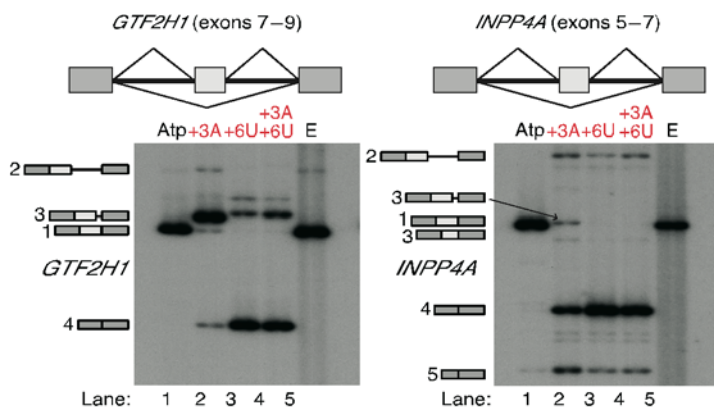


## SUPPLEMENTARY MATERIAL

TITLE: Recognition of atypical 5' splice sites by shifted base-pairing to U1 snRNA

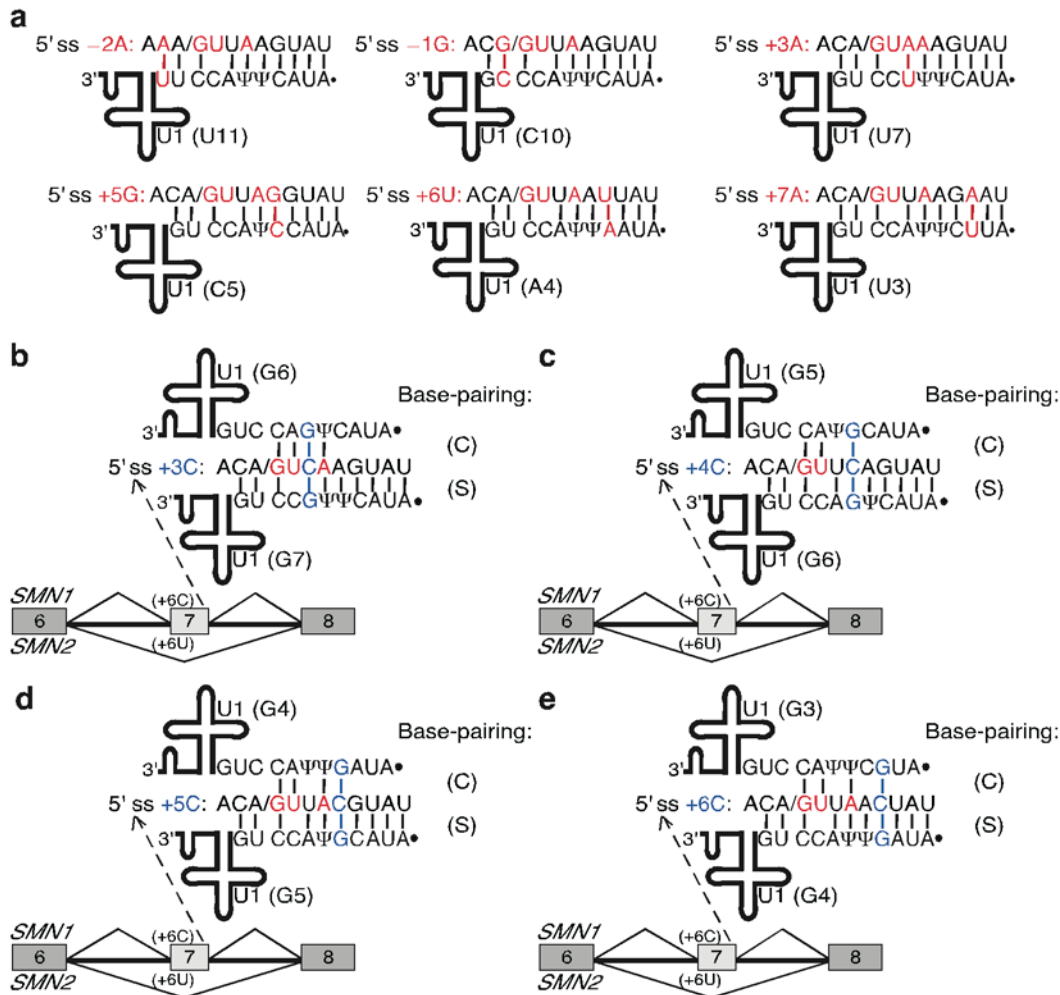
AUTHORS: Xavier Roca and Adrian R. Krainer

## SUPPLEMENTARY FIGURES



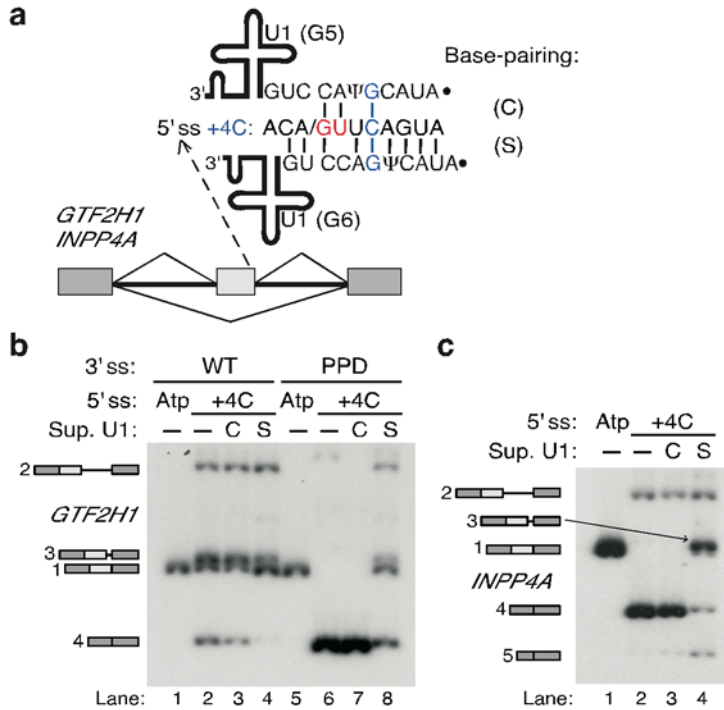
**Supplementary Figure 1 Detection of endogenous and transfected *GTF2H1* and *INPP4A* mRNAs.** cDNAs from HeLa cells transfected with the various minigenes (lanes 1-4) or mock-transfected (lane 5) were amplified using primers located in the exons flanking the exon containing the atypical 5' ss. The RT-PCR products from the mock-transfected samples correspond to the endogenous *GTF2H1* and *INPP4A* transcripts (E, lane 5), confirming the efficient use of the atypical 5' ss. A small amount of intron 8 retention was seen for endogenous *GTF2H1* (left panel, lane 5). Although the primers can amplify the transcripts from both transfected and endogenous genes, in transiently transfected cells the RT-PCR products are derived mostly from the minigenes, which are expressed at higher levels. In the mock-transfected samples (E, lane 5), ~ 50-fold more total cDNA template was added to the reaction. The identity of the RT-PCR products, schematically depicted on

