 (ATATGG) ISREs (Figure 2.8a). This model is further supported by a recent study describing the juxtaposition of an ESS and ESE that results in hnRNP H and F competing for binding with SF2/ASF⁴⁷. A second model is based on results from the class 1 (ACCTCC), 2 (GTAGAA), and 4 (GCTGGG) ISREs and proposes that direct binding of a specific factor is involved in the extensive recruitment of or is itself recruited by several other regulatory factors, thereby resulting in a recruitment pathway for ISRE regulation (Figure 2.8b). Previous work has suggested that splicing factors may be components of larger regulatory complexes in which binding selectivity is dictated by protein-protein interactions^{47,48}. The third model is supported by results from the class 5 (ATATGG) ISRE and is based on an interaction with an agonist factor, where the level of that factor may affect the levels of other splicing factors that play a role in the regulation of the ISRE (Figure 2.8c).

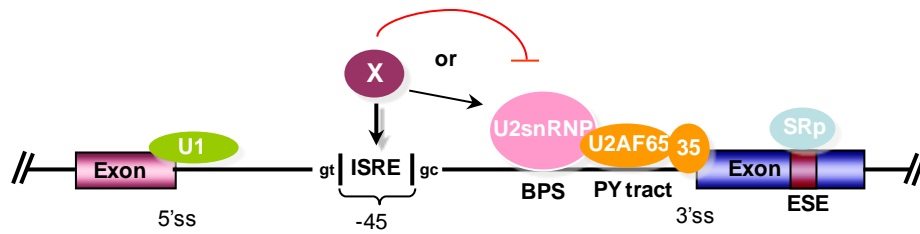
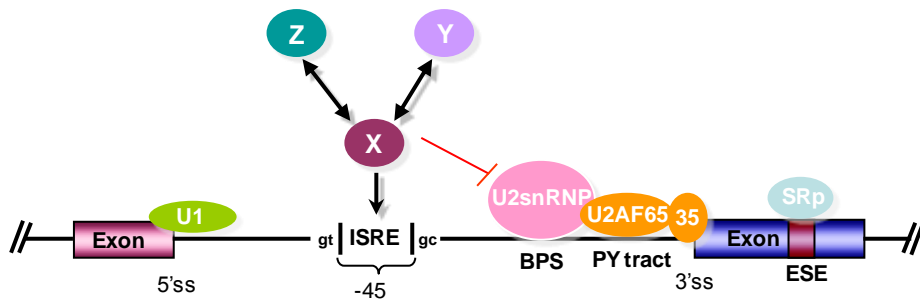
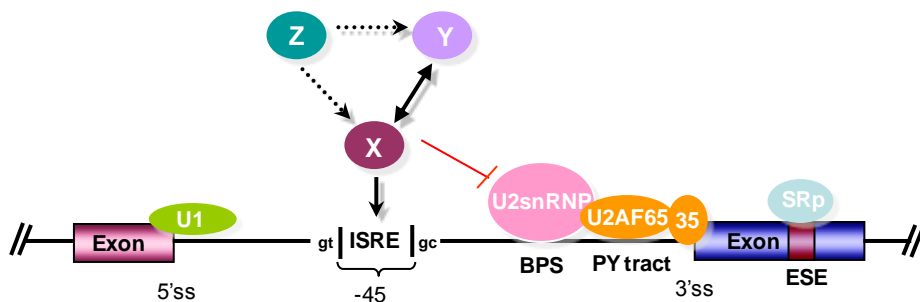
a**Direct/Competitive Binding : Classes I, II and V****b****Recruitment Pathway: Classes I, II and IV****c****Agonist Interaction: Class V**

Figure 2.8. Three models for ISRE regulation. **(a)** A direct binding / competition model for ISRE regulation. Factor X, represents a protein that directly binds to the ISREs, factors Y and Z represent proteins that do not directly bind to the ISRE sequence. We present our ISRE mechanistic models based on a simplified model of ESE-dependent 3'

ss activation by SR proteins involving the recruitment of U2AF65 to the PY tract and subsequent recruitment of the U2 snRNP complex. Selected ISEs stabilize the interactions of general splicing factors to the 3' ss, whereas ISSs destabilize this process. In the direct binding / competition model, factor X binds to the selected ISRE and either enhances (black) or suppresses (red) recognition of the 3' ss. **(b)** A recruitment pathway model for ISRE regulation. This model is based on the direct binding of factor X to the ISRE leading to the recruitment of additional regulatory factors Y and Z. **(c)** An agonist interaction model for ISRE regulation. This model is based on an interaction between the splicing factors and an agonist factor Z, where the level of factor Z may affect the levels of the splicing factors (X and Y) that play a role in the regulation of the ISRE.

2.3. Discussion

Significant advances in our understanding of the mechanisms that guide splice site selection and the distributions of regulatory elements has aided the formulation of an early version of a 'splicing code'⁹. However, currently missing from this draft is a thorough understanding of the sequence characteristics and function of ISREs. The ISREs obtained from our *in vivo* selection provide a diverse composition of sequences that correlate to an array of splicing activities, which enable tuning of alternative splicing, and reveal motifs that resemble binding sites for many known and novel trans-activating factors that are enriched in introns of naturally-occurring spliced genes throughout the human genome. The identified motifs offer a rich dataset to expand the splicing code, determine the extent of single nucleotide polymorphisms (SNPs) that modulate splicing

through ISREs and refine bioinformatic search algorithms for genome-wide identification of intronic regulators, which will facilitate the diagnosis and treatment of disease.

The majority of the tested ISREs retained function when tested in a second cell type, but not a second transcript. These results support that although differences in cellular constituents can lead to differential splicing patterns (i.e., tissue specific splicing), a dominant factor guiding function is likely the immediate transcript context in which these cellular factors bind. Therefore, refinements of the splicing code will likely benefit from identification of co-regulatory sequences that may be identified in functional screens examining pair-wise or combinatorial motifs. Such combinatorial motifs were identified by SPLICE and experimental characterization supports the properties of context dependent and combinatorial regulation. The results from our RNAi silencing study highlight the complexity of the SRNs associated with ISRE function, suggesting a role for multiple splicing factors influencing regulation at a single motif and supporting models for ISRE function in which direct binding, recruitment or agonist interactions with upstream factors interfere or enhance the recruitment of the basal splicing complex. The splicing factor depletion studies also provided experimental validation of ISRE regulation of endogenous alternatively spliced transcripts further supporting their biological significance and context dependent function. Our work sets the stage for larger-scale characterization studies of the identified ISREs and associated trans-acting factors, which will further elucidate ISRE regulatory activity and mechanism, including the role of combinatorial control and sequence context in the function of these elements. Our results provide the first large scale analysis of ISREs *in vivo* and highlight that an

understanding of the complex interplay between multiple factors at a single binding site is necessary to further define the splicing code.

2.4. Materials and Methods

2.4.1. Base *SPLICE* constructs

Plasmids were constructed using standard molecular biology techniques⁴⁹. All enzymes, including restriction enzymes and ligases, were obtained through New England Biolabs unless otherwise noted. DNA synthesis was performed by Integrated DNA Technologies, Inc. Ligation products were electroporated into *E. coli* DH10B (Invitrogen) using a GenePulser XP system (BioRAD), and clones verified through colony PCR and restriction mapping. All cloned constructs were sequence verified through Laragen. Primer sequences and plasmid descriptions are available in Tables S2.7 and S2.8, respectively.

The GFP-SMN1 mini-gene fusion construct was constructed through a PCR assembly and site-directed mutagenesis strategy. A region encompassing exons 6 through 8 of the *SMN1* mini-gene was amplified through PCR from template pCISMNxΔ6-wt⁵⁰ with primers Ex6 and Ex8 and PfuUltra high-fidelity DNA polymerase (Stratagene). The *GFP* gene was amplified from the template pKW430⁵¹ with primers GFP1 and GFP2. The GFP-SMN1 gene fusion was constructed by performing PCR assembly on the resulting purified products (Qiagen) as templates and flanking primers GFP1 and Ex8. The resulting gene fusion product was digested with Xho I and Kpn I and ligated into the corresponding restriction sites of the mammalian expression vector pcDNA5/FRT (Invitrogen), resulting in the positive control vector pCS238. A PTC (TAA at position +1

in exon 7) and ISRE insertion sites Eco RV/Cla I in intron 6 (positions -62 and -51 from 3' ss of exon 7, respectively) and Bam HI/Pml I in intron 7 (positions +43 and +59 from 5' ss of exon 7, respectively) were introduced by site-directed mutagenesis with primers ECmutF1/ECmutR1, PmlImutF/PmlImutR, and BamHIImutF/BamHIImutR using a Quickchange II Kit (Stratagene) according to manufacturer's instructions, resulting in the base NMD reporter construct pCS516. ISRE sequences were digested and ligated into the Eco RV and Cla I restriction sites within intron 6 of the base NMD construct.

The GFP-BRCA1 mini-gene fusion was constructed through a PCR assembly and site-directed mutagenesis strategy. A portion of the wild-type human *BRCA1* gene was amplified from HEK-293 genomic DNA as previously described⁴⁵ using reported primers P2, P3, P4, and P5 with the exception of the forward primer for exon 17 (Ex17) and the reverse primer for exon 19 (Ex19). The resulting wild-type BRCA1 mini-gene contains shortened introns and wild-type exons 17, 18, and 19. The *GFP* gene was amplified from template pKW430⁵¹ with primers GFP1 and GFP3. The GFP-BRCA1 gene fusion was PCR assembled using primers GFP1 and Ex19, digested with Xho I and Kpn I and ligated into the corresponding restriction sites of pcDNA5/FRT, resulting in the GFP-BRCA1 positive control construct pCS990. A PTC (TAA at position +3 in exon 18) and ISRE insertion sites Eco RV/Cla I in intron 17 (positions -61 and -50 from 3' ss of exon 18, respectively) were introduced as described above with primers ECmutF2 and ECmutR2, resulting in the base BRCA1-NMD reporter construct pCS1008. ISS sequences were cloned into intron 17 of the base NMD construct as described above.

2.4.2. Cell culture, transfections, stable cell lines and FACS

HEK293 FLP-In cells (Invitrogen) were cultured in D-MEM supplemented with 10% fetal bovine serum (FBS) and 100 $\mu\text{g/ml}$ Zeocin at 37°C in 5% CO_2 . HeLa cells were cultured in MEM media supplemented with 10% FBS. Transfections for all cell lines were carried out with Fugene (Roche) according to the manufacturer's instructions. All cell culture media was obtained from Invitrogen.

HEK-293 FLP-In stable cell lines were generated by co-transfection of the appropriate SMN1 mini-gene construct with a plasmid encoding the Flp recombinase (pOG44) in growth medium without Zeocin according to the manufacturer's instructions (Invitrogen). The library stable selections were carried out in 225 cm^2 flasks containing $\sim 4 \times 10^7$ HEK-293 FLP-In cells where 37 μg of pOG44 and 3.7 μg of the SMN1 ISRE plasmid library (10:1 ratio) were co-transfected. Fresh medium was added to the cells 24 h after transfection. The cells were expanded by a 1:4 dilution and Hygromycin B was added to a final concentration of 200 $\mu\text{g/ml}$ 48 h after transfection. In total, $\sim 450,000$ stable transformants were pooled from 60 transfections. Clones were harvested by trypsinization, pooled and analyzed on a FACS Aria (Becton Dickinson Immunocytometry Systems) 10–14 days after transfection. GFP fluorescence was excited at 488 nm and emission was measured with a FITC filter. Detailed sorting procedures are presented in Figure S2.2. In the first screening round, positive cells were bulk sorted into 96-well plates, where no more than 25,000 cells were collected into a single well. After ~ 1 –2 weeks of growth, positives were re-sorted into 3 fractions (A, B, and C) based on varying fluorescence levels (Figure S2.2b). Positive cells were bulk sorted in the second screening round as described for the first round. Total genomic DNA from bulk sorted

cells was purified using the DNeasy Blood & Tissue total DNA purification kit (Qiagen) according to the manufacturer's instructions and used as a template for amplification of recovered ISRE sequences with primers Lib3 and Lib4. The recovered ISRE fragments were then digested, ligated into the corresponding sites of pCS516 and sequenced verified by Functional Biosciences, Inc.

For transient transfection studies, HEK293 and HeLa cells were seeded in 12-well plates at $\sim 5 \times 10^4$ cells per well 16 to 24 h prior to transfection. Cell lines were transfected with 625 ng of the appropriate GFP-SMN1 or GFP-BRCA1 mini-gene constructs. The cells were harvested by trypsinization, pooled and analyzed on a FACS Aria 48 h after transfection. Experiments were carried out on different days and transfections were completed in duplicate, where the mean GFP fluorescence of the transfected population and the average error between samples is reported. A comparison of FACS gating procedures used in transient and stable assays is presented in Figure S2.8. Cell lines harboring PTC-containing transcripts tend to increase in fluorescence at higher passage numbers (>10), whereas the GFP-SMN1 cell line does not. As such, the fluorescence levels of the enriched cell populations at the time of sorting (Figure 2.1b) do not directly match the expression levels for individual recreated clones (Figure 2.5b). To minimize differences in expression due to such instabilities, the analysis of all stable cell lines was performed at an identical, early passage.

2.4.3. *qRT-PCR analysis*

Total cellular RNA was purified from stably transfected HEK-293 Flp-In cells using GenElute mammalian total RNA purification kit (Sigma) according to the manufacturer's instructions, followed by DNase treatment (Invitrogen). cDNA was synthesized using Superscript III reverse transcriptase (Invitrogen) according to the manufacturer's instructions. qRT-PCR analysis was performed using isoform-specific primers (Tables S2.9 and S2.10). Expression levels of duplicate PCR samples were normalized to the levels of *HPRT* (Hypoxanthine-guanine phosphoribosyltransferase). Fold expression data is reported as the mean expression for each sample divided by the mean NMD expression value \pm the average error.

2.4.4. *siRNA mediated silencing of trans-acting splicing factors*

siRNAs targeting hnRNP H, hnRNP A1, PTB, CUG-BP1, and SF2/ASF and a mock control siRNA were purchased from Dharmacon and are listed in Table S2.11. All duplexes were resuspended in 1X PBS to a concentration of 20 μ M. Briefly, HEK-293 FLP-In cells were plated at $\sim 2 \times 10^5$ cells per well in 6-well plates. After 24 h, the cells were transfected with individual siRNA duplexes to a final concentration of 50 nM using Lipofectamine RNAiMAX (Invitrogen) according to the manufacturer's instructions. Cells were collected for RNA isolation and western blotting 48 h after transfection.

2.4.5. *Western blot analysis*

Whole-cell extracts were prepared from harvested cells using M-PER mammalian protein extraction reagent (Pierce) and equal amounts of protein (50 μ g) were resolved on

4-12 % SDS-PAGE gels (Invitrogen) and transferred onto Protran nitrocellulose membranes (Whatman) using the Trans-Blot SD semi-dry transfer cell (BioRad). After blocking with 5% BSA in TBST, the membranes were incubated with the specified antibodies overnight at 4°C. After incubation, the membranes were washed with TBST and then incubated with the corresponding secondary antibody conjugated with HRP. Signals were detected using the ECL western blotting substrate (Thermo Scientific) according to the manufacturer's protocol. The primary antibody dilutions were 1:500 for goat anti-hnRNP H (N-16), 1:1000 for goat anti-Actin (I-19), 1:200 for goat anti-hnRNP A1 (Y-15), 1:200 for mouse anti-PTB (SH54), 1:200 for mouse anti-SF2/ASF (96) and 1:200 for mouse anti-CUG-BP1 (3B1). The secondary antibody dilutions were 1:10,000 for donkey anti-goat IgG-HRP (sc-2020) and 1:10,000 for goat anti-mouse IgG-HRP (sc-2005). All of the antibodies were purchased from Santa Cruz Biotechnology Inc. The relative band intensities were measured by densitometry analyses using Quantity One (BioRad).

2.4.6. Discovery of sequence motifs enriched in ISRE sequences

A sliding-window count of all n-mers (4–6-nt) within the nonredundant sample set of 125 sequences was performed. Two nucleotides flanking the 5' and 3' ends of the random region were included to account for bias due to the constant sequences. A similar sliding-window count on a set of 450,000 computer generated sequences containing a uniformly random 15-nt region flanked by the same constant nucleotides was performed to calculate the maximum likelihood probabilities for expected occurrences (see

Methods). For both data sets the counts were transformed into probabilities and the enrichment was determined according to the binomial confidence interval method²¹.

2.4.7. Overlap of ISRE sequences with known splicing regulatory elements

The set of pentamers enriched in the ISRE sequences were compared to previously compiled lists of ESEs^{6,7}, ESSs^{7,9}, and ISEs¹⁶ (<http://www.snl.salk.edu/~geneyeo/stuff/papers/supplementary/ISRE/>). These datasets were originally reported as hexamers, such that pentameric equivalents were created by extracting all pentamers that occurred at least one time within the original datasets. The ISRE enriched pentamers were also compared to ISREs³³, CISs²¹, and motifs enriched upstream of weak PY tracts³⁷. Both of these datasets were composed of various length n-mers and were adjusted to pentameric equivalents to achieve independent sampling by extracting all pentamers that occurred at least once. Lastly, the ISRE pentamers were compared to conserved pentamers enriched in intronic regions of exons excluded in NP cells³⁶. Since these were reported as pentamers no adjustments were necessary.

The significance of overlap between datasets was determined using a 2x2 Chi-test of association. Each pentamer was classified according to which of the four ways it could be distributed: (1) in both sets, (2) in set A but not set B, (3) not in set A but in set B, (4) in neither set. The counts for each distribution were then used to calculate the likelihood that this arrangement could have occurred randomly (according to the Chi-distribution with 1-degree of freedom).

2.4.8. Statistical Analysis

Data are expressed as normalized or fold expression \pm average error where applicable. Student's *t*-test and Anova analyses were performed using Microsoft Excel. *P* < .05 were taken to be significant.

2.4.9. ISRE library and ISS controls construction

A random 15-nt ISRE library was generated through PCR using a 47-nt template (ISStemp) with primers Lib1 and Lib2. The library PCR was conducted for 12 cycles in a 100 μ l reaction containing 20 pmol DNA template, 300 pmol each Lib1 and Lib2, 200 μ M each dNTPs, 1.6 mM MgCl₂, and 10 U *Taq* DNA polymerase (Roche). ISS and negative controls were constructed by replacing the random 15-nt region in the above template with previously characterized ISS sequences and scrambled sequences, respectively. The resulting ISRE library, ISS control, and negative control fragments were digested with Eco RV and Cla I and ligated into the corresponding restriction sites within intron 6 of pCS516. Control ISS sequences correspond to previously characterized binding sites for U2AF65 (TTTTTTTTTCCTTTTTTTTCCTTTT; pCS668)³⁹; hnRNP H (TAAATGTGGGACCTAGA; pCS669)³⁸, PTB(1) (TAGCATCAGCCTGGTGCCTACCTTCGGCCCC; pCS670)³⁹; PTB(2) (TCTTCTCTTCTCTTCTCTTC; pCS667)⁴⁰. In addition, 15 scrambled sequences were examined in place of the 15-nt random region as negative control constructs. The base random 15-nt sequence ACCTCAGGCTCTGAA (pCS517) was subsequently used as the negative control for all FACS experiments.

2.4.10. Quantitative RT-PCR analysis

Total cellular RNA was purified from stably transfected HEK-293 Flp-In cells using GenElute mammalian total RNA purification kit (Sigma) according to the manufacturer's instructions, followed by DNase treatment (Invitrogen). cDNA was synthesized using a gene-specific primer for the pcDNA5/FRT vector (SMN1cDNA) and Superscript III reverse transcriptase (Invitrogen) according to the manufacturer's instructions. qRT-PCR analysis was performed using isoform-specific primers (Tables S2.9 and S2.10) where each reaction contained 1 μ L template cDNA, 10 pmol of each primer and 1X iQ SYBR green supermix (BioRAD) to a final volume of 25 μ L. Reactions were carried out using a iCycler iQ system (BioRAD) for 30 cycles (95°C for 15 s, 72°C for 30 s). The purity of the PCR products was determined by melt curve analysis. Data analysis was completed using the iCycler IQ system software v.3.1.7050 (BioRAD). Isoform-specific relative expression was calculated using the Δ Ct (change in cycling threshold) method⁵². Expression levels were normalized to the levels of *HPRT* (Hypoxanthine-guanine phosphoribosyltransferase). Fold expression data is reported as the mean expression for each sample divided by the mean NMD expression value \pm the average error.