the left of each panel, is described in the **Fig. 1** legend. For the *GTF2H1* minigene, the mutations at the exon 8 5' splice site activated cryptic 5' splice sites at positions +20 (GAG/GUGCAG) and +25 (GCA/GUAACU) in intron 8 (band #3). For the *INPP4A* minigene, the analogous mutations activated a cryptic 5' splice site 31 nucleotides upstream of the 5' splice site in exon 6 (CGU/GUAUGA), and to a small extent another cryptic 5' splice site at position +176 (AGG/GUAAGA) in intron 6. The arrow indicates a cryptic 5' splice site at position +5 (AAA/GUAAUG) of *INPP4A* intron 6 that was activated only in the context of the +3A mutation. In addition, band #5 in *INPP4A* corresponds to an mRNA species generated by skipping of exon 6 and use of a cryptic 5' splice site (GCG/GUGAGU) 63 nucleotides upstream of the 5' splice site of exon 5.

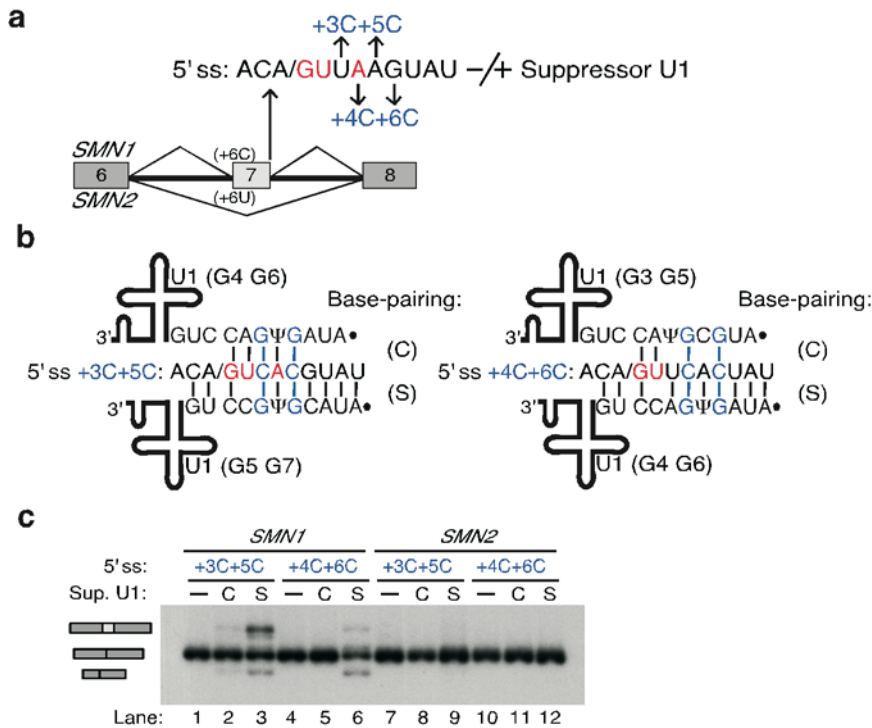


**Supplementary Figure 2 Schematic of the base-pairing profiles for the mutant 5' ss used in Figs. 2 and 3.** **a**, Base-pairing profiles of mutants  $-2A$ ,  $-1G$ ,  $+3A$ ,  $+5G$ ,  $+6U$  and  $+7A$  with their corresponding suppressor U1 snRNAs restoring base-pairing in the shifted register (**Fig. 2**). Mutations at the 5' ss or U1 are highlighted in red. For these mutant 5' ss, there is no suppressor U1 in the canonical register, because the mutant nucleotide already base-pairs to endogenous U1 in the canonical register. **b-e**, Diagrams show the atypical 5' ss carrying the  $+3C$  (**b**),  $+4C$  (**c**),  $+5C$  (**d**) and  $+6C$  (**e**) mutations, along with the corresponding suppressor U1 snRNAs in the canonical (C) or shifted (S) base-pairing registers (**Fig. 3**). The 5' ss mutations and the compensatory U1 mutations are highlighted in blue. Note that some of the suppressor U1 snRNAs were used for more than one 5' ss

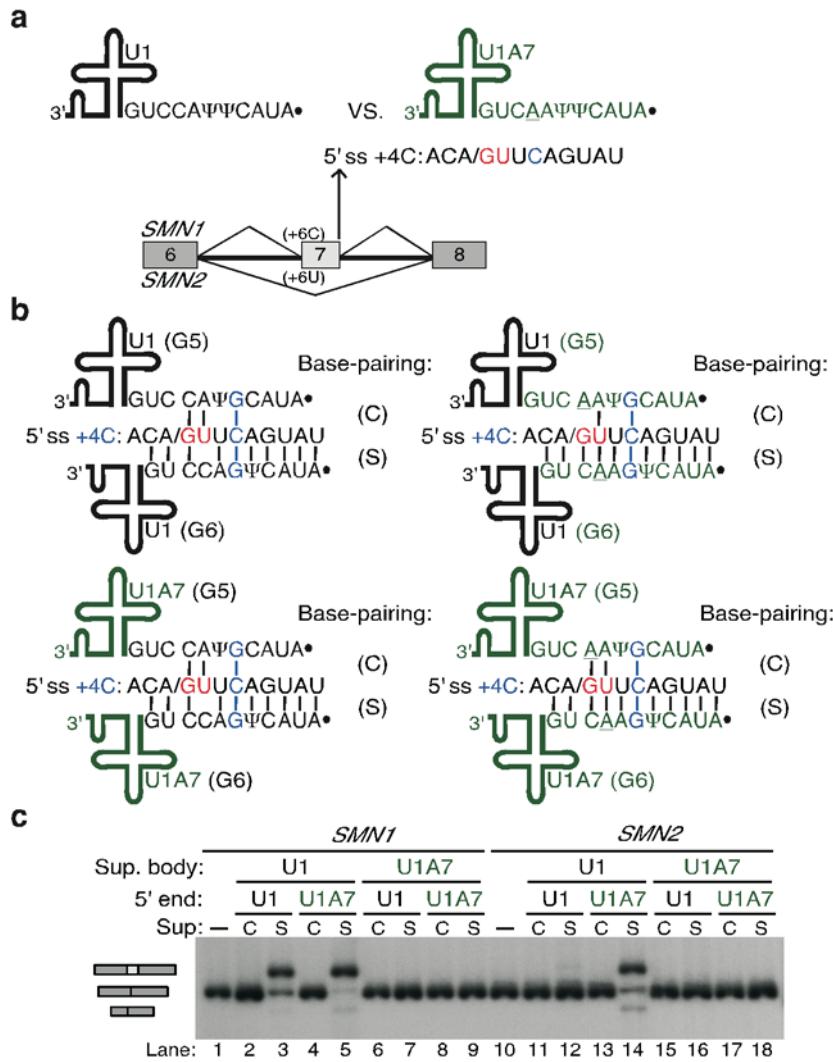
mutation: U1 (G5) compensates for 5' ss mutation +4C in the canonical register, and for mutation +5C in the shifted register. **Fig. 3** shows that U1 (G5) rescued exon 7 inclusion for mutant +5C but not for +4C, and that U1 (G6) rescued mutation +4C (shifted register) but not +3C (canonical register). These results indicate that the distinct effects of the suppressor U1s are not due to their relative expression levels, and further strengthen the demonstration of shifted base-pairing between the atypical 5' ss and U1.



**Supplementary Figure 3 Suppressor U1 snRNA analysis for atypical 5' ss in their natural context.** **a**, Base-pairing profiles of the +4C mutant 5' ss in the *INPP4A* and *GTF2H1* minigenes with the corresponding suppressor U1 snRNAs carrying compensatory mutations in the canonical (C) or shifted (S) registers. Mutant nucleotides and compensatory U1 mutations are highlighted in blue. **b**, Suppressor U1 analysis of the +4C mutant 5' ss in the *GTF2H1* minigene. The top labels indicate the version of the 3' ss upstream exon 8, the version of the atypical 5' ss, and the mock or suppressor U1 snRNA used. In lanes 5-8, the 3' ss upstream of exon 8 was weakened by a point mutation at the polypyrimidine tract ('polypyrimidine tract down mutation' or PPD), so as to compromise exon 8 inclusion. Suppressor U1 in the shifted but not in the canonical register rescued exon 8 inclusion in both the wild-type and the PPD minigenes (lanes 2-4, 6-8). **c**, Suppressor U1 analysis in the *INPP4A* minigene. Whereas the suppressor U1 in the canonical register did not rescue exon 6 inclusion (lane 3), the suppressor in the shifted register did, albeit weakly (lane 4).



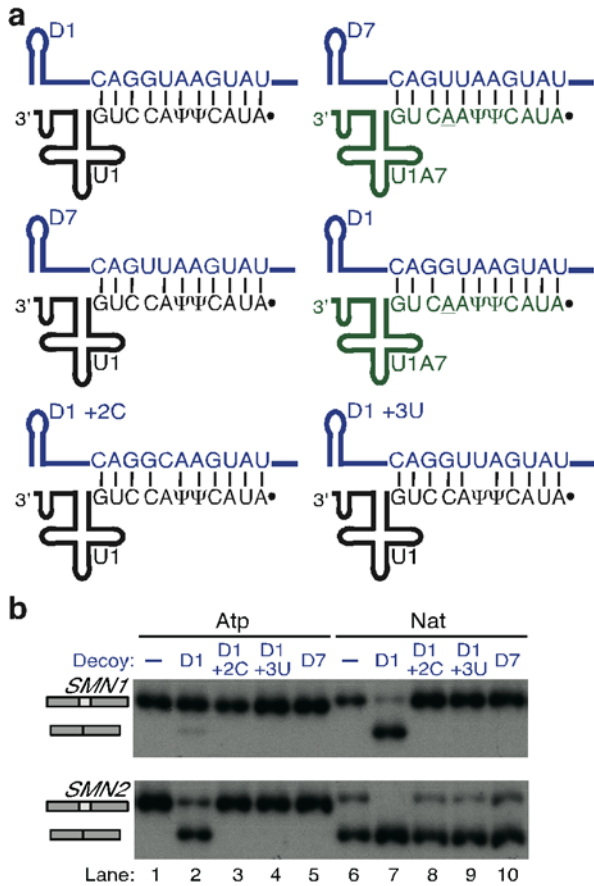
**Supplementary Figure 4 Analysis of double mutations at atypical 5' ss.** **a**, Schematic of the +3C+5C and +4C+6C 5' ss mutations in the *SMN1/2* context. Mutant nucleotides are shown in blue. **b**, Base-pairing profiles of the two double-mutant 5' ss with the corresponding suppressor U1 snRNAs carrying compensatory mutations in the canonical (C) or shifted (S) registers. Compensatory U1 mutations are indicated and highlighted in blue. **c**, Analysis of the mutant 5' ss. The top labels indicate the minigene (*SMN1* or 2), the mutant 5' ss in exon 7, and the mock or suppressor U1 snRNA used. In this case, the suppressor U1 snRNAs rescued exon 7 inclusion only in the *SMN1* context (lanes 1-6). For +3C+5C (lanes 1-3), both suppressors rescued splicing, but the suppressor U1 (G5 G7) in the shifted register was more effective. For +4C+6C (lanes 4-6), only U1 (G4 G6) in the shifted register had an effect. In this experiment, aberrant splicing products were also seen upon use of suppressor U1s, for unknown reasons.