2.4.11. Discovery of sequence motifs enriched in ISRE sequences

Sequence motifs were constructed from the significantly enriched ($P < 0.1$) n-mers using the graph clustering method and software (GCCS)²¹ with the following parameters: minimum cluster size = 4, rounds of clustering = 5, minimum substring

length = 5 (rounds 1–3) and 4 (rounds 4 and 5). GCCS uses the MCL algorithm^{53,54} to find clusters. Parameters were set as follows: MCL inflation = 3 and MCL scheme = 4. The other MCL parameters were set to default values. To validate the enrichment of ISRE motifs, the GCCS analysis was repeated using 5 sets of 125 random 15-mers with the same constant flanking bases. The average number of significantly enriched n-mers observed in the random samples (RS) was only 91 and each yielded an average of 11 clusters.

2.4.12. Hierarchical clustering

A distance matrix for ISRE sequences recovered from SPLICE was produced using the Jukes-Cantor method⁵⁵ in which the distance is defined by the maximum likelihood estimate of the number of nucleotide substitutions between two sequences (Matlab default method) (<http://mathworks.com>). The distance matrix was then used to cluster sequences using the standard average linkage hierarchical clustering implemented in Matlab. 15-nt clusters were defined by using a dissimilarity cutoff of 1.1 in the dendrogram. Sequences within each cluster were then aligned with ClustalX using default parameters⁵⁶.

2.4.13. RNA structural analysis

RNA secondary structure predictions were performed using RNAfold⁵⁷.

Acknowledgments

We thank R. Diamond and D. Perez (Caltech Cell Sorting Facility) for FACS assistance and expert technical advice, A. Krainer for providing the pCISMN Δ 6-wt construct.

References

1. Black, D.L. Mechanisms of alternative pre-messenger RNA splicing. *Annu Rev Biochem* **72**, 291–336 (2003).
2. Blencowe, B.J. Alternative splicing: new insights from global analyses. *Cell* **126**, 37–47 (2006).
3. Wang, E.T. et al. Alternative isoform regulation in human tissue transcriptomes. *Nature* **456**, 470–6 (2008).
4. Castle, J.C. et al. Expression of 24,426 human alternative splicing events and predicted cis regulation in 48 tissues and cell lines. *Nat Genet* **40**, 1416–25 (2008).
5. Pan, Q., Shai, O., Lee, L.J., Frey, B.J. & Blencowe, B.J. Deep surveying of alternative splicing complexity in the human transcriptome by high-throughput sequencing. *Nat Genet* **40**, 1413–5 (2008).
6. Fairbrother, W.G., Yeh, R.F., Sharp, P.A. & Burge, C.B. Predictive identification of exonic splicing enhancers in human genes. *Science* **297**, 1007–13 (2002).
7. Zhang, X.H. & Chasin, L.A. Computational definition of sequence motifs governing constitutive exon splicing. *Genes Dev* **18**, 1241–50 (2004).
8. Coulter, L.R., Landree, M.A. & Cooper, T.A. Identification of a new class of exonic splicing enhancers by in vivo selection. *Mol Cell Biol* **17**, 2143–50 (1997).
9. Wang, Z. et al. Systematic identification and analysis of exonic splicing silencers. *Cell* **119**, 831–45 (2004).
10. Wang, Z. & Burge, C.B. Splicing regulation: from a parts list of regulatory elements to an integrated splicing code. *Rna* **14**, 802–13 (2008).
11. Wang, G.S. & Cooper, T.A. Splicing in disease: disruption of the splicing code and the decoding machinery. *Nat Rev Genet* **8**, 749–61 (2007).

12. Charlet, B.N., Logan, P., Singh, G. & Cooper, T.A. Dynamic antagonism between ETR-3 and PTB regulates cell type-specific alternative splicing. *Mol Cell* **9**, 649–58 (2002).
13. Smith, C.W. & Valcarcel, J. Alternative pre-mRNA splicing: the logic of combinatorial control. *Trends Biochem Sci* **25**, 381–8 (2000).
14. Simard, M.J. & Chabot, B. SRp30c is a repressor of 3' splice site utilization. *Mol Cell Biol* **22**, 4001–10 (2002).
15. Matlin, A.J., Clark, F. & Smith, C.W. Understanding alternative splicing: towards a cellular code. *Nat Rev Mol Cell Biol* **6**, 386–98 (2005).
16. Yeo, G., Hoon, S., Venkatesh, B. & Burge, C.B. Variation in sequence and organization of splicing regulatory elements in vertebrate genes. *Proc Natl Acad Sci U S A* **101**, 15700–5 (2004).
17. Venables, J.P. Downstream intronic splicing enhancers. *FEBS Lett* **581**, 4127–31 (2007).
18. Green, R.E. et al. Widespread predicted nonsense-mediated mRNA decay of alternatively-spliced transcripts of human normal and disease genes. *Bioinformatics* **19 Suppl 1**, i118–21 (2003).
19. Ladd, A.N. & Cooper, T.A. Finding signals that regulate alternative splicing in the post-genomic era. *Genome Biol* **3**, (2002).
20. McCullough, A.J. & Berget, S.M. G triplets located throughout a class of small vertebrate introns enforce intron borders and regulate splice site selection. *Mol Cell Biol* **17**, 4562–71 (1997).
21. Voelker, R.B. & Berglund, J.A. A comprehensive computational characterization of conserved mammalian intronic sequences reveals conserved motifs associated with constitutive and alternative splicing. *Genome Res* **17**, 1023–33 (2007).
22. Burd, C.G. & Dreyfuss, G. RNA binding specificity of hnRNP A1: significance of hnRNP A1 high-affinity binding sites in pre-mRNA splicing. *Embo J* **13**, 1197–204 (1994).
23. Markovtsov, V. et al. Cooperative assembly of an hnRNP complex induced by a tissue-specific homolog of polypyrimidine tract binding protein. *Mol Cell Biol* **20**, 7463–79 (2000).
24. Chan, R.C. & Black, D.L. Conserved intron elements repress splicing of a neuron-specific c-src exon in vitro. *Mol Cell Biol* **15**, 6377–85 (1995).

25. Hui, J., Stangl, K., Lane, W.S. & Bindereif, A. HnRNP L stimulates splicing of the eNOS gene by binding to variable-length CA repeats. *Nat Struct Biol* **10**, 33–7 (2003).
26. Tacke, R. & Manley, J.L. The human splicing factors ASF/SF2 and SC35 possess distinct, functionally significant RNA binding specificities. *Embo J* **14**, 3540–51 (1995).
27. Liu, H.X., Zhang, M. & Krainer, A.R. Identification of functional exonic splicing enhancer motifs recognized by individual SR proteins. *Genes Dev* **12**, 1998–2012 (1998).
28. Cavaloc, Y., Bourgeois, C.F., Kister, L. & Stevenin, J. The splicing factors 9G8 and SRp20 transactivate splicing through different and specific enhancers. *Rna* **5**, 468–83 (1999).
29. Tacke, R., Tohyama, M., Ogawa, S. & Manley, J.L. Human Tra2 proteins are sequence-specific activators of pre-mRNA splicing. *Cell* **93**, 139–48 (1998).
30. Lim, L.P. & Burge, C.B. A computational analysis of sequence features involved in recognition of short introns. *Proc Natl Acad Sci U S A* **98**, 11193–8 (2001).
31. Nasim, M.T., Chernova, T.K., Chowdhury, H.M., Yue, B.G. & Eperon, I.C. HnRNP G and Tra2beta: opposite effects on splicing matched by antagonism in RNA binding. *Hum Mol Genet* **12**, 1337–48 (2003).
32. Pagani, F. et al. A new type of mutation causes a splicing defect in ATM. *Nat Genet* **30**, 426–9 (2002).
33. Yeo, G.W., Nostrand, E.L. & Liang, T.Y. Discovery and analysis of evolutionarily conserved intronic splicing regulatory elements. *PLoS Genet* **3**, e85 (2007).
34. Marquis, J. et al. CUG-BP1/CELF1 requires UGU-rich sequences for high-affinity binding. *Biochem J* **400**, 291–301 (2006).
35. Hovhannisyan, R.H. & Carstens, R.P. Heterogeneous ribonucleoprotein m is a splicing regulatory protein that can enhance or silence splicing of alternatively spliced exons. *J Biol Chem* **282**, 36265–74 (2007).
36. Yeo, G.W. et al. Alternative splicing events identified in human embryonic stem cells and neural progenitors. *PLoS Comput Biol* **3**, 1951–67 (2007).
37. Murray, J.I., Voelker, R.B., Henscheid, K.L., Warf, M.B. & Berglund, J.A. Identification of motifs that function in the splicing of non-canonical introns. *Genome Biol* **9**, R97 (2008).

38. Chen, C.D., Kobayashi, R. & Helfman, D.M. Binding of hnRNP H to an exonic splicing silencer is involved in the regulation of alternative splicing of the rat beta-tropomyosin gene. *Genes Dev* **13**, 593–606 (1999).
39. Singh, R., Valcarcel, J. & Green, M.R. Distinct binding specificities and functions of higher eukaryotic polypyrimidine tract-binding proteins. *Science* **268**, 1173–6 (1995).
40. Perez, I., Lin, C.H., McAfee, J.G. & Patton, J.G. Mutation of PTB binding sites causes misregulation of alternative 3' splice site selection in vivo. *Rna* **3**, 764–78 (1997).
41. Buratti, E., Stuani, C., De Prato, G. & Baralle, F.E. SR protein-mediated inhibition of CFTR exon 9 inclusion: molecular characterization of the intronic splicing silencer. *Nucleic Acids Res* **35**, 4359–68 (2007).
42. Hanamura, A., Caceres, J.F., Mayeda, A., Franza, B.R., Jr. & Krainer, A.R. Regulated tissue-specific expression of antagonistic pre-mRNA splicing factors. *Rna* **4**, 430–44 (1998).
43. Wollerton, M.C. et al. Differential alternative splicing activity of isoforms of polypyrimidine tract binding protein (PTB). *Rna* **7**, 819–32 (2001).
44. Pozzoli, U. & Sironi, M. Silencers regulate both constitutive and alternative splicing events in mammals. *Cell Mol Life Sci* **62**, 1579–604 (2005).
45. Liu, H.X., Cartegni, L., Zhang, M.Q. & Krainer, A.R. A mechanism for exon skipping caused by nonsense or missense mutations in BRCA1 and other genes. *Nat Genet* **27**, 55–8 (2001).
46. Paul, S. et al. Interaction of muscleblind, CUG-BP1 and hnRNP H proteins in DM1-associated aberrant IR splicing. *Embo J* **25**, 4271–83 (2006).
47. Mauger, D.M., Lin, C. & Garcia-Blanco, M.A. hnRNP H and hnRNP F complex with Fox2 to silence fibroblast growth factor receptor 2 Exon IIIc. *Mol Cell Biol* **28**, 5403–19 (2008).
48. Venables, J.P. et al. Multiple and Specific mRNA Processing Targets for the Major Human hnRNP Proteins. *Mol Cell Biol* (2008).
49. Sambrook, J. & Russell, D.W. *Molecular Cloning: a laboratory manual*, (Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, 2001).
50. Cartegni, L. & Krainer, A.R. Disruption of an SF2/ASF-dependent exonic splicing enhancer in SMN2 causes spinal muscular atrophy in the absence of SMN1. *Nat Genet* **30**, 377–84 (2002).

51. Stade, K., Ford, C.S., Guthrie, C. & Weis, K. Exportin 1 (Crm1p) is an essential nuclear export factor. *Cell* **90**, 1041–50 (1997).
52. Livak, K.J. & Schmittgen, T.D. Analysis of relative gene expression data using real-time quantitative PCR and the 2(-Delta Delta C(T)) Method. *Methods* **25**, 402–8 (2001).
53. Dongen, S. University of Utrecht (2000).
54. Enright, A.J., Van Dongen, S. & Ouzounis, C.A. An efficient algorithm for large-scale detection of protein families. *Nucleic Acids Res* **30**, 1575–84 (2002).
55. Jukes, T.a.C., CR. Evolution of protein molecules. in *Mammalian protein metabolism* 21–123 (Academic Press, New York, 1969).
56. Thompson, J.D., Gibson, T.J., Plewniak, F., Jeanmougin, F. & Higgins, D.G. The CLUSTAL_X windows interface: flexible strategies for multiple sequence alignment aided by quality analysis tools. *Nucleic Acids Res* **25**, 4876–82 (1997).
57. Hofacker, I.L. et al. Fast Folding and Comparison of RNA Secondary Structures. *Monatshefte f. Chemie* **125**, 167–188 (1994).
58. Cartegni, L., Hastings, M.L., Calarco, J.A., de Stanchina, E. & Krainer, A.R. Determinants of exon 7 splicing in the spinal muscular atrophy genes, SMN1 and SMN2. *Am J Hum Genet* **78**, 63–77 (2006).
59. Hentze, M.W. & Kulozik, A.E. A perfect message: RNA surveillance and nonsense-mediated decay. *Cell* **96**, 307–10 (1999).
60. Singh, R. & Valcarcel, J. Building specificity with nonspecific RNA-binding proteins. *Nat Struct Mol Biol* **12**, 645–53 (2005).
61. Buratti, E. & Baralle, F.E. Influence of RNA secondary structure on the pre-mRNA splicing process. *Mol Cell Biol* **24**, 10505–14 (2004).
62. Hiller, M., Zhang, Z., Backofen, R. & Stamm, S. Pre-mRNA Secondary Structures Influence Exon Recognition. *PLoS Genet* **3**, e204 (2007).
63. Chou, M.Y., Rooke, N., Turck, C.W. & Black, D.L. hnRNP H is a component of a splicing enhancer complex that activates a c-src alternative exon in neuronal cells. *Mol Cell Biol* **19**, 69–77 (1999).
64. Hutchison, S., LeBel, C., Blanchette, M. & Chabot, B. Distinct sets of adjacent heterogeneous nuclear ribonucleoprotein (hnRNP) A1/A2 binding sites control 5' splice site selection in the hnRNP A1 mRNA precursor. *J Biol Chem* **277**, 29745–52 (2002).

65. Miriami, E., Margalit, H. & Sperling, R. Conserved sequence elements associated with exon skipping. *Nucleic Acids Res* **31**, 1974–83 (2003).
66. Chou, M.Y., Underwood, J.G., Nikolic, J., Luu, M.H. & Black, D.L. Multisite RNA binding and release of polypyrimidine tract binding protein during the regulation of c-src neural-specific splicing. *Mol Cell* **5**, 949–57 (2000).
67. Itoh, H., Washio, T. & Tomita, M. Computational comparative analyses of alternative splicing regulation using full-length cDNA of various eukaryotes. *Rna* **10**, 1005–18 (2004).
68. Paronetto, M.P., Achsel, T., Massiello, A., Chalfant, C.E. & Sette, C. The RNA-binding protein Sam68 modulates the alternative splicing of Bcl-x. *J Cell Biol* **176**, 929–39 (2007).
69. Lim, L.P. & Burge, C.B. A computational analysis of sequence features involved in recognition of short introns. *Proceedings of the National Academy of Sciences of the United States of America* **98**, 11193–11198 (2001).

Supplementary Information

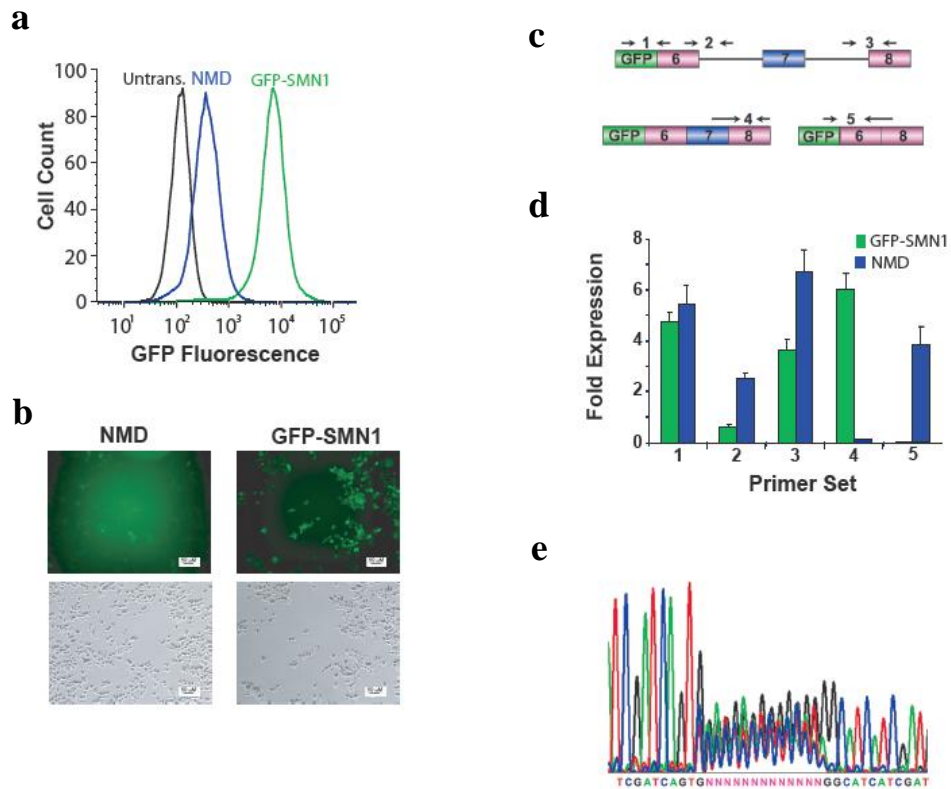


Figure S2.1. Fluorescence expression/analysis of SPLICE control constructs and library sequence bias. **(a)** Flow cytometry histograms of the stable cell lines expressing the control constructs. An untransfected HEK-293 FLP-In cell population (Untrans.) was also analyzed for reference. **(b)** Microscope images of stable cell lines expressing the negative (NMD) and positive (GFP-SMN1) control constructs. Upper panels: GFP fluorescence, lower panels: phase contrast images. **(c)** Schematic representing the relative locations of primer set binding for transcript isoform analysis by qRT-PCR. **(d)** qRT-PCR analysis of the NMD and GFP-SMN1 control cell lines supports decay of the PTC harboring isoform. The observed high level of exon 7 inclusion for the GFP-SMN1 control are in line with previous observations for the splicing of the SMN1 mini-gene⁵⁸. Our transcript isoform analysis also reveals that levels of exon 7 exclusion are elevated in the NMD control compared to the GFP-SMN1 control, suggesting that the PTC may have a

secondary effect of increasing the levels of the exon excluded transcript. Such observations have been previously observed and may be the result of nonsense-associated altered splicing⁵⁹. Expression levels were normalized to the levels of *HPRT* (Hypoxanthine-guanine phosphoribosyltransferase). Data presented is the mean expression of duplicate PCR samples \pm the average error. (e) DNA sequencing analysis of purified genomic DNA from HEK-293 cell lines harboring the library constructs. The transfected library exhibits a slight sequence bias at positions 1 and 15, but all other positions are free of bias. A comparison of sequencing results from two independent transfections supports that the sequence bias at these positions is minimal, indicating that the ISRE library represents an essentially random pool (data not shown).