**Supplementary Figure 5 Use of U1 and U1A7 snRNA suppressors to rescue splicing via mutant atypical 5' ss. a**, Schematic of the U1 (black) and U1A7 (green) snRNAs and the mutant atypical (+4C) 5' ss in the *SMN1/2* context. U1 and U1A7 differ by one nucleotide at their 11-nucleotide 5' ends (underlined in U1A7), and have several additional nucleotide differences in the snRNA body. **b**, Base-pairing of the four different combinations of suppressor U1/U1A7: the U1 snRNA body with the U1 snRNA 5' end; the U1 body with the U1A7 5' end; the U1A7 body with the U1 5' end; and the U1A7 body with its 5' end. For each combination, the compensatory mutation was introduced in the canonical (G5, (C)) or shifted (G6, (S)) register, indicated in blue. **c**, Analysis of the

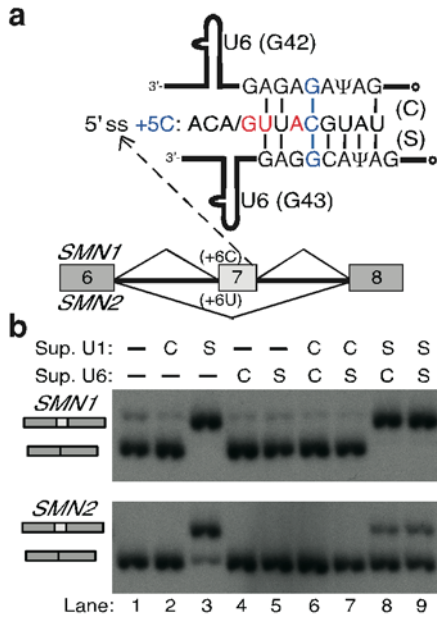
+4C mutant 5' ss with U1/U1A7 suppressors. The top label indicates the minigene (*SMN1* or 2), the suppressor body (U1 or U1A7), the 5' end (U1 or U1A7), and the compensatory mutation (control (-), C or S). None of the suppressors with the U1A7 body rescued exon 7 inclusion, but both 5' ends of U1 and U1A7 in the U1 body did.





**Supplementary Figure 6 Specificity of the U1/U1A7 snRNA decoys.** **a**, Base-pairing profiles of the different decoys used with the 5' end of U1/U1A7. The D1 and D7 decoys have perfect (11 bp) complementarity to their cognate snRNAs (top diagrams). In addition, the D1 decoy has one mismatch (10 bp) to the U1A7 snRNA in the canonical register, and the D7 decoy has one less potential base pair (10 bp) to U1 in the shifted register (middle diagrams). To assess the specificity of the D1 and D7 decoys, two additional decoys were constructed, D1+2C and D1+3U, which can have the same number of base pairs (10 bp) with U1 as the D7 decoy (bottom diagrams). **b**, Decoys need perfect complementarity (11 bp) to their cognate snRNA to have a strong effect. The top label indicates the 5' ss in *SMN1/2* exon 7 and the decoy used. As in b, expression of the D1 decoy resulted in exon 7 skipping. In contrast, the D1+2C, D1+3U and D7 decoys, which have one nucleotide

mismatch (10 bp) when base-pairing with U1, did not affect recognition of the natural or the atypical 5' ss. These results indicate that the snRNA decoys need perfect complementarity to their cognate snRNA, and explain why the D7 decoy does not significantly affect exon 7 inclusion.



**Supplementary Figure 7 U6 snRNA does not base-pair to atypical 5' ss in a shifted register.** **a**, Schematic of the suppressor U6 experiment. In this case, the suppressor U6 snRNAs carry only a point mutation to compensate for the +5C mutation in the atypical 5' ss. The U6 G42 and G43 mutations rescue base-pairing to the +5C 5' ss in the canonical (C) or shifted (S) register, respectively. Note that the G43 suppressor U6 has much stronger predicted base-pairing to the mutant 5' ss than the G42 suppressor. **b**, RT-PCR analysis of co-transfections of the *SMN1/2* constructs carrying +5C 5' ss and the various suppressor U1 and U6 snRNAs. The suppressor U1s (C or S) used are shown in **Supplementary Fig. 2d** online. As in **Fig. 3**, the suppressor U6s alone did not show any effect on exon 7 inclusion. Combined with suppressor U1, none of the suppressor U6 enhanced the levels of exon 7 inclusion.

## SUPPLEMENTARY TABLES

	<i>H. sapiens</i>	<i>M. musculus</i>	<i>D. melanogaster</i>	<i>C. elegans</i>	<i>A. thaliana</i>
Search					
NNHGTYRAGT <sup>i</sup>	40	35	10	11	55
NYGGTYRAGT	7	5	3	4	17
NYAGTRRAGT	4	5	0	3	5
NYAGTYYYAGT	1	2	1	1	7
NYAGTYRBGT	1	3	2	14	4
NYAGTYRAHT	1	1	2	23	9
NYAGTYRAGV	4	5	2	7	18
Homology <sup>ii</sup>	1	3			
Total	59	59	20	63	115
Human-mouse orthologous 5' ss					
Conserved	27	27			
Not conserved	13	8			
Intron not in ortholog <sup>iii</sup>	1	7			
No ortholog <sup>iv</sup>	10	0			
Not in databases <sup>v</sup>	8	17			
Total	59	59			

**Supplementary Table 1 Summary of the in silico searches for atypical 5' ss.** <sup>i</sup>See Methods for details about the SpliceRack<sup>1</sup> searches. <sup>ii</sup>Number of 5' ss that were added based on homology to the gene from the other mammalian species. <sup>iii</sup>The other mammalian counterpart does not have an intron mapped at the same location. <sup>iv</sup>Genes for which we were not able to find an ortholog in the other mammalian species. <sup>v</sup>Genes that we were unable to find in the ENSEMBL or UCSC genome resources, and that correspond to putative genes awaiting confirmation. See Supplementary Table 3 for a complete list of atypical 5' ss in all five species.