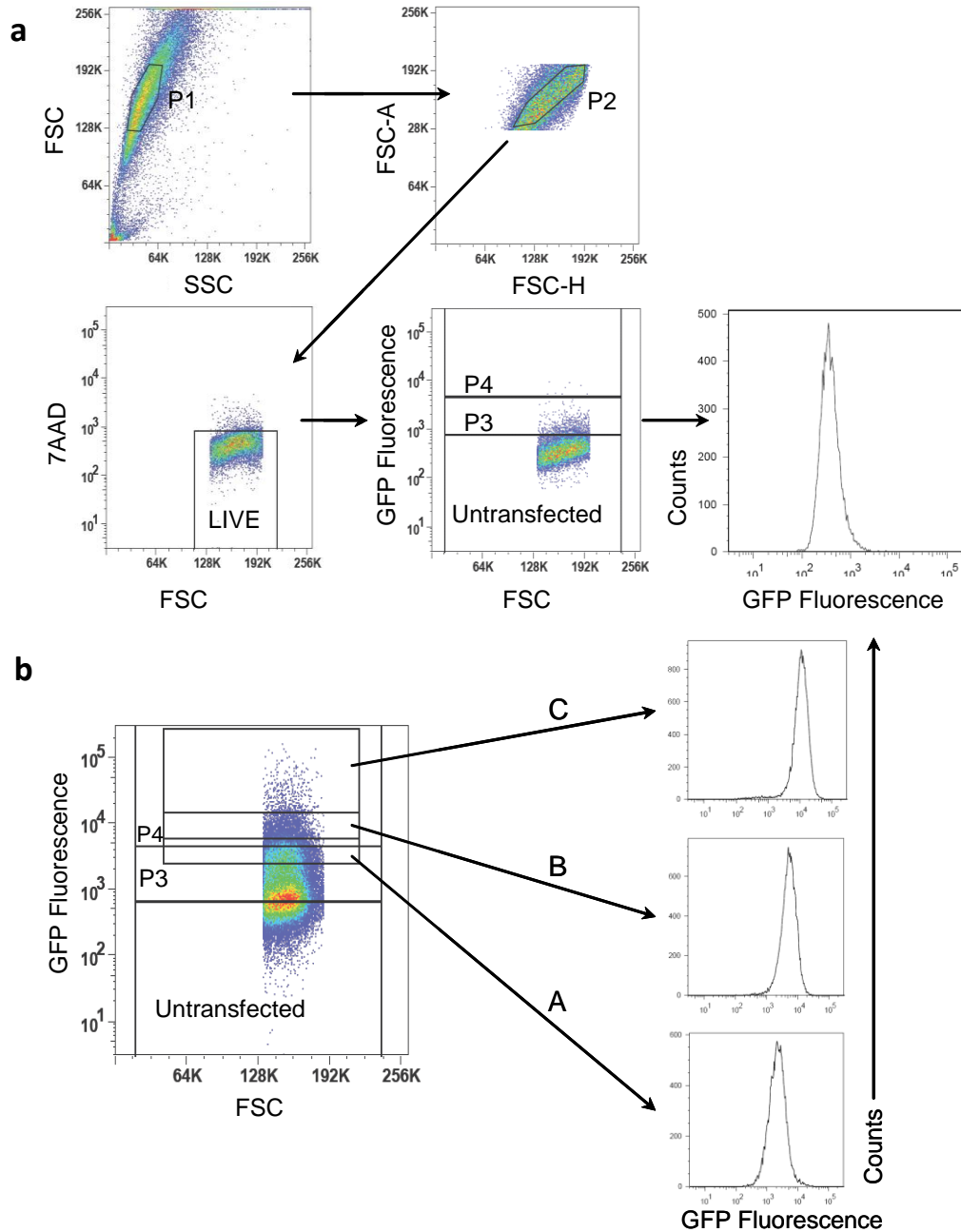


Figure S2.2. FACS analysis and ISRE library sorting scheme. **(a)** FACS analysis and gating procedure for all HEK-293 FLP-In cells. As an example, flow cytometry data from the stable NMD control is presented. Dot plots show initial gating of stable cells (P1), followed by P2 gating for cell uniformity (i.e., to remove cell aggregates) and finally the selection of live cells using 7-Amino-Actinomycin D (7AAD) staining. The P3 gate

reflects the GFP positive cells and the P4 gate is drawn to indicate the upper GFP fluorescence limit of the NMD control population. P4 was used as the gate for the selection of ISS positive cells. The histogram reports the intensity of GFP fluorescence in the NMD control population. **(b)** FACS analysis of ISS positive stable cells after one round of sorting. Cells from gates A, B, and C were sorted and the resulting histograms indicate the intensity of GFP fluorescence after 1 week in culture.

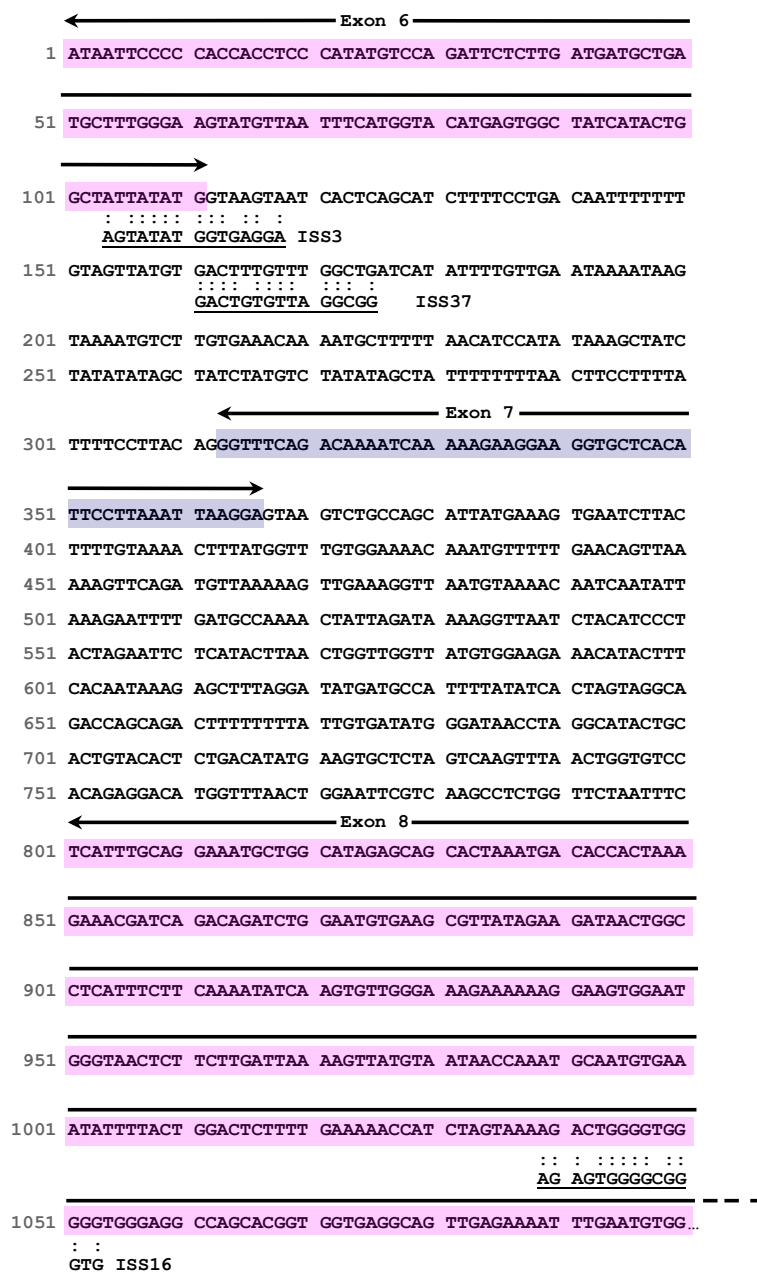


Figure S2.3. SPLICE-generated ISREs contain elements similar to sequences within SMN1. The sequence of the wild-type SMN1 mini-gene used in this study is shown. SPLICE-generated sequences similar to portions of the SMN1 mini-gene are underlined and labeled accordingly. Sequences ISS3, ISS16, and ISS37 are similar to sections of intron 6 and exon 8 and regions spanning the exonic and intronic portion of the 5' splice site.

exon 6. Splicing repression may be a result of cooperative repressor binding to multiple silencer elements creating a 'zone of silencing' (at one site or at overlapping sites) between splice sites or nucleation that causes the looping out of RNA between repressor elements^{15,60}. The sequence composition of ISS3, 16 and 37 supports a model of cooperative assembly of repressor elements to the SMN1 transcript in regulating splicing repression. The silencer activity of ISS3 and ISS16 has been confirmed (see Figure 2.5b).

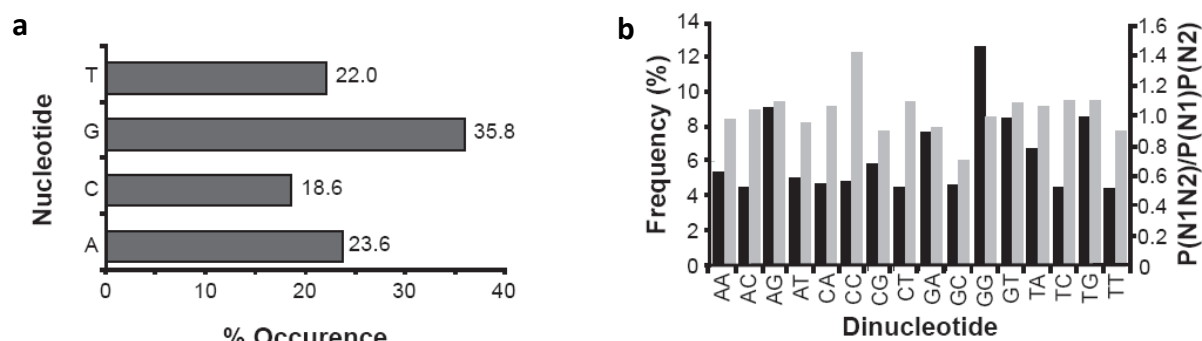


Figure S2.4. Overall compositional features of recovered ISRE sequences. **(a)** Overall nucleotide composition of the recovered ISREs. The occurrence value of each nucleotide prior to enrichment for ISRE activity was 25%. SPLICE-generated sequences have a higher level of G (35.8%) and reduced levels of T (22%), C (18.6%), and A (23.6%). **(b)** Dinucleotide frequency (black) and odds ratios (gray) in recovered ISRE. The occurrence of dinucleotides within all 125 ISRE sequences was calculated (black). The odds ratio for each dinucleotide was determined by dividing the probability of a dinucleotide occurrence within the 125 SPLICE selected sequences by the probability of an individual nucleotide occurrence within selected sequences ($P(N1N2)/P(N1)P(N2)$) (gray). The dinucleotide CC is overrepresented in the ISRE dataset, while others, such as AC, AG, CA, GT, TA, TC, and TG, are only slightly enriched.

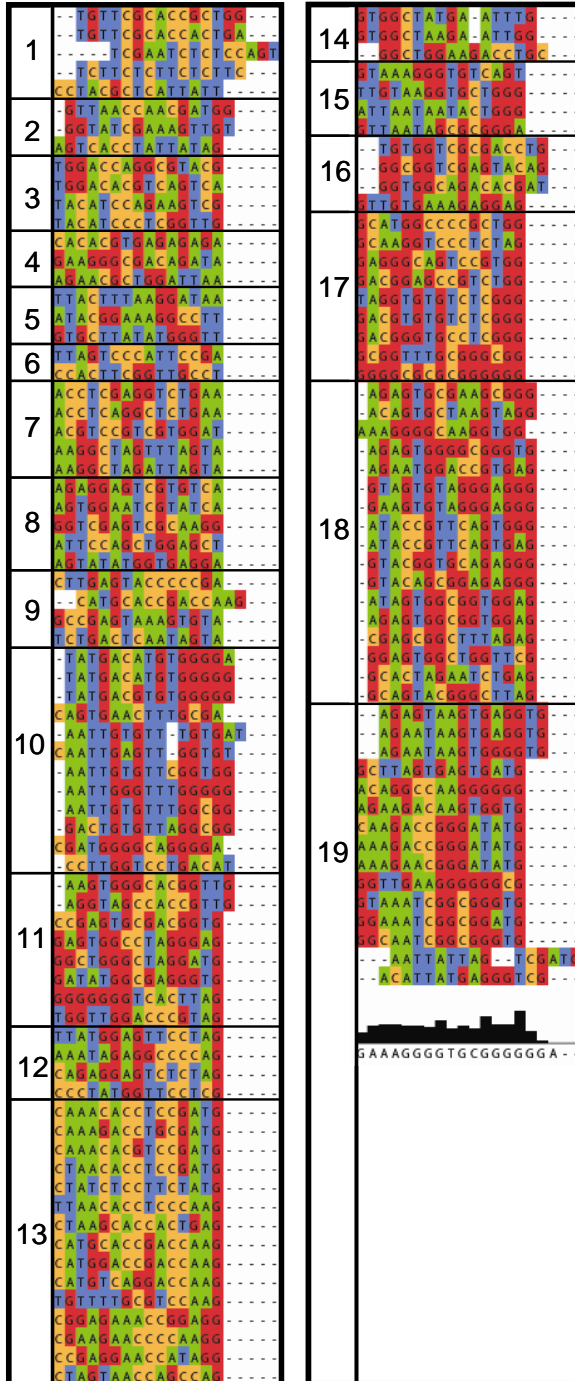


Figure S2.5. ISRE hierarchical clusters and sequence alignment. ISRE hierarchical clusters and sequence alignment of individual clusters. Sequences were aligned using ClustalX⁵⁶.

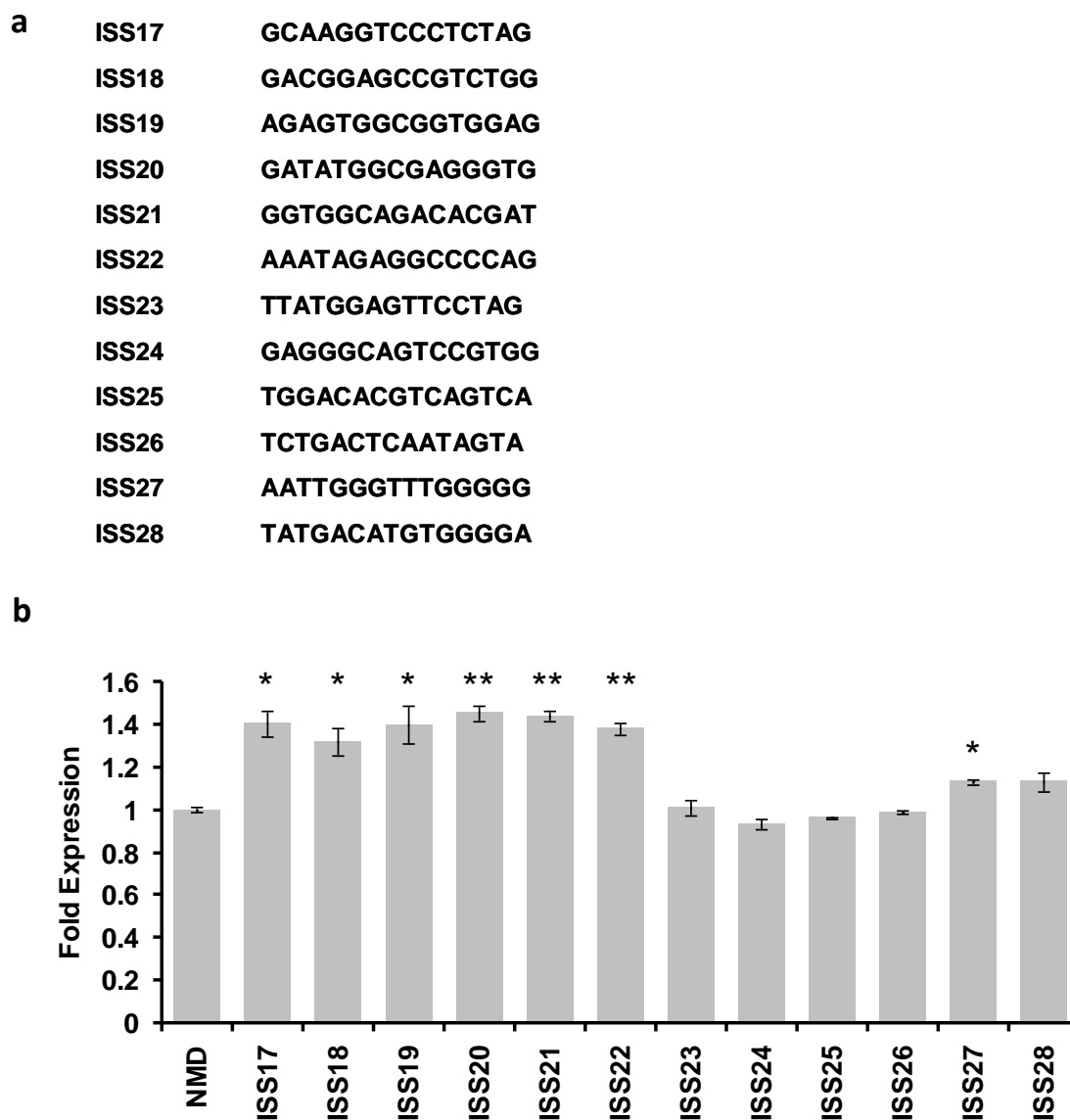


Figure S2.6. The activity of additional recovered ISRE sequences is validated by stable cell line assays. **(a)** Additional recovered ISRE sequences examined for regulatory activity. **(b)** Flow cytometry analysis of HEK-293 FLP-In stable cell lines generated for each recovered ISRE sequence and control construct. Mean GFP levels from two independent experiments were determined and normalized to the NMD control. The fold

expression of each sample relative to NMD and average error are reported. Resulting P -values in comparison to the NMD control: * $P < 0.03$ and ** $P < 0.01$.

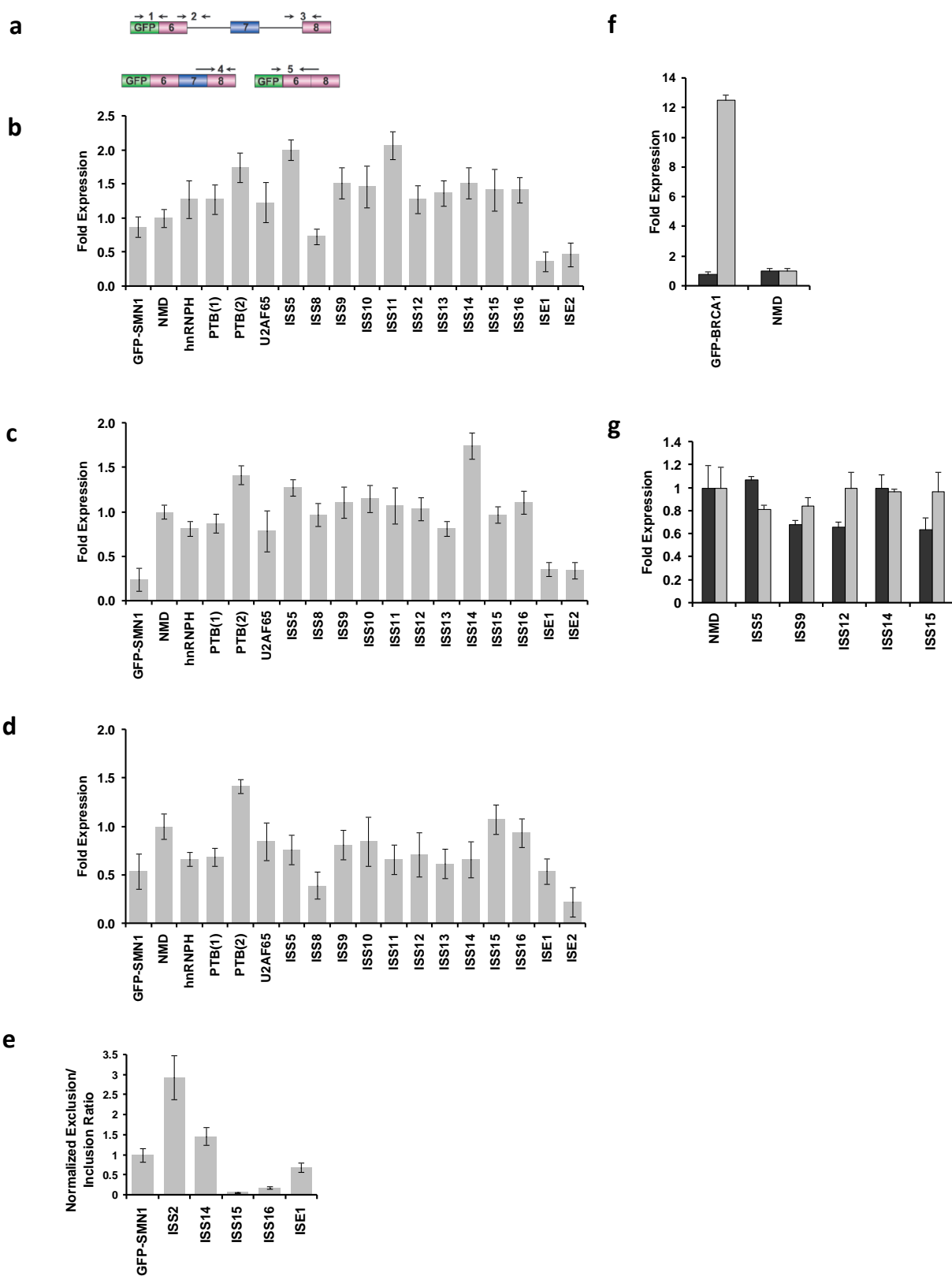


Figure S2.7. Additional qRT-PCR isoform analysis of recovered ISREs and control constructs. **(a)** Schematic representing the relative locations of primer set binding on the reporter system for transcript isoform analysis. **(b)** qRT-PCR analysis with primer set 1. Results demonstrate that overall transcript levels for the GFP-SMN1, ISS controls, ISSs and ISEs did not significantly differ from the NMD control ($P = 0.20$). For all subsequent analyses, expression levels of duplicate PCR samples were normalized to the levels of *HPRT*. Fold expression data is reported as the mean expression for each sample divided by the mean NMD expression value \pm the average error. **(c)** qRT-PCR analysis with primer set 2. The levels of intron 6 retained in transcripts containing the selected and control ISS sequences are similar to the NMD control ($P = 0.48$). In contrast, intron 6 retention in ISE transcripts are similar to the GFP-SMN1 control ($P = 0.74$) and different from the NMD control ($P < 0.05$), suggesting that intron 6 in the GFP-SMN1 control and ISEs are processed similarly by the general splicing machinery. The retention level of intron 6 for the GFP-SMN1 control is statistically different from the NMD control ($P < 0.05$). **(d)** qRT-PCR analysis with primer set 3. The levels of intron 7 retention for the recovered and control ISS sequences and the GFP-SMN1 are similar to the NMD control ($P = 0.23$). The intron 7 retention levels in ISE transcripts are significantly different from the NMD control ($P < .05$). **(e)** qRT-PCR analysis with primer sets 4 and 5 on selected ISREs inserted in the non-NMD-based GFP-SMN1 control construct. Stable cell lines containing ISS2, ISS14 and ISE1 maintained selected ISRE function; however, ISS15 and ISS16 displayed significant enhancer activity ($P < 0.05$). The results suggest that ISS15 and ISS16 may exhibit enhanced fluorescence levels in the context of the NMD reporter due to the evasion of the NMD process. Data is reported as the expression ratio

of the mean expression of the exon excluded isoform to the exon included isoform normalized to the ratio for the GFP-SMN1 control \pm the average error. **(f)** qRT-PCR analysis of the NMD and GFP-BRACA1 constructs with primer sets specific for exon 18 excluded (black bars) and included (gray bars) products, supports decay of the PTC harboring isoform. **(g)** qRT-PCR analysis of selected sequences inserted into the BRCA1-NMD construct with primer sets specific for exon 18 excluded (black bars) and included

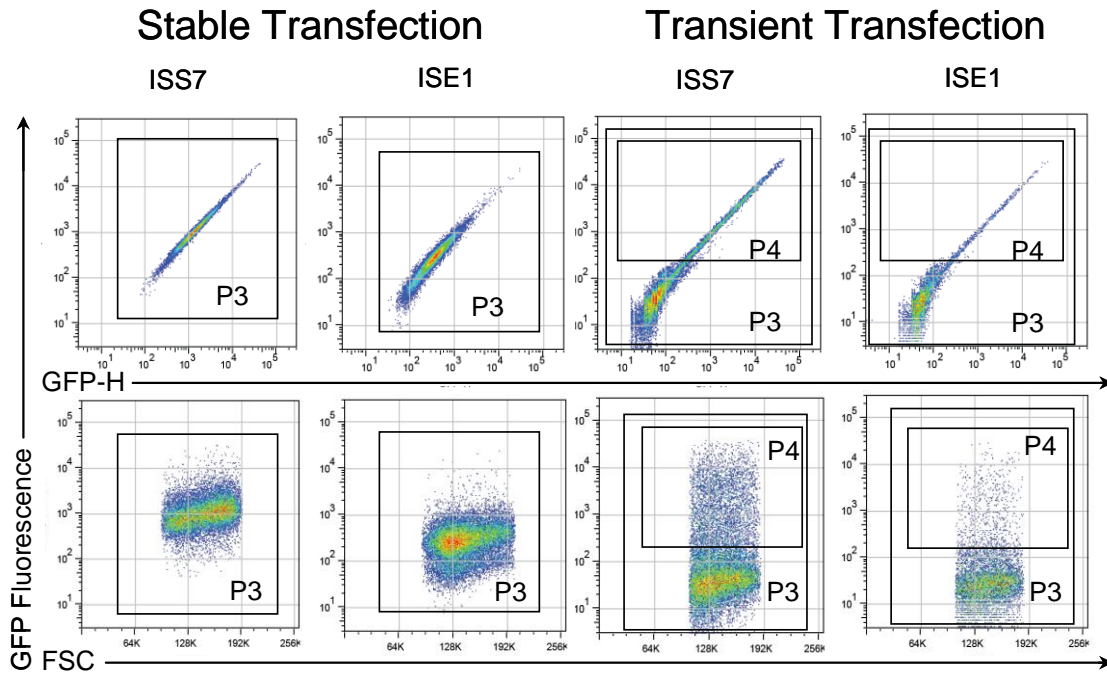


Figure S2.8. Assessment of splicing regulatory activity through stable and transient transfection assays. Sixteen recovered ISS sequences (ISS1-ISS16) and 1 recovered ISE sequence (ISE1) were examined for regulatory activity in both transient and stable transfection assays. Examples of assay results for two recovered sequences (ISS7, ISE1) are shown. For the stable cell line assays, mean GFP fluorescence levels were determined using gate P3. For the transient transfection assays, the P3 gate represents the untransfected cell population and the P4 gate represents the GFP-positive cells. The results of an ANOVA analysis applied to data from the transient and stable assays indicate that the two methods are not statistically similar ($P = 0.27$).

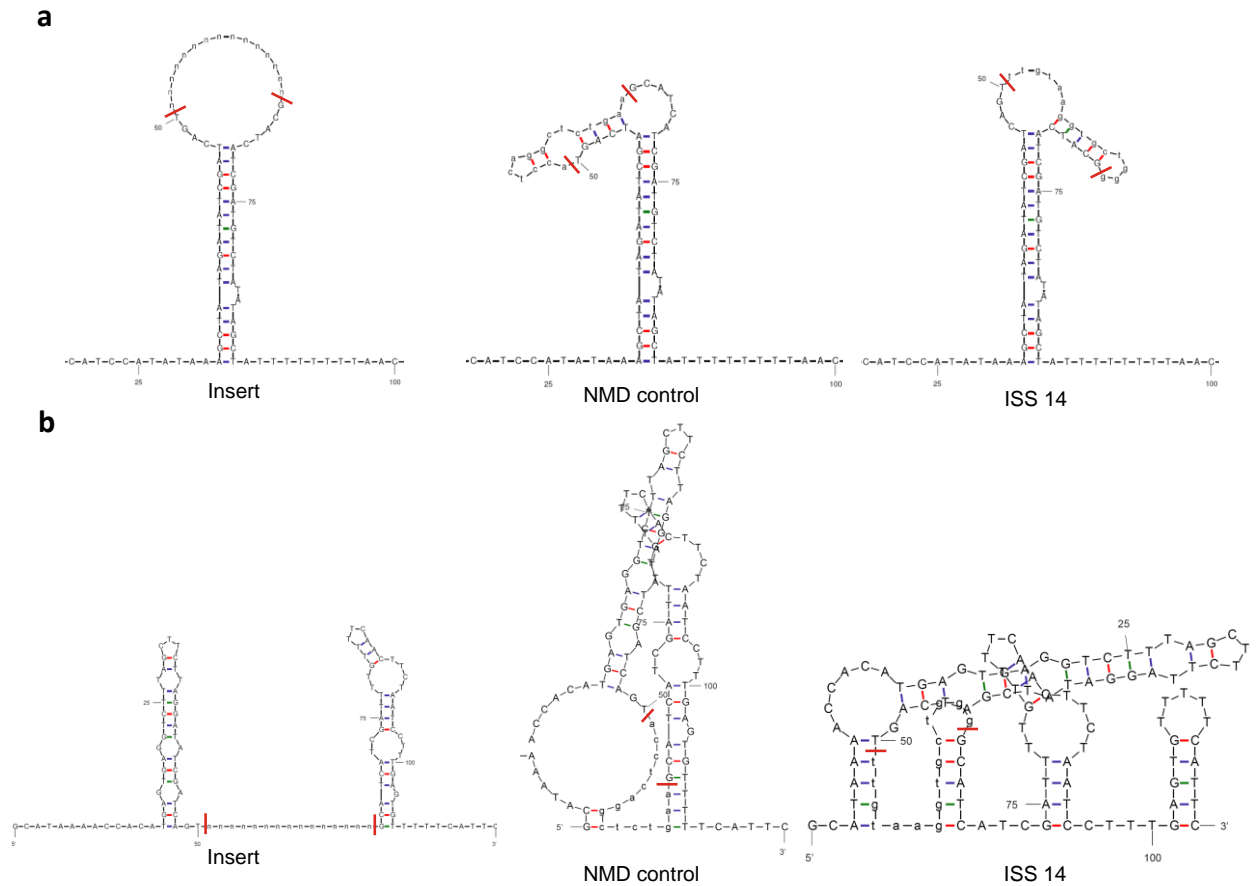


Figure S2.9. Predicted secondary structure for ISREs in the SMN1 and BRCA1 mini-genes. RNA secondary structures have been proposed to play a role in alternative splicing⁶¹. It is possible that the insertion of sequences within intron 6 of the SMN1 mini-gene may have resulted in a secondary structure that disrupts the binding of trans-acting factors. We examined structures of a region -50- to + 50-nt from the 15-mers within the SMN1 mini-gene using RNAfold⁵⁷. **(a)** Predicted secondary structure for the intronic regions +/- 50-nt of any ISRE sequence, the random 15-nt (NMD control) and ISS14 in SMN1 mini-gene. Inserted sequences are located in position 51 to 65 in lower case letters and denoted by a red dash. The inserted sequences are predicted to be generally located

within a looped region. The local secondary structure around 15-mers within the SMN1 mini-gene is dominated by a hairpin structure. **(b)** Predicted secondary structure for the intronic regions +/- 50-nt of any ISS sequence, the random 15-mer (NMD control) and ISS14 in the BRCA1 mini-gene. Inserted sequences are located in position 51 to 65 in lower case letters and denoted by a red dash. The inserted sequences are predicted to be generally located within single stranded regions. The major structure within the BRCA1 mini-gene is a double hairpin. Splicing motifs are preferentially found in single-stranded contexts⁶². Taken together, the results suggest that the lack of selected ISRE function within the BRCA1 mini-gene may be due to differences in local secondary structure around the ISREs.

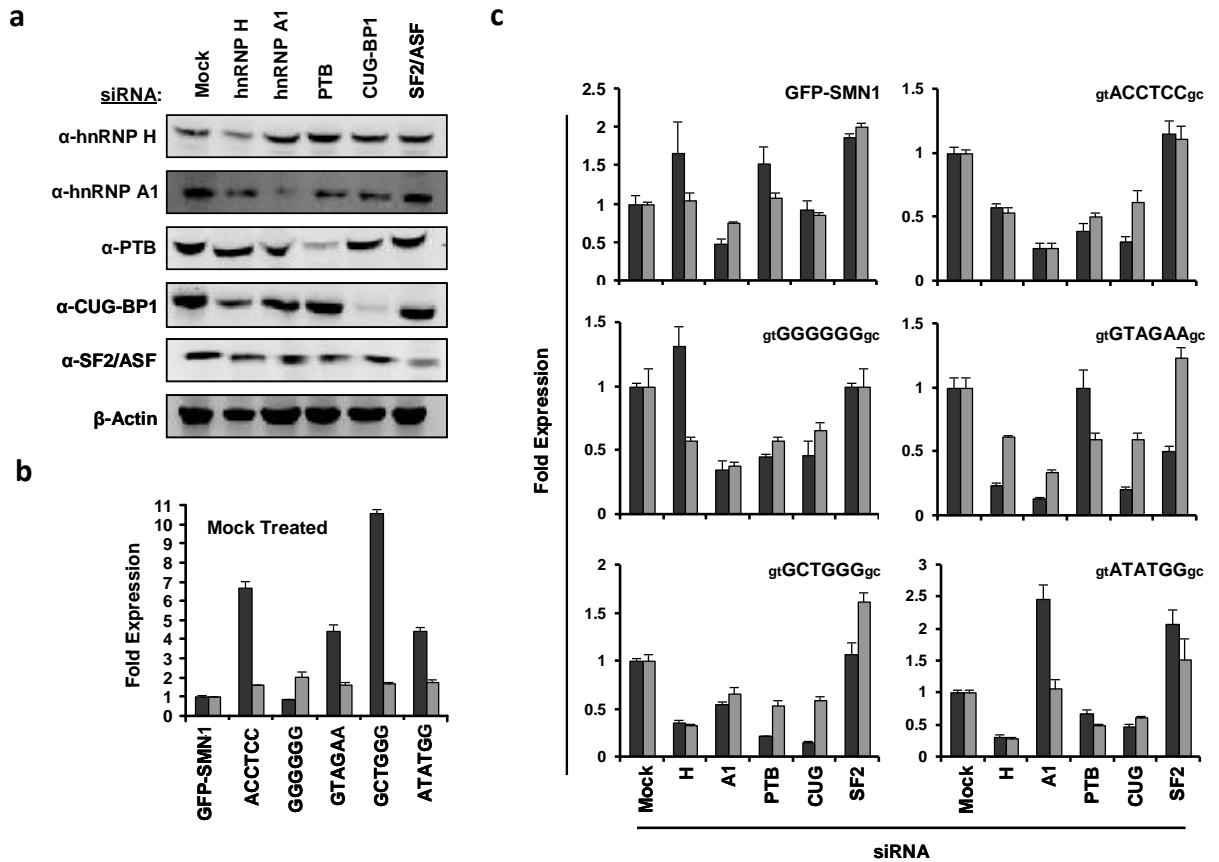


Figure S2.10. The effects of *in vivo* depletion of splicing factors on ISRE regulated splicing patterns. **(a)** Western blot analysis of total cell lysates prepared from the GFP-SMN1 control cell line treated with siRNAs targeted to trans-acting splicing factors and a mock siRNA negative control. β -Actin was used as a loading control for all blots. The results demonstrate that individual siRNAs have minimal to no off-target affects. **(b)** qRT-PCR analysis of the mock treated ISRE hexamer and GFP-SMN1 control cell lines with primer sets specific for exon 7 excluded (black bars) and included (gray bars) products. Expression levels of duplicate PCR samples were normalized to the levels of *HPRT*. Fold expression data is reported as the mean expression for each sample divided by the mean GFP-SMN1 control expression value \pm the average error. **(c)** qRT-PCR analysis of the siRNA treated ISRE hexamer and GFP-SMN1 control cell lines with primer sets specific for exon 7 excluded (black bars) and included (gray bars) products.

Fold expression data is reported as the mean expression for each sample divided by the mean mock siRNA treated cell line control expression value \pm the average error.

Table S2.1. Identified ISRE regulatory sequences

ISS sequences	Name	Tested stably	Tested transiently	Heterologous transcript
GACGTGTGTCTCGGG	ISE1	Y	Y	Y
ATAGTGGCGGTGGAG	ISE2	Y	N	N
TACATCCCTCGGTTG	ISS1	Y	Y	Y
AGAATAAGTGGGGTG	ISS2	Y	Y	Y
AGTATATGGTGAGGA	ISS3	Y	Y	Y
TGTTTTGCGTCCAAG	ISS4	Y	Y	Y
AGAATAAGTGAGGTG	ISS5	Y	Y	Y
CCGAGTGCGACGGTG	ISS6	Y	Y	Y
ACAGGCCAAGGGGGG	ISS7	Y	Y	Y
CAAACACCTCCGATG	ISS8	Y	Y	Y
GGTCGAGTCGCAAGG	ISS9	Y	Y	Y
TAGGTGTGTCTCGGG	ISS10	Y	Y	Y
ACAGTGCTAAGTAGG	ISS11	Y	Y	Y
AAAGACCGGGATATG	ISS12	Y	Y	Y
AGTCACCTATTATAG	ISS13	Y	Y	Y
TTGTAAGGTGCTGGG	ISS14	Y	Y	Y
GGGGCGCGCGGGGGG	ISS15	Y	Y	Y
AGAGTGGGGCGGGTG	ISS16	Y	Y	Y
GCAAGGTCCCTCTAG	ISS17	Y	N	N
GACGGAGCCGTCTGG	ISS18	Y	N	N
AGAGTGGCGGTGGAG	ISS19	Y	N	N
GATATGGCGAGGGTG	ISS20	Y	N	N
GGTGGCAGACACGAT	ISS21	Y	N	N
CCGAGGAACCATAGG	ISS23	Y	N	N
AAATAGAGGCCCCAG	ISS24	Y	N	N
TTATGGAGTTCCTAG	ISS25	Y	N	N
GAGGGCAGTCCGTGG	ISS26	Y	N	N
TGGACACGTCAGTCA	ISS27	Y	N	N
TCTGACTCAATAGTA	ISS28	Y	N	N
AATTGGGTTTGGGGG	ISS29	Y	N	N
TATGACATGTGGGGA	ISS30	Y	N	N
CCCTATGGTTCCTCG	ISS31			
GACGGGTGCCTCGGG	ISS32			
GGCTGGAAGACCTGC	ISS33			
GGAGTGGCTGGTTCG	ISS34			
GGCTGGGCTAGGATG	ISS35			
ACCTCAGGCTCTGAA	ISS36			
GACTGTGTTAGGCGG	ISS37			
AAAGAACGGGATATG				
TCGAATCTCTCCAGT				
CCTACGCTCATTATT				
TCTTCTTCTCTTCTC				
TGTTTCGCACCGCTGG				
TGTTTCGCACCACTGA				
GTTAACCAACGATGG				
GGTATCGAAAGTTGT				

TACATCCAGAAAGTCG
TGGACCAGGCGTACG
CACACGTGAGAGAGA
GAAGGGCGACAGATA
AGAACGCTGGATTAA
TTACTTTAAGGATAA
ATACGGAAAGGCCTT
GTGCTTATATGGGTT
TTAGTCCCATTCCGA
CCACTTCGGTTGCCT
ACGTCCGTCGTGGAT
ACCTCGAGGTCTGAA
AAGGCTAGTTTAGTA
AAGGCTAGATTAGTA
AGAGGAGTCGTGTCA
AGTGGAATCGTATCA
ATTCCAGCTGGAGCT
GCCGAGTAAAGTGTA
CTTGAGTACCCCCGA
CATGCACCGACCAAG
AATTGTGTTTGTGAT
GACTGTGTTAGGCGG
AATTGTGTTTGGCGG
TATGACGTGTGGGGG
TATGACATGTGGGGG
CAATTGAGTTGGTGT
CGATGGGGCAGGGGA
CAGTGAACCTTGCGA
CCTTGCTCCTGACAT
GAGTGGCCTAGGGAG
GGCTGGGCTAGGATG
AAGTGGGCACGGTTG
AGGTAGCCACCGTTG
GGGGGGGTCACCTAG
TGGTTGGACCCGTAG
CCCTATGGTTCCTCG
CTAGTAACCAGCCAG
CTAAGCACCAGTGA
CATGTCAGGACCAAG
CATGGACCGACCAAG
TATGCCTCCCCGATA
CGAAGAACCCCAAGG
CGGAGAAACCGGAGG
CTATCTCCTTCTATG
TTAACACCTCCCAAG
CAAAGACCTGCGATG
CAAACACGTCCGATG
CTAACACCTCCGATG
GTGGCTATGAATTTG
GTGGCTAAGAATTGG
GGCTGGAAGACCTGC
GTAAAGGGTGTCAGT
ATTAATAATACTGGG
GTTAATAGCGCGGGA
TGTGGTCGCGACCTG
GGCGGTTCGAGTACAG

GTTGTGAAAGAGGAG
 GCGGTTTGCGGGCGG
 GACGGGTGCCTCGGG
 GCATGGCCCCGCTGG
 GCACTAGAATCTGAG
 GCAGTACGGGCTTAG
 CGAGCGGCTTTAGAG
 AGAATGGACCGTGAG
 GGAGTGGCTGGTTCG
 GTACAGCGGAGAGGG
 GTACGGTGCAGAGGG
 GTAGTGTAGGGAGGG
 GAAGTGTAGGGAGGG
 ATACCGTTCAGTGGG
 ATACCGTTCAGTGAG
 AAAGGGGCAAGGTGG
 AGAGTGCGAAGCGGG
 GTAAATCGGCGGGTG
 GGAAATCGGCGGATG
 GGCAATCGGCGGGTG
 AAAGAACGGGATATG
 CAAGACCGGGATATG
 AATTATTAGTCGATG
 GCTTAGTGAGTGATG
 AGAAGACAAGTGGTG
 GGTGAAGGGGGGCG
 ACATTATGAGGGTCG
 AGAGTAAGTGAGGTG

The splicing activity of the first 30 sequences was assessed by stable transfection assays.