HUMAN		<i>Homo sapiens</i>							
Search		NNHGYRAGT	40						
Accession	Gene Name	Int. #	Int. Length	Sequence	Orthologs	Ortholog sequence			
H1	NM_006314	<i>CNKSR1</i>	15	170	AAGAAATAgtaaagctctgggtgctgctggggg	M32			
H2	NM_006252	<i>PRKAA2</i>	1	28900	AGTGAAAGgttgagtagcgcctccgaccctgtgctgg	NO, XM_131633	AGTGAAGAgtagtgcagcggcggcgggcgggcag		
H3	NM_024911	<i>GPR177</i>	2	34708	CCAAATCAgtaagtgtaactcctctcatcctctct	NO, NM_026582	CCAAATCAgtaagtagcccaactgctatgctcaactct		
H4	NM_015485	<i>RWDD3</i>	2	1827	AGTTAAATgttgagtagaatctgctatttctgctt	NO, NM_028456	ACATCAAGgtgctgaaactcttttgaatgaat		
H5	NM_006135	<i>CAPZA1</i>	4	4351	TCTGTACTgttaagtatctgactgctctatcatgag	M15			
H6	NM_012134	<i>LMOD1</i>	2	840	CTCAAGAAgtaagtagttgtggcatgagcaacagg	NO, NM_053106	TGAAGAAAgtaaaaggtggcagcgaggctaggggct		
H7	NM_025145	<i>C10orf79</i>	17	2434	AACACACAgttgagttcctatgactagcttctgctc	Intron not in mouse ortholog, <i>D19Erid652e</i>			
H8	NM_198795	<i>TDRD1</i>	25	3361	GTAAAAAgttaagtaagtaaaactgatgtttcc	NO, NM_031387	GAAGCTAAgtaagttcccgtttccctctgtgatgc		
H9	NM_005316	<i>GTF2H1</i>	8	117	AGAAAAACgttaagataaagtcagaggtgcagtaact	M19			
H10	NM_033150	<i>COL2A1</i>	22	371	GAGAAAGAgtaagtgaaatgggagctccatccatg	NO, NM_031163	GAGAGAGAgtaagtagcaagagctccagctcacagg		
H11	NM_020818	<i>KIAA1409</i>	26	9379	CAAAATTCgttgagtagctctctcatctccctgaa	M27			
H12	NM_173805	<i>FLJ38723</i>	6	479	CAGGCCAgttgagtagctctccctggggaagg	NO ORTHOLOG			
H13	NM_198830	<i>ACLY</i>	5	442	AAGAAAGAgttgagtagagaggacgtagggaggtgc	M5			
H14	NM_002265	<i>KPNB1</i>	19	1607	CGTACACgttgagtagatacaaccagttccctctg	M4			
H15	NM_014906	<i>PPM1E</i>	1	198971	TACAAATAgtaagtagctggcctgctgctgctggg	M26			
H16	NM_032853	<i>MUM1</i>	3	1300	GAGAAAAgttaagttgtgttttaataactgttt	YES, <i>Mum1</i>	AAGAAAAgttaagtagctgttttttttttttt		
H17	NM_022737	<i>LPPR2</i>	8	342	CTCGAAAgttaagttggcgcaggtaaagggggccgg	M21			
H18	NM_017814	<i>TMEM161A</i>	8	104	ATGTGCAgtaagtaggctgctgctgggggggga	YES, NM_145597	CTGCTACAgtagggctctccaggtagtagtga		
H19	NM_003419	<i>ZNF345</i>	1	1213	AAGAGGAAgttaagtagcaagagtggggattgttac	NO ORTHOLOG			
H20	NM_012344	<i>NTSR2</i>	2	1834	GGTTCTCAgttaagtagctctgtaacatcccccaag	M6			
H21	NM_004027	<i>INPP4A</i>	6	284	AGGGAACAgtaagtagtagtctgtagtctgtgggatt	M1			
H22	NM_004543	<i>NEB</i>	11	671	TTAGCACAgtaagtagggaaggatgctgctgtgtca	NO, XM_130232	TCAGCTCAgtaagccggatccggccagatgctctg		
H23	NM_198851	<i>C22orf34</i>	4	187	CTAAAATAgttgagtagcagagcagaggggtgggaa	NO ORTHOLOG			
H24	NM_016173	<i>HEMK1</i>	1	378	CCAAAACAgtaagtagtagtctcaagagcagaga	M24			
H25	NM_198196	<i>CD96</i>	11	12792	CAAAAAATgttaagtagtagtctgctgctgagttcc	M9			
H26	NM_145026	<i>SPATS1</i>	1	9492	GACCTACAgttgagttcaagaatccagctgtggag	NO, NM_027649	TTCCAAGgttagtaggaagcttccccaggggaagtag		
H27	NM_033481	<i>FBXO9</i>	3	2848	GAAGAAAAgttaagtagtagatattgtaacaaatta	M23			
H28	NM_020320	<i>RARSL</i>	2	5285	CAAAAAGAgtaagtagataactgtaacatctat	M16			
H29	NM_004956	<i>ETV1</i>	10	1540	GAACAGAAgttaagtagtccaccagaaattccgtgt	M7			
H30	NM_001735	<i>C5</i>	14	979	TGGAAGAgtaagtagtagtctgctgctgctgctgtt	M13			
H31	NM_024657	<i>MORC4</i>	9	402	AAGAAACAgttgagtagtagatagaggttctctac	M25			
H32	XM_378381	<i>DKFZp686F0839</i>	9	142	ATGGAGCAgttaagtagtagtagtagtagtagtagtag	NO ORTHOLOG			
H33	XM_290670	<i>LOC23117</i>	15	21700	AAATTAAGgttaagtagtagtagtagtagtagtagtag	NO ORTHOLOG			
H34	XM_377716	<i>LOC402057</i>	1	87	ATCCTGTTgttaagtagtagtagtagtagtagtagtag	NO ORTHOLOG			
H35	XM_497955	<i>NO NAME</i>	2	10117	AGAACTACgttaagtagtagtagtagtagtagtagtag	ENSEMBL, UCSC NOT FOUND IN ENSEMBL, UCSC			
H36	XM_379482	<i>LOC401320</i>	2	4835	CAAAAACAgtaagtagtagtagtagtagtagtagtag	ENSEMBL, UCSC			
H37	XM_059830	<i>LOC136242</i>	3	1225	TATTTACCgttagtagtagtagtagtagtagtagtag	NO, <i>1700016G05Rik</i> NOT FOUND IN ENSEMBL, UCSC	TATTTACCgttagtagtagtagtagtagtagtagtag		
H38	XM_374277	<i>NO NAME</i>	1	466	CCCCAAAgtaagtagtagtagtagtagtagtagtag	Maps different in ENSEMBL, UCSC			
H39	XM_088578	<i>LOC375748</i>	1	3764	CAGAATGAgttgagtagtagtagtagtagtagtagtag	ENSEMBL, UCSC			
H40	XM_088525	<i>FAM125B</i>	1	13498	GGGGAACAgtaagtagtagtagtagtagtagtagtag	M12			
Search		NYGGTYRAGT	7						
H41	NM_144569	<i>SPOCD1</i>	10	2355	ATCTGGTGgttaagtagtagtagtagtagtagtagtag	NO ORTHOLOG			
H42	NM_001225	<i>CASP4</i>	8	1775	TGAAAATGgttaagtagtagtagtagtagtagtagtag	NO, NM_007609	TGAGAACAgtaagtagtagtagtagtagtagtagtag		
H43	NM_017828	<i>COMMD4</i>	2	222	CCAAGATGgttagtagtagtagtagtagtagtagtag	M22			
H44	NM_014705	<i>DOCK4</i>	4	4972	TCTATGTGgttaagtagtagtagtagtagtagtagtag	NO, XM_358307	TCTATGTGgttaagtagtagtagtagtagtagtagtag		
H45	NM_022444	<i>SLC13A1</i>	11	2334	AGAAATTGgttagtagtagtagtagtagtagtagtag	M17			
H46	XM_496074	<i>LOC440287</i>	5	221	CCAAGATGgttagtagtagtagtagtagtagtagtag	NO ORTHOLOG			
H47	XM_496078	<i>LOC440292</i>	2	222	CCAAGATGgttagtagtagtagtagtagtagtagtag	NO ORTHOLOG			
Search		NYAGTRRAGT	4						
H48	NM_001851	<i>COL9A1</i>	19	2534	GACCTCCAgtaagtagtagtagtagtagtagtagtag	M41			
H49	XM_372840	<i>LOC391209</i>	1	160	TCCCTGCAgttagtagtagtagtagtagtagtagtag	ENSEMBL, UCSC NOT FOUND IN ENSEMBL, UCSC			
H50	XM_114158	<i>NO NAME</i>	9	94	AGGAGCCAgttagtagtagtagtagtagtagtagtag	ENSEMBL, UCSC			
H51	XM_379433	<i>NO NAME</i>	4	1689	ATACCTTAgtaagtagtagtagtagtagtagtagtag	NOT FOUND IN			

							ENSEMBL, UCSC
	Search	NYAGTYYAGT	1				
H52	NM_006256	<i>PKN2</i>	7	1280	CTTGTCCAgttcagtaaccagatcttttaaaatcatgt	M46	
	Search	NYAGTYRBTG	1				
H53	NM_139025	<i>ADAMTS13</i>	25	466	GCAGTCCAgttatgtcctctcctcctcctcagcgca	NO, NM_001001322	AGAGCATGgttaggtttcctctttcctgtggga
	Search	NYAGTYRAHT	1				
H54	XM_370939	LOC388221	5	3534	AATTTGTAgttaattgacctcagtcacctgactaca	NOT FOUND IN ENSEMBL, UCSC	
	Search	NYAGTYRAGV	4				
H55	NM_004476	<i>FOLH1</i>	8	882	TCTACACAgttaagagactatttaatttaactctt	M52	
H56	NM_153696	<i>PSMAL</i>	3	887	TCTACACAgttaagagactatttaatttaactctt	NO ORTHOLOG	
H57	NM_024721	<i>ZFHX4</i>	4	9276	GAGCAGCAgttgagatcagctcaggtaatggtccta	NO, NM_030708	GGCCTCAGgttaaatgtattctctcagaagctagcc
H58	XM_375358	<i>Klkb14</i>	5	1216	CATCTGCAgttaaggcattcctcccagagaaggct	YES, XM_134585	CGTCTGCAgttaatgcattcctctctagaggagct
Manually added human/mouse orthologs							
H59	NM_015199	<i>ANKRD28 (KIAA0379)</i>	21	5662	CTACAATGgttaagtatacaacaacaatgcatatcat	M39	
MOUSE <i>Mus musculus</i>							
	Search	NNHGTIRAGT	35				
Accession	Gene Name	Int. #	Int. Length	Sequence	Orthologs		
M1	NM_030266	<i>Inpp4a</i>	4	260	AGGGAACAgttaagtatcgtcagtcggcaacagtct	H21	
M2	NM_007495	<i>Astn1</i>	17	8557	CATCATTAgtaagtcagaatgtctctcctcctct	Intron not in human ortholog, NM_004319	
M3	NM_007622	<i>Cbx1</i>	5	1420	TTAGCCCTgttgatcaccgccctctcaactctgact	Intron not in human ortholog, NM_006807	
M4	NM_008379	<i>Kpnb1</i>	19	697	CGTACACCgttgatataagccagctcccaaatcc	H14	
M5	NM_134037	<i>Acly</i>	5	447	AAGAAAGAgttgagttctgagagggcaacgctgggca	H13	
M6	NM_008747	<i>Ntsr2</i>	2	1495	GGTTCTCAgttaagtcctcaccacaatccctgagg	H20	
M7	NM_007960	<i>Erv1</i>	9	1873	GAACAGAAgttaagtcatttataagattctataa	H29	
M8	NM_030225	<i>Dlst</i>	4	113	TGTGTGCAgttaagtacctgattcttctgaatggact	NO, NM_001933	TGTATGCAgttaagctctctcttctggaaatgaatt
M9	NM_032465	<i>Cd96</i>	10	1677	CAAGAACTgttgatattgtcatggctcctccattcat	H25	
M10	NM_009985	<i>Ctsw</i>	6	86	CAACAACAgttgagtcaccagcctgtggaggtcagt	NO, NM_001335	CAACAACAgttgagtcaccagcctgtggaggtcagt
M11	NM_028412	<i>Ciz1</i>	1	507	CAGGATTAgtaagtgccgttggctatctgcgaggg	NO, NM_012127	CCGACAGgtgtgtgtgtctgctggaggccgtgcc
M12	NM_175184	<i>2610528K11Rik</i>	1	12860	GGGGAACAgttaagtcagcccgcgggaagacgggc	H40	
M13	NM_010406	<i>Hc</i>	14	1361	TGCAAAAGgttaagtgaattccgtgtgacctgggg	H30	
M14	NM_011978	<i>Slc27a2</i>	5	2964	TACAAAGAgtaagtaccgcaagatgacaaatgtg	NO, NM_003645	TACAGAAAgtaagtacattgaaatgagagcatacgt
M15	NM_009797	<i>Capz1</i>	5	4924	TCTGTACTgttaagtatctcctgtccctcgagaaca	H5	
M16	NM_181406	<i>Rars1</i>	2	4784	AAAAAAGAgtaagtatcctgtagtatctcaaatg	H28	
M17	NM_019481	<i>Slc13a1</i>	11	2557	AGAAATTAgttgagttatctttgaaactgaattat	H45	
M18	NM_019715	<i>Kemf1</i>	3	3286	CCTAGAAAgttgagtaagcagtcacccagagcacagg	Intron not in human ortholog, NM_020122	
M19	NM_008186	<i>Gtf2h1</i>	8	103	AGGAAACAgttaagtacatgatgaagctgtcactga	H9	
M20	NM_133214	<i>BC017612</i>	2	20057	TGCATACAgttgagtcaccattagtctctctgtga	NO, <i>C11orf75</i>	GCGGGCAGgtggctgggagcgcgcccggcgcgcc
M21	NM_144935	<i>BC018242</i>	7	329	CTCGAAAgttaagtgtgctcagtaagggggccgg	H17	
M22	NM_025417	<i>Comm4</i>	1	227	CCAAAATTgttaagtgcactgtttactagggttggg	H43	
M23	NM_023605	<i>Fbxo9</i>	3	1684	GAAGAAAgttaagtctgaagtcctcaacaataa	H27	
M24	NM_133984	<i>Hemk1</i>	2	338	CCAAAACAgttaagttaagttagggagacaggag	H24	
M25	NM_029413	<i>Morc4</i>	12	498	AAGCAACAgttgagttacacgggctgtttctccat	H31	
M26	XM_484060	<i>Ppm1e</i>	1	109399	TACAATAgttaagtgcgggctcgcgcccgggtg	H15	
M27	XM_484173	<i>9030205A07Rik</i>	21	9472	CAAAATCCgttgagttatctctctctctctgtgac	H11	
M28	XM_488661	<i>NO NAME</i>	5	3011	CCAGCACTgttaagtggaaatcataatcttctgg	NOT FOUND IN ENSEMBL, UCSC	
M29	XM_358492	<i>NO NAME</i>	1	15613	GGCATTCAgttaagtaattctactgagtcacacagt	NOT FOUND IN ENSEMBL, UCSC	
M30	XM_485305	<i>NO NAME</i>	2	205	CCAAAACAgttaagttataagtcataaagtgggttt	NOT FOUND IN ENSEMBL, UCSC	
M31	XM_485339	<i>170008P02Rik</i>	4	1296	AAATCAATgttgagttatatacaataatattagaat	Intron not in human ortholog	