Additional sequences validated through transiently transfection assays are indicated (Y=Tested and N=Not tested) in addition to those tested in a second transcript.

Table S2.2. Significantly enriched ISRE n-mers

Field	Description										
n-mer	The n-mer										
Length	N-mer length										
Count(ISS)	Counts observed in ISS sample										
Count(RS)	Counts observed in RS sample										
N(ISS)	Total counts performed in ISS sample										
N(RS)	Total counts performed in RS sample										
P(ISS)	Probability of n-mer in ISS sample										
P(RS)	Probability of n-mer in RS sample										
CI(low)	"Lower cutoff for confidence interval (alpha= 0.02, two tailed)"										
CI(high)	"Upper cutoff for confidence interval (alpha= 0.02, two tailed)"										
Z	Z-score										
P(Z)	P-value based on Z-score										

n-mer	length	count(ISS)	count(RS)	N(ISS)	N(RS)	P(ISS)	P(RS)	CI(low)	CI(high)	Z	P(Z)
AAGG	4	19	28158	500	1799620	0.038	0.015647	0.009956	0.024508	4.02621	5.67E-05
AAGT	4	10	20912	500	1799620	0.02	0.01162	0.006882	0.019556	1.74801	0.080462
AGAA	4	11	20909	750	2699430	0.014667	0.007746	0.004585	0.013058	2.16143	0.030662
AGAG	4	14	28151	1000	3599240	0.014	0.007821	0.004971	0.012286	2.21742	0.026594
AGGC	4	32	49223	500	1799620	0.064	0.027352	0.019439	0.03836	5.02257	5.10E-07
AGGG	4	21	28200	500	1799620	0.042	0.01567	0.009975	0.024537	4.73884	2.15E-06
AGTA	4	11	20932	750	2699430	0.014667	0.007754	0.004591	0.013069	2.15759	0.03096
AGTG	4	24	28201	500	1799620	0.048	0.015671	0.009975	0.024537	5.8182	5.95E-09
ATGG	4	22	28223	500	1799620	0.044	0.015683	0.009985	0.024552	5.09436	3.50E-07
CCGA	4	11	21079	500	1799620	0.022	0.011713	0.006951	0.019672	2.13739	0.032566
CCTC	4	10	21415	750	2699430	0.013333	0.007933	0.004724	0.013292	1.66665	0.095584
CGGG	4	18	28018	500	1799620	0.036	0.015569	0.009896	0.024414	3.68909	0.000225
CGGT	4	10	21095	500	1799620	0.02	0.011722	0.006958	0.019683	1.71939	0.085544
GACC	4	12	21062	500	1799620	0.024	0.011704	0.006944	0.01966	2.55586	0.010593
GAGG	4	26	28015	750	2699430	0.034667	0.010378	0.006591	0.016305	6.56052	5.36E-11
GAGT	4	16	21335	500	1799620	0.032	0.011855	0.007058	0.019849	4.16024	3.18E-05
GATG	4	13	28209	750	2699430	0.017333	0.01045	0.006647	0.016392	1.85334	0.063834
GCGG	4	18	28223	750	2699430	0.024	0.010455	0.006651	0.016398	3.64572	0.000267
GGAG	4	19	28339	750	2699430	0.025333	0.010498	0.006685	0.016451	3.98487	6.75E-05
GGCG	4	15	28442	750	2699430	0.02	0.010536	0.006715	0.016497	2.53766	0.01116
GGCT	4	11	21106	500	1799620	0.022	0.011728	0.006962	0.019691	2.13292	0.032931
GGGC	4	43	49067	500	1799620	0.086	0.027265	0.019367	0.038259	8.06106	7.56E-16
GGGG	4	42	28105	2000	7198480	0.021	0.003904	0.002479	0.006143	12.2506	1.67E-34
GGGT	4	11	21111	500	1799620	0.022	0.011731	0.006965	0.019694	2.1321	0.032999
GGTG	4	21	27984	750	2699430	0.028	0.010367	0.006583	0.016291	4.76593	1.88E-06
GGTT	4	10	21126	500	1799620	0.02	0.011739	0.006971	0.019705	1.71456	0.086425
GTAG	4	24	55917	750	2699430	0.032	0.020714	0.015022	0.028501	2.16957	0.03004
GTGA	4	20	49127	500	1799620	0.04	0.027299	0.019395	0.038298	1.74258	0.081408
GTGG	4	44	56695	750	2699430	0.058667	0.021003	0.015265	0.028833	7.19059	6.45E-13
GTGT	4	26	49154	1000	3599240	0.026	0.013657	0.009687	0.019221	3.36222	0.000773
TAGA	4	14	27883	500	1799620	0.028	0.015494	0.009837	0.024323	2.26367	0.023595
TAGG	4	16	35047	500	1799620	0.032	0.019475	0.012986	0.02911	2.02634	0.04273
TAGT	4	13	27981	750	2699430	0.017333	0.010366	0.006582	0.016289	1.88362	0.059616
TATG	4	15	35246	500	1799620	0.03	0.019585	0.013075	0.029241	1.68025	0.09291
TGAG	4	15	35211	500	1799620	0.03	0.019566	0.013059	0.029218	1.6842	0.092142
TGGA	4	15	28342	500	1799620	0.03	0.015749	0.010036	0.024633	2.55883	0.010502
TGGC	4	40	56418	500	1799620	0.08	0.03135	0.022796	0.042973	6.24044	4.36E-10
TGGG	4	24	35134	500	1799620	0.048	0.019523	0.013025	0.029167	4.60088	4.21E-06
TGTG	4	17	34998	1000	3599240	0.017	0.009724	0.006474	0.014582	2.34428	0.019064
AAAGA	5	4	4874	500	1799620	0.008	0.002708	0.000946	0.007727	2.2758	0.022858
AACAC	5	4	4855	375	1349720	0.010667	0.003597	0.001255	0.010269	2.28582	0.022265
AAGAC	5	5	4829	375	1349720	0.013333	0.003578	0.001245	0.01024	3.16238	0.001565
AAGGC	5	8	11863	375	1349720	0.021333	0.008789	0.004406	0.017456	2.60165	0.009278
AAGGG	5	7	6614	375	1349720	0.018667	0.0049	0.001968	0.012151	3.81559	0.000136
AAGTG	5	7	6506	375	1349720	0.018667	0.00482	0.001922	0.012037	3.8693	0.000109

AATTG	5	6	6471	375	1349720	0.016	0.004794	0.001907	0.012	3.14001	0.001689
ACCAA	5	4	4776	500	1799620	0.008	0.002654	0.000918	0.007646	2.32261	0.0202
ACCGT	5	4	4779	375	1349720	0.010667	0.003541	0.001225	0.010186	2.32219	0.020223
ACCTC	5	5	4900	375	1349720	0.013333	0.00363	0.001272	0.010318	3.12256	0.001793
AGAAT	5	5	4745	375	1349720	0.013333	0.003516	0.001213	0.010149	3.21047	0.001325
AGACC	5	4	4807	375	1349720	0.010667	0.003561	0.001236	0.010216	2.30871	0.02096
AGAGG	5	7	6553	375	1349720	0.018667	0.004855	0.001942	0.012087	3.84579	0.00012
AGAGT	5	4	4903	375	1349720	0.010667	0.003633	0.001273	0.010321	2.26321	0.023623
AGGGC	5	7	11795	375	1349720	0.018667	0.008739	0.004372	0.01739	2.06499	0.038924
AGGGG	5	8	6582	375	1349720	0.021333	0.004877	0.001954	0.012117	4.57194	4.83E-06
AGGTG	5	5	6563	375	1349720	0.013333	0.004863	0.001946	0.012097	2.35725	0.018411
AGTGA	5	6	4826	500	1799620	0.012	0.002682	0.000932	0.007687	4.02655	5.66E-05
AGTGG	5	10	6699	375	1349720	0.026667	0.004963	0.002004	0.01224	5.97609	2.29E-09
ATATG	5	6	6591	375	1349720	0.016	0.004883	0.001958	0.012127	3.08677	0.002023
ATGGC	5	14	11824	375	1349720	0.037333	0.00876	0.004387	0.017418	5.93424	2.95E-09
ATTAT	5	4	4997	625	2249520	0.0064	0.002221	0.000785	0.006267	2.21807	0.02655
CAAGG	5	10	6657	375	1349720	0.026667	0.004932	0.001986	0.012196	6.00339	1.93E-09
CACCT	5	4	4692	375	1349720	0.010667	0.003476	0.001192	0.01009	2.36473	0.018043
CAGTG	5	6	6656	375	1349720	0.016	0.004931	0.001985	0.012195	3.05845	0.002225
CCAAG	5	7	6510	375	1349720	0.018667	0.004823	0.001924	0.012041	3.86729	0.00011
CCGAG	5	5	6559	375	1349720	0.013333	0.00486	0.001944	0.012093	2.35878	0.0018335
CCGAT	5	4	4808	375	1349720	0.010667	0.003562	0.001237	0.010217	2.30823	0.020987
CCTCC	5	4	4871	625	2249520	0.0064	0.002165	0.000756	0.006184	2.2766	0.02281
CGAGT	5	4	4886	375	1349720	0.010667	0.00362	0.001266	0.010303	2.27118	0.023136
CGATG	5	8	6496	375	1349720	0.021333	0.004813	0.001918	0.012027	4.61974	3.84E-06
CGGGA	5	4	4850	375	1349720	0.010667	0.003593	0.001253	0.010263	2.28819	0.022127
CGGGG	5	5	6625	375	1349720	0.013333	0.004908	0.001972	0.012162	2.33353	0.01962
CGGGT	5	4	4712	375	1349720	0.010667	0.003491	0.0012	0.010112	2.35486	0.01853
CGGTG	5	5	6503	375	1349720	0.013333	0.004818	0.001921	0.012034	2.38046	0.017291
CGGTT	5	4	4834	375	1349720	0.010667	0.003582	0.001246	0.010246	2.2958	0.021688
CTAGG	5	5	6483	375	1349720	0.013333	0.004803	0.001912	0.012013	2.38826	0.016928
CTCGG	5	5	6722	375	1349720	0.013333	0.00498	0.002014	0.012264	2.29696	0.021621
CTGGG	5	6	6289	375	1349720	0.016	0.00466	0.001831	0.011808	3.22319	0.001268
GACCA	5	4	4763	375	1349720	0.010667	0.003529	0.001219	0.010168	2.32994	0.019809
GAGGA	5	5	4770	625	2249520	0.008	0.00212	0.000733	0.006118	3.19377	0.001404
GAGGC	5	10	11892	375	1349720	0.026667	0.008811	0.00442	0.017485	3.69855	0.000217
GAGGG	5	8	6569	500	1799620	0.016	0.00365	0.001461	0.009089	4.57631	4.73E-06
GAGTA	5	4	4835	375	1349720	0.010667	0.003582	0.001247	0.010247	2.29532	0.021715
GAGTG	5	7	6725	500	1799620	0.014	0.003737	0.001511	0.009212	3.7592	0.00017
GATAT	5	4	4827	375	1349720	0.010667	0.003576	0.001244	0.010238	2.29914	0.021497
GATGG	5	10	6628	500	1799620	0.02	0.003683	0.00148	0.009136	6.01865	1.76E-09
GCACC	5	4	4837	375	1349720	0.010667	0.003584	0.001248	0.010249	2.29437	0.021769
GCGGG	5	10	6456	500	1799620	0.02	0.003587	0.001425	0.009	6.1336	8.59E-10
GCGGT	5	4	4843	375	1349720	0.010667	0.003588	0.00125	0.010256	2.29151	0.021934
GCTGG	5	8	6494	500	1799620	0.016	0.003609	0.001437	0.00903	4.61807	3.87E-06
GGACC	5	5	4794	375	1349720	0.013333	0.003552	0.001231	0.010202	3.18228	0.001461
GGAGG	5	7	6609	625	2249520	0.0112	0.002938	0.001179	0.007303	3.8143	0.000137
GGAGT	5	4	4919	375	1349720	0.010667	0.003644	0.001279	0.010339	2.25573	0.024087
GGATA	5	4	4790	375	1349720	0.010667	0.003549	0.00123	0.010198	2.31688	0.02051
GGCGG	5	11	6706	625	2249520	0.0176	0.002981	0.001204	0.007364	6.69827	2.11E-11
GGCTA	5	5	4800	375	1349720	0.013333	0.003556	0.001233	0.010209	3.17886	0.001479
GGGAG	5	6	6735	500	1799620	0.012	0.003742	0.001514	0.00922	3.02259	0.002506
GGGCA	5	4	4761	375	1349720	0.010667	0.003527	0.001219	0.010166	2.33091	0.019758
GGGCG	5	5	6655	500	1799620	0.01	0.003698	0.001488	0.009157	2.32071	0.020302
GGGGC	5	21	11858	375	1349720	0.056	0.008786	0.004404	0.017451	9.78905	1.25E-22
GGGGG	5	17	6589	1875	6748580	0.009067	0.000976	0.000391	0.002435	11.2025	3.96E-29
GGGTG	5	7	6537	500	1799620	0.014	0.003632	0.001451	0.009064	3.85142	0.000117
GGTGG	5	14	6731	625	2249520	0.0224	0.002992	0.00121	0.00738	8.87407	7.05E-19
GGTTG	5	6	6641	500	1799620	0.012	0.00369	0.001484	0.009146	3.06306	0.002191
GTAAA	5	7	11813	375	1349720	0.018667	0.008752	0.004381	0.017407	2.06066	0.039336
GTAGA	5	10	11703	375	1349720	0.026667	0.008671	0.004327	0.0173	3.75724	0.000172
GTGAG	5	9	13633	500	1799620	0.018	0.007575	0.003971	0.014404	2.68747	0.0072
GTGGC	5	20	19065	375	1349720	0.053333	0.014125	0.008168	0.024319	6.4307	1.27E-10
GTGGG	5	12	13696	500	1799620	0.024	0.007611	0.003995	0.01445	4.21516	2.50E-05
GTGTA	5	8	11859	375	1349720	0.021333	0.008786	0.004404	0.017452	2.6027	0.009249
GTGTT	5	7	11954	375	1349720	0.018667	0.008857	0.004451	0.017546	2.027	0.042663
GTTGG	5	8	13697	500	1799620	0.016	0.007611	0.003996	0.014451	2.15776	0.030947
TAAGT	5	5	6561	500	1799620	0.01	0.003646	0.001459	0.009083	2.35658	0.018444
TAATT	5	5	6619	500	1799620	0.01	0.003678	0.001477	0.009129	2.33437	0.019576
TAGAA	5	6	6486	375	1349720	0.016	0.004805	0.001913	0.012016	3.13329	0.001729