M32	XM_110525	<i>Cnkr1</i>	15	161	AAAAAATAgttaagtcttgaatgagatggattgggg	H1	
M33	XM_357710	<i>LOC384557</i>	11	941	TGTAAACAgttaagtactctcaatggaactctgttgc	Different in ENSEMBL NOT FOUND IN ENSEMBL, UCSC	
M34	XM_134048	<i>NO NAME</i>	4	14935	TAAAAACAgttaagtactgcatgtttccacatactt	NOT FOUND IN ENSEMBL, UCSC	
M35	XM_487089	<i>NO NAME</i>	8	241	AGTTTGGAgttgagtgaggagggcaagagcaggataga	NOT FOUND IN ENSEMBL, UCSC	
	Search	NYGGTYRAGT	5				
M36	NM_145930	<i>AW549877</i>	1	92	CCAACATGgttgagtgtctagtaaagttctaagtgt	Intron not in human ortholog	
M37	NM_153795	<i>BC032204</i>	9	85	ATCTTCCGgttaagtatagggcctgctgaggagtg	<i>NO, URP2_HUMAN</i>	ATTCTTCGgttgagttggggccagatgagcagccctg
M38	NM_010576	<i>Itga4</i>	16	2149	GAAAAATGgttgagtatttctgtatttaaacataa	<i>NO, NM_000885</i>	AAAAAACAgtaggaatatttctttatcaaatatt
M39	XM_127673	<i>Ankrd28</i>	21	1747	CTACAATGgttaagtatcagataaaacctgctttgag	<i>YES, ANKRD28</i> <i>(KIAA0379)</i>	CTACAATGgttaagtatacaaacacaatgcatatcat
M40	XM_131242	<i>4930579F01Rik</i>	7	1282	GTCACACGgttaagtacagttgtgtacttcagatg	<i>NO, C4orf17</i>	ACCAAAAGgttaagtacagtttttggtaactattgagt
	Search	NYAGTRRAGT	5				
M41	NM_007740	<i>Col9a1</i>	19	2944	GACCTCCAgttaagtattttttattgtttaccaag	H49	
M42	XM_483979	<i>NO NAME</i>	2	128	TATTTCTAgttgagtggtgtccattgggccccaaatc	NOT FOUND IN ENSEMBL, UCSC	
M43	XM_283061	<i>Hectd1</i>	7	886	TCAGGTCAgttaagtgcgtgtgtatcatttagtttc	Intron not in human ortholog	
M44	XM_484630	<i>NO NAME</i>	1	1470	TAGCACCAGtaaaagtctcttccatgcatgccact	NOT FOUND IN ENSEMBL, UCSC	
M45	XM_142197	<i>Rga1</i>	2	620	GTCTCTCAgttaagtggtaaaagcatcagaatgtcc	Intron not in human ortholog, NM_020769	
	Search	NYAGTYYAGT	2				
M46	NM_178654	<i>Pkn2</i>	6	873	CTTGTCAGttcagtaaccatactttaaagtcttgc	H53	
M47	XM_486735	<i>No NAME</i>	2	16518	TTAGGTTAgttttagtgatgtaaatccctgtccctac	NOT FOUND IN ENSEMBL, UCSC	
	Search	NYAGTYRBT	3				
M48	NM_010655	<i>Kpna2</i>	12	42411	AAAGACCAgttcgggtctacagatgagttccaggacag	Intron not in human ortholog, NM_002266	
M49	XM_139844	<i>EG240055</i>	1	111	TTCAGCCAgttgggtagcatgatcccaagccacacc	Does not map to the gene in ENSEMBL, UCSC	
M50	XM_489266	<i>NO NAME</i>	1	759	TTCCCTAgttatgtccacatcttataaaaaatcca	NOT FOUND IN ENSEMBL, UCSC	
	Search	NYAGTYRAHT	1				
M51	NM_177011	<i>9330154K18Rik</i>	6	639	TGCAAGCAgttgaaatggaaccctaagtctctctg	Intron not in ENSEMBL, gene not found in UCSC	
	Search	NYAGTYRAGV	5				
M52	NM_016770	<i>Folh1</i>	8	588	TCAACACAgttgagatattttctaagtattctgtca	H56	
M53	XM_488594	<i>LOC432562</i>	1	2429	CAAGGCCAgttcaagcactgacctgaaagtctttgtg	NOT FOUND IN ENSEMBL, UCSC	
M54	XM_358474	<i>NO NAME</i>	3	2725	TGTTGCCAgttcaagaaaaaacactgattccacaa	NOT FOUND IN ENSEMBL, UCSC	
M55	XM_488085	<i>NO NAME</i>	5	103	ATGGCATAgtcaaggtctgagaatgggaagcttatcatc	NOT FOUND IN ENSEMBL, UCSC	
M56	XM_358170	<i>NO NAME</i>	1	119	ACTTCACAgttgagccctgaacaacactgtggagcac	NOT FOUND IN ENSEMBL, UCSC	
	Manually added human/mouse orthologs		3				
M57	NM_023431	<i>Mum1</i>	3	976	AAGGAAAgtttaagtgtgctgttttttttttttt	H16	
M58	NM_145597	<i>Tmem161a</i>	8	87	CTGCTACAgttaggcgggtctcccaggtaggattgga	H18	
M59	XM_134585	<i>Klkb14</i>	5	857	CGTCTGCAGttaatgattcctctctagaggaggtct	H58	
	FRUITFLY	<i>Drosophila melanogaster</i>					
	Search	NNHGTYRAGT	10				
	Accession	Gene Name	Int. #	Int. Length	Sequence		
D1	NM_057884	<i>ppk</i>	2	185	GATGAGGAgttgatttccagcatatacttttag		
D2	NM_137551	<i>CG15117</i>	4	61	GATACAGAgttgatttccgacacctgaaaggatt		
D3	NM_139436	<i>CG15822</i>	1	1762	AAAACCAgttcagatttctacatgaacaacacgca		
D4	NM_168190	<i>NO NAME</i>	1	283	GCCAAAGCgtcaactgtggcgtgactcaaatgt		