TAGAG	5	7	8315	375	1349720	0.018667	0.006161	0.002716	0.013912	3.09377	0.001976
TAGGC	5	8	13605	375	1349720	0.021333	0.01008	0.005283	0.01915	2.18095	0.029187
TAGTA	5	5	6516	625	2249520	0.008	0.002897	0.001155	0.007244	2.37311	0.017639
TATGA	5	5	6628	375	1349720	0.013333	0.004911	0.001974	0.012166	2.33239	0.01968
TATGG	5	9	8411	375	1349720	0.024	0.006232	0.00276	0.01401	4.37004	1.24E-05
TCAGT	5	5	6748	500	1799620	0.01	0.00375	0.001518	0.009231	2.28584	0.022264
TCCGA	5	6	6654	375	1349720	0.016	0.00493	0.001985	0.012193	3.05931	0.002218
TGAGG	5	9	8379	375	1349720	0.024	0.006208	0.002745	0.013977	4.38416	1.16E-05
TGGAC	5	5	6631	375	1349720	0.013333	0.004913	0.001975	0.012169	2.33125	0.01974
TGGCT	5	5	6581	500	1799620	0.01	0.003657	0.001465	0.009099	2.34889	0.018829
TGGGC	5	10	13583	375	1349720	0.026667	0.010064	0.005271	0.019129	3.22006	0.001282
TGGGG	5	12	8309	375	1349720	0.032	0.006156	0.002713	0.013906	6.39366	1.62E-10
TGTGG	5	6	8242	375	1349720	0.016	0.006106	0.002683	0.013837	2.45835	0.013958
TGTGT	5	6	6606	1000	3599240	0.006	0.001835	0.000736	0.004568	3.07548	0.002102
TTAGT	5	5	6596	500	1799620	0.01	0.003665	0.00147	0.009111	2.34315	0.019122
TTATG	5	6	8434	375	1349720	0.016	0.006249	0.00277	0.014033	2.39546	0.0166
TTGTG	5	7	8242	375	1349720	0.018667	0.006106	0.002683	0.013837	3.12078	0.001804
AACACC	6	3	1051	375	1349720	0.008	0.000779	0.000104	0.005807	5.00613	5.55E-07
AAGACC	6	4	1099	375	1349720	0.010667	0.000814	0.000112	0.005869	6.67681	2.44E-11
AAGGGC	6	3	2802	375	1349720	0.008	0.002076	0.000539	0.007955	2.51905	0.011767
AAGGGG	6	3	1567	375	1349720	0.008	0.001161	0.000208	0.006465	3.88538	0.000102
AAGTGG	6	3	1570	375	1349720	0.008	0.001163	0.000208	0.006469	3.88041	0.000104
AATCGG	6	3	1493	375	1349720	0.008	0.001106	0.000191	0.006372	4.01211	6.02E-05
AATTGT	6	3	1098	375	1349720	0.008	0.000814	0.000112	0.005868	4.87458	1.09E-06
ACACCT	6	3	1070	375	1349720	0.008	0.000793	0.000107	0.005832	4.95198	7.35E-07
ACACGT	6	3	1118	375	1349720	0.008	0.000828	0.000116	0.005894	4.82098	1.43E-06
ACCAAG	6	3	1566	375	1349720	0.008	0.00116	0.000207	0.006464	3.88703	0.000101
ACCGTT	6	3	1118	375	1349720	0.008	0.000828	0.000116	0.005894	4.82098	1.43E-06
ACCTCC	6	3	1103	375	1349720	0.008	0.000817	0.000113	0.005875	4.86105	1.17E-06
AGAGGA	6	3	1082	375	1349720	0.008	0.000802	0.000109	0.005847	4.91846	8.72E-07
AGAGTG	6	3	1607	375	1349720	0.008	0.001191	0.000217	0.006515	3.82024	0.000133
AGGGAG	6	3	1551	500	1799620	0.006	0.000862	0.000153	0.004839	3.91151	9.17E-05
AGGGGC	6	5	2965	375	1349720	0.013333	0.002197	0.00059	0.008145	4.60244	4.18E-06
AGGTGG	6	3	1554	375	1349720	0.008	0.001151	0.000205	0.006449	3.90704	9.34E-05
AGTAGC	6	3	2827	375	1349720	0.008	0.002095	0.000547	0.007984	2.50008	0.012416
AGTGAG	6	4	1587	500	1799620	0.008	0.000882	0.000159	0.004873	5.35549	8.53E-08
AGTGGC	6	4	2920	375	1349720	0.010667	0.002163	0.000576	0.008093	3.54163	0.000398
AGTGGG	6	4	1585	375	1349720	0.010667	0.001174	0.000212	0.006487	5.36048	8.30E-08
AGTGTA	6	3	1106	375	1349720	0.008	0.000819	0.000114	0.005879	4.85298	1.22E-06
ATATGG	6	6	1546	375	1349720	0.016	0.001145	0.000203	0.006439	8.48791	2.10E-17
ATCGGC	6	3	2780	375	1349720	0.008	0.00206	0.000533	0.007929	2.53593	0.011215
ATTGTG	6	3	1470	375	1349720	0.008	0.001089	0.000186	0.006343	4.05327	5.05E-05
CAAGGC	6	5	2932	375	1349720	0.013333	0.002172	0.00058	0.008107	4.63833	3.51E-06
CAAGGG	6	3	1539	375	1349720	0.008	0.00114	0.000201	0.00643	3.93234	8.41E-05
CACCTC	6	3	1030	375	1349720	0.008	0.000763	0.0001	0.00578	5.06762	4.03E-07
CCAAGG	6	7	1610	375	1349720	0.018667	0.001193	0.000217	0.006519	9.78205	1.34E-22
CCAGGC	6	3	2816	375	1349720	0.008	0.002086	0.000544	0.007971	2.5084	0.012128
CCGATG	6	3	1597	375	1349720	0.008	0.001183	0.000214	0.006502	3.83631	0.000125
CCTCGG	6	3	1594	375	1349720	0.008	0.001181	0.000214	0.006499	3.84116	0.000122
CGATGG	6	7	1576	375	1349720	0.018667	0.001168	0.00021	0.006476	9.90067	4.14E-23
CGTGGG	6	3	1512	375	1349720	0.008	0.00112	0.000195	0.006396	3.97876	6.93E-05
CGGGAT	6	3	1101	375	1349720	0.008	0.000816	0.000113	0.005872	4.86645	1.14E-06
CGGGGC	6	4	2873	375	1349720	0.010667	0.002129	0.000561	0.008038	3.585	0.000337
CGGGTG	6	4	1488	375	1349720	0.010667	0.001102	0.00019	0.006366	5.57367	2.49E-08
CGGTGG	6	4	1557	375	1349720	0.010667	0.001154	0.000205	0.006452	5.42011	5.96E-08
CGGTTG	6	3	1558	375	1349720	0.008	0.001154	0.000206	0.006454	3.90035	9.61E-05
CTAGGC	6	3	2821	375	1349720	0.008	0.00209	0.000545	0.007977	2.50462	0.012258
CTCGGG	6	3	1591	375	1349720	0.008	0.001179	0.000213	0.006495	3.84602	0.00012
CTGGGC	6	4	2641	375	1349720	0.010667	0.001957	0.000491	0.007766	3.81387	0.000137
GACATG	6	3	1553	375	1349720	0.008	0.001151	0.000204	0.006447	3.90871	9.28E-05
GACCAA	6	3	1126	375	1349720	0.008	0.000834	0.000117	0.005904	4.79991	1.59E-06
GACCTG	6	3	1547	375	1349720	0.008	0.001146	0.000203	0.00644	3.91881	8.90E-05
GAGTGG	6	4	1582	375	1349720	0.010667	0.001172	0.000211	0.006484	5.3668	8.01E-08
GATATG	6	4	1550	375	1349720	0.010667	0.001148	0.000204	0.006444	5.43525	5.47E-08
GATGGC	6	8	2823	375	1349720	0.021333	0.002092	0.000546	0.00798	8.14456	3.81E-16
GCAAGG	6	3	1535	375	1349720	0.008	0.001137	0.0002	0.006425	3.93915	8.18E-05
GCTGGA	6	3	1087	375	1349720	0.008	0.000805	0.00011	0.005854	4.90465	9.36E-07
GCTGGG	6	4	1480	375	1349720	0.010667	0.001097	0.000188	0.006356	5.59212	2.24E-08
GGAGGC	6	4	2891	375	1349720	0.010667	0.002142	0.000567	0.008059	3.56828	0.000359
GGAGGG	6	3	1564	500	1799620	0.006	0.000869	0.000155	0.004851	3.88979	0.0001

GGATAT	6	3	1087	375	1349720	0.008	0.000805	0.00011	0.005854	4.90465	9.36E-07
GGCGGG	6	6	1527	500	1799620	0.012	0.000849	0.000149	0.004816	8.54717	1.26E-17
GGCGGT	6	3	1116	375	1349720	0.008	0.000827	0.000116	0.005892	4.82628	1.39E-06
GGCTAG	6	3	1543	375	1349720	0.008	0.001143	0.000202	0.006435	3.92556	8.65E-05
GGCTGG	6	3	1475	500	1799620	0.006	0.00082	0.00014	0.004767	4.04369	5.26E-05
GGGAGG	6	3	1494	500	1799620	0.006	0.00083	0.000144	0.004785	4.00978	6.08E-05
GGGATA	6	3	1099	375	1349720	0.008	0.000814	0.000112	0.005869	4.87187	1.11E-06
GGGCGG	6	3	1606	500	1799620	0.006	0.000892	0.000162	0.004891	3.82127	0.000133
GGGGCG	6	3	1516	375	1349720	0.008	0.001123	0.000196	0.006401	3.97181	7.13E-05
GGGGGC	6	6	2780	375	1349720	0.016	0.00206	0.000533	0.007929	5.94795	2.72E-09
GGGGGG	6	10	1546	1750	6298670	0.005714	0.000245	4.35E-05	0.001385	14.5575	5.23E-48
GGGTGG	6	5	1591	500	1799620	0.01	0.000884	0.00016	0.004877	6.84782	7.50E-12
GGTGGC	6	10	2938	375	1349720	0.026667	0.002177	0.000581	0.008114	10.1586	3.03E-24
GGTTGG	6	3	1544	500	1799620	0.006	0.000858	0.000152	0.004832	3.92331	8.73E-05
GTAAGG	6	5	3287	375	1349720	0.013333	0.002435	0.000693	0.008517	4.27842	1.88E-05
GTAATT	6	5	2886	375	1349720	0.013333	0.002138	0.000565	0.008053	4.68928	2.74E-06
GTACAG	6	4	3303	375	1349720	0.010667	0.002447	0.000699	0.008535	3.21956	0.001284
GTAGAA	6	5	2839	375	1349720	0.013333	0.002103	0.000551	0.007998	4.74249	2.11E-06
GTAGAG	6	5	3255	375	1349720	0.013333	0.002412	0.000683	0.00848	4.30867	1.64E-05
GTATAC	6	3	2834	375	1349720	0.008	0.0021	0.000549	0.007992	2.49481	0.012602
GTGAGG	6	6	3343	375	1349720	0.016	0.002477	0.000712	0.008581	5.26377	1.41E-07
GTGGCT	6	5	2866	375	1349720	0.013333	0.002123	0.000559	0.00803	4.71178	2.46E-06
GTGGGC	6	4	4664	375	1349720	0.010667	0.003456	0.001182	0.01006	2.37863	0.017377
GTGGGG	6	8	3311	375	1349720	0.021333	0.002453	0.000701	0.008544	7.38201	1.56E-13
GTGTAG	6	4	3242	375	1349720	0.010667	0.002402	0.000679	0.008465	3.26747	0.001085
GTTATG	6	4	3288	375	1349720	0.010667	0.002436	0.000694	0.008518	3.23123	0.001233
GTTGGA	6	3	2833	375	1349720	0.008	0.002099	0.000549	0.007991	2.49556	0.012576
TAAGTG	6	4	1947	375	1349720	0.010667	0.001443	0.000299	0.006934	4.70162	2.58E-06
TAATTG	6	4	1923	375	1349720	0.010667	0.001425	0.000293	0.006905	4.73989	2.14E-06
TAGAAT	6	4	1532	375	1349720	0.010667	0.001135	0.0002	0.006421	5.47463	4.38E-08
TAGAGG	6	3	1937	375	1349720	0.008	0.001435	0.000296	0.006922	3.35564	0.000792
TAGAGT	6	4	1550	375	1349720	0.010667	0.001148	0.000204	0.006444	5.43525	5.47E-08
TAGGGA	6	3	1529	375	1349720	0.008	0.001133	0.000199	0.006417	3.9494	7.83E-05
TAGTAG	6	3	1945	625	2249520	0.0048	0.000865	0.000179	0.004167	3.34476	0.000824
TATGAC	6	3	1566	375	1349720	0.008	0.00116	0.000207	0.006464	3.88703	0.000101
TATGGC	6	5	3316	375	1349720	0.013333	0.002457	0.000703	0.00855	4.25135	2.12E-05
TCAGTG	6	4	2000	375	1349720	0.010667	0.001482	0.000312	0.006999	4.61938	3.85E-06
TCCGAG	6	3	1961	375	1349720	0.008	0.001453	0.000302	0.006951	3.32607	0.000881
TCCGAT	6	3	1595	375	1349720	0.008	0.001182	0.000214	0.0065	3.83954	0.000123
TCGGCG	6	3	1999	375	1349720	0.008	0.001481	0.000312	0.006998	3.28023	0.001037
TCGGGG	6	3	1959	375	1349720	0.008	0.001451	0.000302	0.006949	3.32851	0.000873
TGACAT	6	3	1546	375	1349720	0.008	0.001145	0.000203	0.006439	3.92049	8.84E-05
TGAGGC	6	4	3245	375	1349720	0.010667	0.002404	0.00068	0.008468	3.26508	0.001094
TGGACC	6	4	1539	375	1349720	0.010667	0.00114	0.000201	0.00643	5.45924	4.78E-08
TGGCGG	6	4	1988	375	1349720	0.010667	0.001473	0.000309	0.006984	4.63774	3.52E-06
TGGCTG	6	3	1996	500	1799620	0.006	0.001109	0.000233	0.005252	3.2832	0.001026
TGGGGC	6	6	3240	375	1349720	0.016	0.002401	0.000678	0.008463	5.3766	7.59E-08
TGGGGG	6	4	2007	375	1349720	0.010667	0.001487	0.000314	0.007007	4.60874	4.05E-06
TGTCAG	6	3	1914	375	1349720	0.008	0.001418	0.00029	0.006894	3.38444	0.000713
TGTGGG	6	3	1953	375	1349720	0.008	0.001447	0.0003	0.006942	3.33587	0.00085
TGTGTT	6	4	1581	375	1349720	0.010667	0.001171	0.000211	0.006482	5.36891	7.92E-08
TGTTTCG	6	3	2005	375	1349720	0.008	0.001486	0.000314	0.007005	3.2731	0.001064
TTATGA	6	4	1553	375	1349720	0.010667	0.001151	0.000204	0.006447	5.42875	5.68E-08
TTGGAC	6	3	1535	375	1349720	0.008	0.001137	0.0002	0.006425	3.93915	8.18E-05
TTGTGT	6	3	1550	375	1349720	0.008	0.001148	0.000204	0.006444	3.91375	9.09E-05
TTTGCG	6	3	2005	375	1349720	0.008	0.001486	0.000314	0.007005	3.2731	0.001064

Table S2.3. GCCS clusters derived from the ISRE enriched n-mers

Field	Description									
n-mer	The n-mer (4-6mers)									
clustID	ClusterID									
GCS	Greatest Common Substring									
Len	n-mer length									
aligned	Aligned n-mers									
wWeight	Edge weight									
count	Count of n-mer in ISS dataset									
Zscore	Z-score for n-mer									
round	Clustering round in which produced cluster									
vDegree	Vertex degree (number of other vertices attached)									
TA	Association score									

n-mer	clustID	GCS	len	aligned	wWeight	count	Zscore	round	vDegree	TA
AGGTG	1	GGTG	5	'AGGTG--	2.35725	5	2.35725	1	1	0
AGGTGG	1	GGTG	6	'AGGTGG-	3.90704	3	3.90704	1	5	0.6
CGGTG	1	GGTG	5	'CGGTG--	2.38046	5	2.38046	1	1	0
CGGTGG	1	GGTG	6	'CGGTGG-	5.42011	4	5.42011	1	5	0.6
GGGTG	1	GGTG	5	'GGGTG--	3.85142	7	3.85142	1	1	0
GGGTGG	1	GGTG	6	'GGGTGG-	6.84782	5	6.84782	1	5	0.6
GGTGG	1	GGTG	5	'-GGTGG-	8.87407	14	8.87407	1	4	1
GGTGCC	1	GGTG	6	'-GGTGCC	10.1586	10	10.1586	1	4	1
AAGGC	2	AAGG	5	'--AAGGC-	2.60165	8	2.60165	1	1	0
AAGGG	2	AAGG	5	'--AAGGG-	3.81559	7	3.81559	1	2	1
AAGGGC	2	AAGG	6	'--AAGGGC-	2.51905	3	2.51905	1	2	1
CAAGG	2	AAGG	5	'-CAAGG--	6.00339	10	6.00339	1	4	1
CAAGGC	2	AAGG	6	'-CAAGGC-	4.63833	5	4.63833	1	5	0.6
CAAGGG	2	AAGG	6	'-CAAGGG-	3.93234	3	3.93234	1	6	0.47
CCAAGG	2	AAGG	6	'CCAAGG--	9.78205	7	9.78205	1	4	1
GCAAGG	2	AAGG	6	'GCAAGG--	3.93915	3	3.93915	1	4	1
CGGCGG	3	GGCGG	6	'CGGCGG-	3.80021	3	3.80021	1	5	1
GGCGG	3	GGCGG	5	'-GGCGG-	6.69827	11	6.69827	1	5	1
GGCGGG	3	GGCGG	6	'-GGCGGG-	8.54717	6	8.54717	1	5	1
GGCGGT	3	GGCGG	6	'-GGCGGT-	4.82628	3	4.82628	1	5	1
GGGCGG	3	GGCGG	6	'GGGCGG-	3.82127	3	3.82127	1	5	1
TGGCGG	3	GGCGG	6	'TGGCGG-	4.63774	4	4.63774	1	5	1
AGGGGC	4	GGGGC	6	'AGGGGC-	4.60244	5	4.60244	1	4	1
CGGGGC	4	GGGGC	6	'CGGGGC-	3.585	4	3.585	1	4	1
GGGGC	4	GGGGC	5	'-GGGGC-	9.78905	21	9.78905	1	4	1
GGGGCG	4	GGGGC	6	'-GGGGCG-	3.97181	3	3.97181	1	4	1
TGGGGC	4	GGGGC	6	'TGGGGC-	5.3766	6	5.3766	1	4	1
CGCTGG	5	GCTGG	6	'CGCTGG-	3.97876	3	3.97876	1	4	1
GCTGG	5	GCTGG	5	'-GCTGG-	4.61807	8	4.61807	1	4	1
GCTGGA	5	GCTGG	6	'-GCTGGA-	4.90465	3	4.90465	1	4	1
GCTGGG	5	GCTGG	6	'-GCTGGG-	5.59212	4	5.59212	1	4	1
GGCTGG	5	GCTGG	6	'GGCTGG-	4.04369	3	4.04369	1	4	1
AGTGGG	6	GTGGG	6	'AGTGGG-	5.36048	4	5.36048	1	4	1
GTGGG	6	GTGGG	5	'-GTGGG-	4.21516	12	4.21516	1	4	1
GTGGGC	6	GTGGG	6	'-GTGGGC-	2.37863	4	2.37863	1	4	1
GTGGGG	6	GTGGG	6	'-GTGGGG-	7.38201	8	7.38201	1	4	1
TGTGGG	6	GTGGG	6	'TGTGGG-	3.33587	3	3.33587	1	4	1
GTAGAG	7	TAGAG	6	'GTAGAG-	4.30867	5	4.30867	1	3	1
TAGAG	7	TAGAG	5	'-TAGAG-	3.09377	7	3.09377	1	3	1
TAGAGG	7	TAGAG	6	'-TAGAGG-	3.35564	3	3.35564	1	3	1
TAGAGT	7	TAGAG	6	'-TAGAGT-	5.43525	4	5.43525	1	3	1
ATTGTG	8	TGTG	6	'ATTGTG--	4.05327	3	4.05327	1	2	1
TGTGT	8	TGTG	5	'--TGTGT-	3.07548	6	3.07548	1	2	1

TGTGTT	8	TGTG	6	'-TGTGTT	5.36891	4	5.36891	1	2	1
TTGTG	8	TGTG	5	'-TTGTG--	3.12078	7	3.12078	1	2	1
TTGTGT	8	TGTG	6	'-TTGTGT-	3.91375	3	3.91375	1	4	0.33
GGGGG	9	GGGGG	5	'-GGGGG-	11.2025	17	11.2025	1	3	1
GGGGGC	9	GGGGG	6	'-GGGGGC	5.94795	6	5.94795	1	3	1
GGGGGG	9	GGGGG	6	'-GGGGGG	14.5575	10	14.5575	1	3	1
TGGGGG	9	GGGGG	6	'TGGGGG-	4.60874	4	4.60874	1	3	1
AAGTGG	10	AGTG	6	'-AAGTGG	3.88041	3	3.88041	2	2	1
AGAGTG	10	AGTG	6	'AGAGTG-	3.82024	3	3.82024	2	2	1
AGTGG	10	AGTG	5	'--AGTGG	5.97609	10	5.97609	2	2	1
GAGTG	10	AGTG	5	'-GAGTG-	3.7592	7	3.7592	2	2	1
GAGTGG	10	AGTG	6	'-GAGTGG	5.3668	4	5.3668	2	4	0.33
AGAGG	11	AGG	5	'-AGAGG---	3.84579	7	3.84579	4	13	1
AGAGGA	11	AGG	6	'-AGAGGA--	4.91846	3	4.91846	4	13	1
AGGC	11	AGG	4	'---AGGC--	5.02257	32	5.02257	4	6	1
AGGG	11	AGG	4	'---AGGG--	4.73884	21	4.73884	4	4	1
AGGGAG	11	AGG	6	'---AGGGAG	3.91151	3	3.91151	4	8	0.86
CCAGGC	11	AGG	6	'-CCAGGC--	2.5084	3	2.5084	4	6	1
CTAGGC	11	AGG	6	'-CTAGGC--	2.50462	3	2.50462	4	6	1
GAGG	11	AGG	4	'--GAGG---	6.56052	26	6.56052	4	13	1
GAGGA	11	AGG	5	'--GAGGA--	3.19377	5	3.19377	4	13	1
GAGGAG	11	AGG	6	'--GAGGAG-	5.56098	4	5.56098	4	14	0.92
GAGGC	11	AGG	5	'--GAGGC--	3.69855	10	3.69855	4	17	0.68
GAGGG	11	AGG	5	'--GAGGG--	4.57631	8	4.57631	4	15	0.83
GAGGGG	11	AGG	6	'--GAGGGG-	5.4483	4	5.4483	4	15	0.83
GGAGG	11	AGG	5	'-GGAGG---	3.8143	7	3.8143	4	14	0.92
GGAGGC	11	AGG	6	'-GGAGGC--	3.56828	4	3.56828	4	18	0.64
GGAGGG	11	AGG	6	'-GGAGGG--	3.88979	3	3.88979	4	15	0.83
GGGAGG	11	AGG	6	'GGGAGG---	4.00978	3	4.00978	4	14	0.92
TAGGC	11	AGG	5	'--TAGGC--	2.18095	8	2.18095	4	6	1
TGAGG	11	AGG	5	'-TGAGG---	4.38416	9	4.38416	4	13	1
TGAGGC	11	AGG	6	'-TGAGGC--	3.26508	4	3.26508	4	17	0.68
ATATG	12	TATG	5	'-ATATG--	3.08677	6	3.08677	4	10	1
ATATGG	12	TATG	6	'-ATATGG-	8.48791	6	8.48791	4	10	1
GATATG	12	TATG	6	'GATATG--	5.43525	4	5.43525	4	10	1
GTTATG	12	TATG	6	'GTTATG--	3.23123	4	3.23123	4	10	1
TATG	12	TATG	4	'--TATG--	1.68025	15	1.68025	4	10	1
TATGA	12	TATG	5	'--TATGA-	2.33239	5	2.33239	4	10	1
TATGAC	12	TATG	6	'--TATGAC	3.88703	3	3.88703	4	10	1
TATGG	12	TATG	5	'--TATGG-	4.37004	9	4.37004	4	10	1
TATGGC	12	TATG	6	'--TATGGC	4.25135	5	4.25135	4	10	1
TTATG	12	TATG	5	'-TTATG--	2.39546	6	2.39546	4	10	1
TTATGA	12	TATG	6	'-TTATGA-	5.42875	4	5.42875	4	10	1
AAGACC	13	ACC	6	'AAGACC---	6.67681	4	6.67681	4	6	1
ACCAA	13	ACC	5	'---ACCAA--	2.32261	4	2.32261	4	3	1
ACCAAG	13	ACC	6	'---ACCAAG	3.88703	3	3.88703	4	3	1
AGACC	13	ACC	5	'-AGACC---	2.30871	4	2.30871	4	6	1
GACC	13	ACC	4	'--GACC---	2.55586	12	2.55586	4	6	1
GACCA	13	ACC	5	'--GACCA--	2.32994	4	2.32994	4	8	0.64
GACCAA	13	ACC	6	'--GACCAA-	4.79991	3	4.79991	4	8	0.64
GACCTG	13	ACC	6	'--GACCTG-	3.91881	3	3.91881	4	6	1
GGACC	13	ACC	5	'-GGACC---	3.18228	5	3.18228	4	6	1
AATCGG	14	CGG	6	'AATCGG--	4.01211	3	4.01211	4	5	1
ATCGGC	14	CGG	6	'-ATCGGC-	2.53593	3	2.53593	4	5	1
CGGGG	14	CGG	5	'---CGGGG	2.33353	5	2.33353	4	2	1
CTCGG	14	CGG	5	'-CTCGG--	2.29696	5	2.29696	4	5	1
CTCGGG	14	CGG	6	'-CTCGGG-	3.84602	3	3.84602	4	6	0.73
TCGGCG	14	CGG	6	'--TCGGCG	3.28023	3	3.28023	4	5	1
TCGGGG	14	CGG	6	'--TCGGGG	3.32851	3	3.32851	4	6	0.73
AGAA	15	AGAA	4	'--AGAA-	2.16143	11	2.16143	4	4	1
AGAAT	15	AGAA	5	'--AGAAT	3.21047	5	3.21047	4	4	1
GTAGAA	15	AGAA	6	'GTAGAA-	4.74249	5	4.74249	4	4	1
TAGAA	15	AGAA	5	'-TAGAA-	3.13329	6	3.13329	4	4	1
TAGAAT	15	AGAA	6	'-TAGAAT	5.47463	4	5.47463	4	4	1
AATTG	16	AATT	5	'--AATTG-	3.14001	6	3.14001	4	4	1
AATTGT	16	AATT	6	'--AATTGT	4.87458	3	4.87458	4	4	1
GTAATT	16	AATT	6	'GTAATT--	4.68928	5	4.68928	4	4	1
TAATT	16	AATT	5	'-TAATT--	2.33437	5	2.33437	4	4	1
TAATTG	16	AATT	6	'-TAATTG-	4.73989	4	4.73989	4	4	1
ACCTC	17	CCTC	5	'-ACCTC--	3.12256	5	3.12256	4	5	1








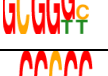








ACCTCC	17	CCTC	6	'-ACCTCC-	4.86105	3	4.86105	4	5	1
CACCTC	17	CCTC	6	'CACCTC--	5.06762	3	5.06762	4	5	1
CCTC	17	CCTC	4	'--CCTC--	1.66665	10	1.66665	4	5	1
CCTCC	17	CCTC	5	'--CCTCC-	2.2766	4	2.2766	4	5	1
CCTCGG	17	CCTC	6	'--CCTCGG	3.84116	3	3.84116	4	5	1
AGGGC	18	GGGC	5	'-AGGGC-	2.06499	7	2.06499	4	5	1
CTGGGC	18	GGGC	6	'CTGGGC-	3.81387	4	3.81387	4	5	1
GGG	18	GGGC	4	'--GGGC-	8.06106	43	8.06106	4	5	1
GGGCA	18	GGGC	5	'--GGGCA	2.33091	4	2.33091	4	5	1
GGGCG	18	GGGC	5	'--GGGCG	2.32071	5	2.32071	4	5	1
TGGGC	18	GGGC	5	'TGGGC-	3.22006	10	3.22006	4	5	1
AGTA	19	AGTA	4	'-AGTA--	2.15759	11	2.15759	4	4	1
AGTAGC	19	AGTA	6	'-AGTAGC	2.50008	3	2.50008	4	4	1
GAGTA	19	AGTA	5	'GAGTA--	2.29532	4	2.29532	4	4	1
TAGTA	19	AGTA	5	'TAGTA--	2.37311	5	2.37311	4	4	1
TAGTAG	19	AGTA	6	'TAGTAG-	3.34476	3	3.34476	4	4	1
AGTG	20	GTG	4	'-AGTG--	5.8182	24	5.8182	4	3	1
AGTGA	20	GTG	5	'-AGTGA-	4.02655	6	4.02655	4	5	0.6
AGTGAG	20	GTG	6	'-AGTGAG	5.35549	4	5.35549	4	5	0.6
CAGTG	20	GTG	5	'CAGTG--	3.05845	6	3.05845	4	3	1
GTGA	20	GTG	4	'--GTGA-	1.74258	20	1.74258	4	3	1
GTGAG	20	GTG	5	'--GTGAG	2.68747	9	2.68747	4	3	1
CGGTT	21	GGTT	5	'CGGTT--	2.2958	4	2.2958	4	4	1
CGGTTG	21	GGTT	6	'CGGTTG-	3.90035	3	3.90035	4	4	1
GGTT	21	GGTT	4	'-GGTT--	1.71456	10	1.71456	4	4	1
GGTTG	21	GGTT	5	'-GGTTG-	3.06306	6	3.06306	4	4	1
GGTTGG	21	GGTT	6	'-GGTTGG	3.92331	3	3.92331	4	4	1
AGTGTA	22	GTGT	6	'AGTGTA-	4.85298	3	4.85298	4	4	1
GTGT	22	GTGT	4	'-GTGT--	3.36222	26	3.36222	4	4	1
GTGTA	22	GTGT	5	'-GTGTA-	2.6027	8	2.6027	4	4	1
GTGTAG	22	GTGT	6	'-GTGTAG	3.26747	4	3.26747	4	4	1
GTGTT	22	GTGT	5	'-GTGTT-	2.027	7	2.027	4	4	1
CGGG	23	CGGG	4	'CGGG--	3.68909	18	3.68909	4	3	1
CGGGA	23	CGGG	5	'CGGGA-	2.28819	4	2.28819	4	3	1
CGGGT	23	CGGG	5	'CGGGT-	2.35486	4	2.35486	4	3	1
CGGGTG	23	CGGG	6	'CGGGTG	5.57367	4	5.57367	4	3	1
GCGG	24	GCGG	4	'GCGG--	3.64572	18	3.64572	4	4	1
GCGGG	24	GCGG	5	'GCGGG-	6.1336	10	6.1336	4	4	1
GCGGGC	24	GCGG	6	'GCGGGC	3.67273	4	3.67273	4	4	1
GCGGGT	24	GCGG	6	'GCGGGT	5.00325	3	5.00325	4	4	1
GCGGT	24	GCGG	5	'GCGGT-	2.29151	4	2.29151	4	4	1
AACAC	25	ACAC	5	'AACAC--	2.28582	4	2.28582	4	3	1
AACACC	25	ACAC	6	'AACACC-	5.00613	3	5.00613	4	3	1
ACACCT	25	ACAC	6	'-ACACCT	4.95198	3	4.95198	4	3	1
ACACGT	25	ACAC	6	'-ACACGT	4.82098	3	4.82098	4	3	1
AAGT	26	AAGT	4	'-AAGT--	1.74801	10	1.74801	4	3	1
AAGTG	26	AAGT	5	'-AAGTG	3.8693	7	3.8693	4	3	1
TAAGT	26	AAGT	5	'TAAGT-	2.35658	5	2.35658	4	3	1
TAAGTG	26	AAGT	6	'TAAGTG	4.70162	4	4.70162	4	3	1
AGTGGC	27	GGC	6	'AGTGGC---	3.54163	4	3.54163	5	7	1
ATGGC	27	GGC	5	'-ATGGC---	5.93424	14	5.93424	5	7	1
GATGGC	27	GGC	6	'GATGGC---	8.14456	8	8.14456	5	7	1
GGCT	27	GGC	4	'---GGCT--	2.13292	11	2.13292	5	5	1
GGCTA	27	GGC	5	'---GGCTA-	3.17886	5	3.17886	5	5	1
GGCTAG	27	GGC	6	'---GGCTAG	3.92556	3	3.92556	5	5	1
GTGGC	27	GGC	5	'-GTGGC---	6.4307	20	6.4307	5	7	1
GTGGCT	27	GGC	6	'-GTGGCT--	4.71178	5	4.71178	5	10	0.67
TGGC	27	GGC	4	'--TGGC---	6.24044	40	6.24044	5	7	1
TGGCT	27	GGC	5	'--TGGCT--	2.34889	5	2.34889	5	10	0.67
TGGCTG	27	GGC	6	'--TGGCTG-	3.2832	3	3.2832	5	10	0.67
CCGA	28	CGA	4	'-CCGA---	2.13739	11	2.13739	5	6	1
CCGAG	28	CGA	5	'-CCGAG--	2.35878	5	2.35878	5	6	1
CCGAT	28	CGA	5	'-CCGAT--	2.30823	4	2.30823	5	8	0.71
CCGATG	28	CGA	6	'-CCGATG-	3.83631	3	3.83631	5	8	0.71
CGATG	28	CGA	5	'--CGATG-	4.61974	8	4.61974	5	4	1
CGATGG	28	CGA	6	'--CGATGG	9.90067	7	9.90067	5	4	1
TCCGA	28	CGA	5	'TCCGA---	3.05931	6	3.05931	5	6	1
TCCGAG	28	CGA	6	'TCCGAG--	3.32607	3	3.32607	5	6	1
TCCGAT	28	CGA	6	'TCCGAT--	3.83954	3	3.83954	5	8	0.71
AAGGGG	29	GGG	6	'AAGGGG	3.88538	3	3.88538	5	3	1

AGGGG	29	GGG	5	'-AGGGG	4.57194	8	4.57194	5	3	1
CTGGG	29	GGG	5	'CTGGG-	3.22319	6	3.22319	5	2	1
GGGG	29	GGG	4	'--GGGG	12.2506	42	12.2506	5	3	1
TGGG	29	GGG	4	'-TGGG-	4.60088	24	4.60088	5	2	1
TGGGG	29	GGG	5	'-TGGGG	6.39366	12	6.39366	5	5	0.4
GTTGGA	30	TGGA	6	'GTTGGA--	2.49556	3	2.49556	5	4	1
TGGA	30	TGGA	4	'--TGGA--	2.55883	15	2.55883	5	4	1
TGGAC	30	TGGA	5	'--TGGAC-	2.33125	5	2.33125	5	4	1
TGGACC	30	TGGA	6	'-TGGACC	5.45924	4	5.45924	5	4	1
TTGGAC	30	TGGA	6	'-TTGGAC-	3.93915	3	3.93915	5	4	1

Table S2.4. Summary of the enriched ISRE n-mers and GCCS clustering performance

	ISRE sequences	Random Sample
Total n-mers	5376	5376
Probability $>Cl_{high}^{21}$	241	91
Clustered	193	64
% clustered	80.1%	70.33%
Number of Clusters	30	11

Table S2.5. Detailed comparison of GCCS clusters consensus motifs to known trans-acting factor binding sites

Class	Pictogram	Similar To
1		hnRNP F/H consensus binding site (GGGGG) ²³ , which functions as either a splicing enhancer or silencer ⁶³ . Contains a G-triplet, a known ISE sequence ²⁰ that is abundant in mammalian introns ¹⁶ .
1		High affinity hnRNPA1 binding site (TAGGG) identified by SELEX ²² . Contains a G-triplet, a known ISE sequence ²⁰ that is abundant in mammalian introns ¹⁶ .
1		Exhibits similarity to the hnRNP F/H and hnRNPA1 binding sites and may represent a weak binding site for these factors. Contains a G-triplet, a known ISE sequence ²⁰ that is abundant in mammalian introns ¹⁶ .
1		hnRNP A1 binding site (TAGAGT) ⁶⁴
1		High affinity hnRNP L binding site (CA-rich) identified by SELEX and an ISE element comprised of variable-length CA repeats ²⁵ . A/C-rich ESSs ⁸ .
1		CTCC and CCTCCC repeats identified by computational analysis of introns flanking skipped exons ⁶⁵ . CT-rich intronic sequences that act as PTB binding sites ^{24,66} .
1		Exhibits similarity to the hnRNP F/H and hnRNPA1 binding sites and may represent a weak binding site for these factors. Contains a G-triplet, a known ISE sequence ²⁰ that is abundant in mammalian introns ¹⁶ .
1		Exhibits similarity to the hnRNP F/H and hnRNPA1 binding sites and may represent a weak binding site for these factors.
1		Exhibits similarity to the hnRNP F/H and hnRNPA1 binding sites and may represent a weak binding site for these factors.
1		High affinity hnRNP A1 binding site (TAGGG) identified by SELEX ²² . Contains a G-triplet, a known ISE sequence ²⁰ that is abundant in mammalian introns ¹⁶ .
1		Exhibits similarity to the hnRNP F/H and hnRNPA1 binding sites and may represent a weak binding site for these factors. Contains a G-triplet, a known ISE sequence ²⁰ that is abundant in mammalian introns ¹⁶ .
2		SRp30c recognition sequence (CTGGATT) that is critical for binding ¹⁴ .
2		SRp40 binding site (ACAAG) ²⁷
2		SC35 binding site (AGGAGAT) ²⁸ . A purine-rich element (AGGG) identified in introns flanking skipped exons ⁶⁵ .
2		Sam68 binding site (TAAA) ^{67,68} .
2		9G8 high-affinity binding site (GAC) identified by SELEX ²⁸ .

2		SF2/ASF high-affinity binding site (GAAGAA) identified by SELEX ²⁶ . Tra2β high-affinity binding site (GAA) _n identified by SELEX ²⁹ .
3		Major 5' ss consensus sequence (GT[A/G]AGT) ⁶⁹ .
3		Major 5' ss consensus sequence (GT[A/G]AGT) ⁶⁹ .
3		Major 5' ss consensus sequence (GT[A/G]AGT) ⁶⁹ .
3		Major 5' ss consensus sequence (GT[A/G]AGT) ⁶⁹ . hnRNP G binding motif (AAGT) ³¹ .
4		CELF/Bruno-like family of proteins that bind GT repeats with high affinity ³⁴ . CUG-BP1 binding sites consisting of TGT-repeats ³⁴ . hnRNP M binding sites consisting of poly(G) and poly(T) homopolymers ³⁵ .
4		CELF/Bruno-like family of proteins that bind GT repeats with high affinity ³⁴ .
4		CELF/Bruno-like family of proteins that bind GT repeats with high affinity ³⁴ . CUG-BP1 binding sites consisting of TGT-repeats ³⁴ .
4		CELF/Bruno-like family of proteins that bind GT repeats with high affinity ³⁴ .
5		N/A. A novel regulatory element.
5		N/A. A novel regulatory element.
5		CELF binding site (GT repeats) ³⁴ .
5		hnRNPA1 binding site (TAGGG) ²² .
5		N/A. A novel regulatory element.

Table S2.6. Overlap of enriched hexamers with extended recovered ISRE sequences

Extended ISS sequence	Enriched Hexamers
GTTCGAATCTCTCCAGTGC	
GTCCTACGCTCATTATTGC	
GTTCTTCTCTTCTCTTCGC	
GTTGTTTCGCACCGCTGGGC	CGCTGG CTGGGC GCTGGG TGTTTCG
GTTGTTTCGCACCACTGAGC	TGTTTCG
GTAGTCACCTATTATAGGC	
GTGTTAACCAACGATGGGC	CGATGG
GTGGTATCGAAAGTTGTGC	
GTTACATCCAGAAGTCGGC	
GTTACATCCCTCGGTTGGC	CCTCGG CGGTTG GGTTCG
GTTGGACCAGGCGTACGGC	CCAGGC GTTGGA TGGACC TTGGAC
GTTGGACACGTCAGTCAGC	ACACGT GTTGGA TTGGAC
GTCACACGTGAGAGAGAGC	ACACGT
GTGAAGGGCGACAGATAGC	AAGGGC
GTAGAACGCTGGATTAAGC	CGCTGG GCTGGA GTAGAA
GTTTACTTTAAGGATAAGC	
GTATACGGAAAGGCCTTGC	GTATAC
GTGTGCTTATATGGGTTGC	ATATGG
GTTTAGTCCCATTCCGAGC	TCCGAG
GTCCACTTCGGTTCCTGC	CGGTTG
GTACGTCCGTCGTGGATGC	
GTACCTCGAGGTCTGAAGC	
GTACCTCAGGCTCTGAAGC	
GTAAGGCTAGTTTAGTAGC	AGTAGC GGCTAG
GTAAGGCTAGATTAGTAGC	AGTAGC GGCTAG
GTAGAGGAGTCGTGTCAGC	AGAGGA GTAGAG TAGAGG TGTCAG
GTAGTGGAATCGTATCAGC	

GTGGTCGAGTCGCAAGGGC	AAGGGC	CAAGGG	GCAAGG			
GTATTCCAGCTGGAGCTGC	GCTGGA					
GTAGTATATGGTGAGGAGC	ATATGG	GTGAGG				
GTGCCGAGTAAAGTGTAGC	AGTGTA	GTAAAG	GTGTAG			
GTTCTGACTCAATAGTAGC	AGTAGC					
GTCTTGAGTACCCCCGAGC						
GTCATGCACCGACCAAGGC	ACCAAG	CAAGGC	CCAAGG	GACCAA		
GTAATTGTGTTTGTGATGC	AATTGT	ATTGTG	GTAATT	TAATTG	TGTGTT	TTGTGT
GTGACTGTGTTAGGCGGGC	GGCGGG	TGTGTT				
GTAATTGGGTTTGGGGGGC	GGGGGC	GGGGGG	GTAATT	TAATTG	TGGGGG	
GTAATTGTGTTCCGTGGGC	AATTGT	ATTGTG	CGGTGG	GTAATT	GTGGGC	TAATTG
GTAATTGTGTTTGGCGGGC	TGTGTT	TGTTTCG	TTGTGT			
GTTATGACATGTGGGGAGC	AATTGT	ATTGTG	GGCGGG	GTAATT	TAATTG	TGGCGG
GTTATGACCTGTGGGGGGC	TGTGTT	TTGTGT				
GTTATGACATGTGGGGGGC	GACATG	GTGGGG	GTTATG	TATGAC	TGACAT	TGTGGG
GTTATGACATGTGGGGGGC	TTATGA					
GTTATGACATGTGGGGGGC	GGGGGC	GGGGGG	GTGGGG	GTTATG	TATGAC	TGGGGG
GTTATGACATGTGGGGGGC	TGTGGG	TTATGA				
GTTATGACATGTGGGGGGC	GACATG	GGGGGC	GGGGGG	GTGGGG	GTTATG	TATGAC
GTTATGACATGTGGGGGGC	TGACAT	TGGGGG	TGTGGG	TTATGA		
GTCAATTGAGTTGGTGTGC						
GTCGATGGGGCAGGGGAGC	CGATGG	TGGGGC				
GTCAGTGAACTTTGCGAGC	TCAGTG	TTTGCG				
GTCCTTGCTCCTGACATGC	GACATG	TGACAT				
GTCCGAGTGCGACGGTGGC	CGGTGG	GGTGGC	TCCGAG			
GTGAGTGGCCTAGGGAGGC	AGGGAG	AGTGGC	GAGTGG	GGAGGC	GGGAGG	TAGGGA
GTGGCTGGGCTAGGATGGC	TAGTAG					
GTGATATGGCGAGGGTGGC	CTGGGC	GATGGC	GCTGGG	GGCTAG	GGCTGG	GTGGCT
GTGATATGGCGAGGGTGGC	TGGCTG					
GTAAGTGGGCACGGTTGGC	ATATGG	GATATG	GGGTGG	GGTGGC	TATGGC	
GTAAGTGGGCACGGTTGGC	AAGTGG	AGTGGG	CGGTTG	GGTTGG	GTGGGC	TAAGTG
GTAGGTAGCCACCGTTGGC	ACCGTT					
GTGGGGGGGTCACTTAGGC	GTGGGG	TGGGGG				
GTTGGTTGGACCCGTAGGC	GGTTGG	GTTGGA	TGGACC	TTGGAC		
GTCCTATGGTTCCTCGGC	CCTCGG					
GTCAGAGGAGTCTCTAGGC	AGAGGA	CTAGGC				

GTTTATGGAGTTCCTAGGC	CTAGGC					
GTAAATAGAGGCCCCAGGC	CCAGGC	TAGAGG				
GTCTAGTAACCAGCCAGGC	CCAGGC					
GTCTAAGCACCCTGAGGC	TGAGGC					
GTTGTTTTGCGTCCAAGGC	CAAGGC	CCAAGG	TTTGCG			
GTCATGTGACGACCAAGGC	ACCAAG	CAAGGC	CCAAGG	GACCAA	TGTCAG	
GTCATGGACCGACCAAGGC	ACCAAG	CAAGGC	CCAAGG	GACCAA	TGGACC	
GTTATGCCTCCCCGATAGC	GTTATG					
GTCTGAAGAACCCCAAGGGC	AAGGGC	CAAGGG	CCAAGG			
GTCTGAGAAACCGAGGGC	GGAGGG					
GTCCGAGGAACCATAGGGC	TCCGAG					
GTCTATCTCCTTCTATGGC	TATGGC					
GTTTAACACCTCCCAAGGC	AACACC	ACACCT	ACCTCC	CAAGGC	CACCTC	CCAAGG
GTCAAAGACCTGCGATGGC	AAGACC	CGATGG	GACCTG	GATGGC		
GTCAAACACGTCCGATGGC	ACACGT	CCGATG	CGATGG	GATGGC	TCCGAT	
GTCTAACACCTCCGATGGC	AACACC	ACACCT	ACCTCC	CACCTC	CCGATG	CGATGG
GTCAAACACCTCCGATGGC	AACACC	ACACCT	ACCTCC	CACCTC	CCGATG	CGATGG
GTGTGGCTATGAATTTGGC	GATGGC	TCCGAT				
	GTGGCT					
GTGTGGCTAAGAATTGGGC	GTGGCT					
GTGGCTGGAAGACCTGCGC	AAGACC	GACCTG	GCTGGA	GGCTGG	GTGGCT	TGGCTG
GTGTAAAGGGGTGTCAGTGC	GTAAAG	TCAGTG	TGTCAG			
GTATTAATAATACTGGGGC	TGGGGC					
GTGTTAATAGCGCGGGAGC						
GTTTGTAAGGTGCTGGGGC	GCTGGG	TGGGGC				
GTTGTGGTCGCGACCTGGC	GACCTG					
GTGGCGGTGTCAGTACAGGC	GGCGGT	GTACAG	TGGCGG			
GTGTTGTGAAAGAGGAGGC	AGAGGA	GGAGGC				
GTGGTGGCAGACACGATGC	GGTGGC					

GTGCGGTTTGCGGGCGGGC	GGCGGG	GGGCGG	TTTGCG			
GTGGGGCGCGCGGGGGGGC	GGGGCG	GGGGGC	GGGGGG	GTGGGG	TGGGGC	
GTGAGGGCAGTCCGTGGGC	GTGAGG	GTGGGC				
GTGACGGGTGCCTCGGGGC	CCTCGG	CGGGGC	CGGGTG	CTCGGG	TCGGGG	
GTTAGGTGTGTCTCGGGGC	CGGGGC	CTCGGG	TCGGGG			
GTGACGTGTGTCTCGGGGC	CGGGGC	CTCGGG	TCGGGG			
GTGACGGAGCCGTCTGGGC	CTGGGC					
GTGCATGGCCCCGCTGGGC	CGCTGG	CTGGGC	GCTGGG			
GTGCAAGGTCCCTCTAGGC	CTAGGC	GCAAGG				
GTGCACTAGAATCTGAGGC	TAGAAT	TGAGGC				
GTGCAGTACGGGCTTAGGC						
GTCGAGCGGCTTTAGAGGC	TAGAGG					
GTAGAGTGGGGCGGGTGGC	AGAGTG	AGTGGG	CGGGTG	GAGTGG	GGCGGG	GGGCGG
	GGGGCG	GGGTGG	GGTGGC	GTAGAG	GTGGGG	TAGAGT
	TGGGGC					
GTATAGTGGCGGTGGAGGC	AGTGGC	CGGTGG	GGAGGC	GGCGGT	TGGCGG	
GTAGAGTGGCGGTGGAGGC	AGAGTG	AGTGGC	CGGTGG	GAGTGG	GGAGGC	GGCGGT
	TAGAGT	TGGCGG	GTAGAG			
GTAGAATGGACCGTGAGGC	GTAGAA	GTGAGG	TAGAAT	TGAGGC	TGGACC	
GTGGAGTGGCTGTTTCGGC	AGTGGC	GAGTGG	GAGTGG	GGCTGG	GTGGCT	TGGCTG
GTGTACAGCGGAGAGGGGC	AGGGGC	GTACAG				
GTGTACGGTGCAGAGGGGC	AGGGGC					
GTGTAGTGTAGGGAGGGGC	AGGGAG	AGGGGC	AGTGTA	GGAGGG	GGGAGG	GTGTAG
	TAGTAG	TAGGGA				
GTGAAGTGTAGGGAGGGGC	AGGGAG	AGGGGC	AGTGTA	GGAGGG	GGGAGG	GTGTAG
	TAGTAG	TAGGGA				
GTATACCGTTCAGTGGGGC	ACCGTT	AGTGGG	GTATAC	GTGGGG	TCAGTG	TGGGGC
GTATACCGTTCAGTGAGGC	ACCGTT	AGTGAG	GTATAC	GTGAGG	TCAGTG	TGAGGC
GTAAAGGGGCAAGGTGGGC	AAGGGG	AGGGGC	AGGTGG	GCAAGG	GTAAAG	GTGGGC
GTAGAGTGCGAAGCGGGGC	AGAGTG	CGGGGC	GTAGAG	TAGAGT		
GTACAGTGCTAAGTAGGGC	GTACAG					
GTGTAAATCGGCGGGTGGC	AATCGG	ATCGGC	CGGGTG	GGCGGG	GGGTGG	GGTGGC
	TCGGCG	GGTGGC				
GTGGAAATCGGCGGATGGC	AATCGG	ATCGGC	GATGGC	TCGGCG		

GTGGCAATCGGCGGGTGGC	AATCGG	ATCGGC	CGGGTG	GGCGGG	GGGTGG	GGTGGC
	TCGGCG					
GTAAAGAACGGGATATGGC	ATATGG	CGGGAT	GATATG	GGATAT	GGGATA	GTAAAG
	TATGGC					
GTCAAGACCGGGATATGGC	AAGACC	ATATGG	CGGGAT	GATATG	GGATAT	GGGATA
	TATGGC					
GTAAAGACCGGGATATGGC	AAGACC	ATATGG	CGGGAT	GATATG	GGATAT	GGGATA
	TATGGC	GTAAAG				
GTAATTATTAGTCGATGGC	CGATGG	GATGGC	GTAATT			
GTGCTTAGTGAGTGATGGC	AGTGAG	GATGGC				
GTACAGGCCAAGGGGGGGC	AAGGGG	CAAGGG	CCAAGG	GGGGGC	GGGGGG	GTACAG
GTAGAAGACAAGTGGTGGC	AAGTGG	GGTGGC	GTAGAA			
GTGGTTGAAGGGGGGCGGC	AAGGGG	GGGCGG	GGGGCG	GGGGGC	GGGGGG	
GTACATTATGAGGGTCGGC	TTATGA					
GTAGAGTAAGTGAGGTGGC	AGGTGG	AGTGAG	GGTGGC	GTAGAG	GTGAGG	TAAGTG
	TAGAGT					
GTAGAATAAGTGAGGTGGC	AGGTGG	AGTGAG	GGTGGC	GTAGAA	GTGAGG	TAAGTG
	TAGAAT					
GTAGAATAAGTGGGGTGGC	AAGTGG	AGTGGG	GGGTGG	GGTGGC	GTAGAA	GTGGGG
	TAGAAT	TAAGTG				

Table S2.7. Primer and oligonucleotide sequences

Name	Primer Sequence (5' - 3')
Ex6	CATGGACGAGCTGTACGTTAACATAATTCCCCCACCACCTC
Ex8	CGCTCG AGCACATACGCCTCACATACATTTTG
GFP1	GCGGTACCATGGTGAGCAAGGGCG
GFP2	GGTGGTGGGGGAATTATGTTAACGTACAGCTCGTCCATGCC
ECmutF	CTTTTAAACATCCATATAAAGCTATCGATATCTAGCTATCGAT GTCTATATAGCTATTTTTTTTAACT
ECmutR	AGTTAAAAAAATAGCTATATAGACATCGATAGCTAGATATCG ATAGCTTTATATGGATGTAAAAAAG
PmlImutF	CATTATGAAAGTGAATCTTACTTTTGTAACACGTGATGGTTTG TGGAAAACAAATGTTTTTGAA
PmlImutR	TTCAAAAACATTTGTTTTCCACAAACCATCACGTGTTACAAAAG TAAGATTCACCTTCATAATG
BamHIImutF	CTTTTGTAACACGTGATGGTTTGTGGGATCCAAATGTTTTTGAA CAGTTAAAAAGTTC
BamHIImutR	GAACTTTTTAACTGTTCAAAAACATTTGGATCCCACAAACCATC ACGTGTTACAAAAG
P2	TAAGAAGCTAAAGAGCCTCACTCATGTGGTTTTATGCAGC
P3	TGAGGCTCTTTAGCTTCTTA
P4	AGATAGAGAGGTCAGCGATTTGCAATTCTGAGGTGTTAAA
P5	AATCGCTGACCTCTCTATCT
Ex17	CATGGACGAGCTGTACGTTAACATGCTCGTGTACAAGTTTGCC
Ex19	CGCTCGAGAAGTACTTACCTCATTGAGCATTTTTTC
GFP3	GCAAACCTTGTAACACGAGCATGTTAACGTACAGCTCGTCCATGCC
ECmutF2	TTTAGCTTCTTAGGATATCACTTATCGATTTTGTTTTCAAC
ECmutR2	GTTGAAAACAAAATCGATAAGTGATATCCTAAGAAGCTAAA
ISStemp	GCGCGATATCGATCAGT (N ₁₅) GCATCATCGATGCGC
Lib1	GCGCGATATCGATCAGT
Lib2	GCGCATCGATGATGC
Lib3	GAAACAAAATGCTTTTTAACATCCATA
Lib4	GGAAAATAAAAGGAAGTTAAAAAAAATAGC
SMN1cDNA	TAGAAGGCACAGTCGAGG

Table S2.8. Plasmid constructs used in this work

Name	Description
pCS238	GFP-SMN1. Contains the wild-type SMN1 mini-gene fused to the N-terminus of GFP. Positive control used for all flow cytometry analysis.
pCS516	SMN1 NMD-based reporter construct. Contains the SMN1 mini-gene with a PTC in exon 7 fused to the N-terminus of GFP. Recovered ISREs as well as control ISS were inserted into this construct.
pCS517	SMN1 NMD-based containing random 15-mer. Negative control used for all flow cytometry analysis.
pCS990	GFP-BRCA1. Contains the wild-type BRCA1 mini-gene fused to the N-terminus of GFP. Positive control used for flow cytometry analysis.
pCS1008	BRCA1 NMD-based reporter construct. Contains the BRCA1 mini-gene with a PTC in exon 17 fused to the N-terminus of GFP. Recovered ISREs were inserted into this construct.
pCS668	U2AF65 binding site inserted into pCS516.
pCS669	hnRNP H binding site inserted into pCS516.
pCS670	PTB (1) binding site inserted into pCS516.
pCS667	PTB (2) binding site inserted into pCS516.

Table S2.9. Primer sequences for SMN1 transcript isoform analysis through qRT-PCR

Name	Forward Primer (5' - 3')	Reverse Primer (5' - 3')	Isoform
Pair 1	CTCCCATATGTCCAGATCT	AGCATTTTGTTCACAAGACA	Ex 6 and Int 6
Pair 2	CACTAGTAGGCAGACCAG	CAGTTATCTTCTATAACGCTTCAC	Int 7 and Ex 8
Pair 3	TAAATTAAGGAGAAATGCT	GGTTTTTCAAAAGAGTCCAGTAA	Ex 7/8 and Ex 8
Pair 4	TGAGCAAAGACCCCAA	CCAGCATTTCCATATAATAG	GFP and Ex 6/8
Pair 5	TGAGCAAAGACCCCAA	TGATAGCCACTCATGTACC	GFP and Ex 6
Pair 6	CAAAGATGGTCAAGGTCGCAAG	GGCGATGTCAATAGGACTCC	HPRT

Table S2.10. Primer sequences for endogenous transcript isoform analysis through qRT-PCR

Name	Gene	Hexamer	Sequence (5' – 3')	Isoform	Type of Alternative splicing
Fw.ADD3ex15_16	ADD3	ACCTCC	TGAAAAATTAGAAGAA AACCATGAGC	Exon15/16	cassette
Fw.ADD3ex14_16	ADD3	ACCTCC	GGCC TAG AAGAAA ACCATG AGC	Exon14/16	cassette
Rv.ADD3ex16	ADD3	ACCTCC	CTTCGATTTTCTCTGGA GACT	ADD3 cDNA, Exon15/16, Exon 14/16	cassette
Fw.hnRNPCex1_3	HNRNPC	ACCTCC	CCC CTT CTT GTT TTC GGC TTT	Exon1/3	cassette
Fw.hnRNPCex2_3	HNRNPC	ACCTCC	CTT CAGCTACATTTT C GGCTTT	Exon2/3	cassette
Rv.hnRNPCex3	HNRNPC	ACCTCC	CGAAAAGATTGCCTCC ACAT	hnRNPC cDNA and Exon1/3, Exon 2/3	cassette
Fw.CLK3ex4	CLK3	GGGGGG	CCGTGACAGCGATACA TAC	Exon 4/5, Exon4/6	cassette
Rv.CLK3ex4_5	CLK3	GGGGGG	GTTGGCTTCTCGAGGAG G	Exon 4/5	cassette
Rv.CLK3ex4_6	CLK3	GGGGGG	CCACAATCTCATCGAG GAGG	Exon 4/6	cassette
Rv.CLK3cDNA	CLK3	GGGGGG	CAAGCACTCCACCACCT	CLK3 cDNA	cassette
Fw.CADPSex16	CADPS	GGGGGG	GAAAGATATTGTTACCC CAGT	Exon 16/19, Exon16/18, Exon16/17	mutually exclusive
Rv.CADPSex16_18	CADPS	GGGGGG	CCTTTTGATTCTCTTCG ATTTTG	Exon16/18,	mutually exclusive
Rv.CADPSex16_19	CADPS	GGGGGG	GGCCTACATTTTCTTCG ATTTTG	Exon 16/19	mutually exclusive
Rv.CADPSex16_17	CADPS	GGGGGG	CTCTCTTTTCCCTTCG ATTTTG	Exon16/17	mutually exclusive
Rv.CADPScDNA	CADPS	GGGGGG	AAG CTT TTT GGC AGG AGT GA	CADPS cDNA	mutually exclusive
Fw.c6orf60ex15_16	C6orf60	GTAGAA	CTTTACAAGTGTCATTA GAAGAAATG	Exon 15/16	cassette
Fw.c6orf60ex14_16	C6orf60	GTAGAA	CCA ACA GAT AAG ATT AGA AGA AAT GG	Exon 14/16	cassette
Rv.c6orf60ex16	C6orf60	GTAGAA	GATCTGGTCTCTTTCTG TAAGC	C6orf60 cDNA, Exon 15/16, Exon 14/16	cassette
Fw.RREB1ex11_12	RREB1	GTAGAA	GATAGCACAGACAGTC AGTCG	Exon11/12	cassette
Fw.RREB1ex10_12	RREB1	GTAGAA	ACA CAC ACT GAC AGT CAG TCG	Exon10/12	cassette
Rv.RREB1ex12	RREB1	GTAGAA	CTCCTCCTCCGGCTCAT	RREB1 cDNA, Exon11/12,	cassette

				Exon10/12	
Fw.MADDex35	MADD	GCTGGG	AGTTCCTGTGCGAC	Exon35/36, Exon35/37	cassette
Rv.MADDex35_36	MADD	GCTGGG	TCTATGAAAACCTGATT GTGCA	Exon35/36	cassette
Rv.MADDex35_37	MADD	GCTGGG	TAATTTTCAGGAACTGAT TGTGCA	Exon35/37	cassette
Rv.MADDcDNA	MADD	GCTGGG	TAGTACAGCTCCCGAC ACTT	MADD cDNA	cassette
Fw.CAMK2Gex13_1 4	CAMK2G	GCTGGG	CGGGCAAGCTGCCAAA AG	Exon13/14	cassette
Fw.CAMK2Gex12_1 4	CAMK2G	GCTGGG	GAA CTT CTC AGC TGC CAA AAG	Exon12/14	cassette
Rv.CAMK2Gex14	CAMK2G	GCTGGG	TTGACACCGCCATCCG	CAMK2G cDNA, Exon13/14, Exon12/14	cassette
Fw.A2BP1ex15_17	A2BP1	ATATGG	GCAGACATTTATGGTG GTTATG	Exon15/17	mutually exclusive
Fw.A2BP1ex16_17	A2BP1	ATATGG	TAA ATT GCT GCA GGG TGG TTA TG	Exon16/17	mutually exclusive
Rv.A2BP1ex17	A2BP1	ATATGG	CTGTCACTGTAGGCAGC G	A2BP1 cDNA, Exon15/17, Exon16/17	mutually exclusive
Fw.HNRNPA2B1ex1	HNRNPA2B1	ATATGG	CTCTAGCGGCAGTAGC A	Exon1/2, Exon1/3	cassette
Rv. HNRNPA2B1ex1_2	HNRNPA2B1	ATATGG	GTTTCTAAAGTTTCTC CATCGCG	Exon1/2	cassette
Rv. HNRNPA2B1ex1_3	HNRNPA2B1	ATATGG	GTTCCTTTTCTCTCTCC ATCGC	Exon1/3	cassette
Rv. HNRNPA2B1cDNA	HNRNPA2B1	ATATGG	CCTCAAACCTTTCTTCTG TGG	HNRNPA2B1 cDNA, Exon1/2, Exon1/3	cassette

Table S2.11. siRNA duplex sequences

Targeted mRNA	Target Sequence (5' - 3')
hnRNP H	GAUCCACCACGAAAGCUUA
hnRNP A1	CAACUUCGGUCGUGGAGGA
PTB	CGUCAAAGGAUUCAAGUUC
CUG-BP1	GAGCCAACCUGUUCAUCUA
SF2/ASF	CGUGGAGUUUGUACGGAAA