D5	NM_140439	<i>CG5842-RA</i> ( <i>nan</i> )	4	60	CGAATACAgttaagtattgccttatattataataca
D6	NM_169384	<i>CG31386</i>	3	6050	AAAAACCAgttgagtaattatgtactgcaacaacaaa
D7	NM_142173	<i>CG4210-RA</i> <i>CG7985-RA</i> ( <i>CG7985</i> )	1	63	TGATTCAAgttgagttttataagattcctattttaata
D8	NM_142432	<i>CG1786-RA</i> ( <i>Cyp318a1</i> )	1	11970	CTTTAACAgttaagttcgaatggttgaiaaaccaat
D9	NM_167318	<i>CG3291-RA</i> ( <i>pcm</i> )	5	84	CAGGAACAgttgagtgtagagtcgaagccattagagg
D10	NM_078684				
	Search	NYGGTYRAGT	3		
D11	NM_166710	<i>CG32013</i> <i>CG4795-RA</i> , ( <i>Cpn</i> )	1	89	CGGCGGTGgtcagtcggcaccgactcaacggcggt
D12	NM_169454		2	67	AAACTACGgttgagatcctgtagccgaaagtctctg
D13	NM_169605	<i>CG3631-RB</i>	1	72	AACAACGgttaagtcatagaaaagtatttgcctca
	Search	NYAGTRRAGT	0		
	Search	NYAGTYYAGT	1		
D14	NM_167187	<i>CG17754-RC</i>	6	1279	GTGCTACAgtttagtccaaccaatccacttctggtt
	Search	NYAGTYRBGT	2		
D15	NM_166504	<i>CG11474-RB</i>	1	71	GGTGGTAgttgtgtcacttacctgacctgagc
D16	NM_130679	<i>CG3526-RB</i>	1	1038	TGCGTTTAgttgtgtttatgcttaaatgcgtcatt
	Search	NYAGTYRAHT	2		
D17	NM_168054	<i>CG32259-RB</i>	4	308	CGGCGCCAgttaattgcaactgctcgataacctac
D18	NM_170040	<i>CG31139-RA</i>	1	56	CTCAAAATAgttaactatagaaataataataataa
	Search	NYAGTYRAGV	2		
D19	NM_141447	<i>CG14609-RA</i>	6	58	GGAGAATAgttgagcagctctacaataaaaagtgc
D20	NM_133036	<i>CG15373-RA</i>	5	109	TCGCACCAgttgagcattgattataagaatagaaga

WORM *Caenorhabditis elegans*

Accession	Gene Name	Int. #	Int. Length	Sequence	
		11			
	Search	NNHGYRAGT			
C1	NM_071062	<i>Y46H3D.4</i>	4	43	CTTTAAAAgttgagttttgatcttgaagaattgccaat
C2	NM_071184	<i>srbc-22</i>	3	63	CAAAATCAgttgagttattatccaattgaacatttt
C3	NM_074995	<i>F47H4.2</i>	1	1995	TGAAAAATgttgagttatttcaagtcggaagtaaa
C4	NM_075290	<i>Y80D3A.8</i>	4	119	AGCGATGAgtaagtgattttccgattgagttgtaa
C5	NM_068590	<i>C01G5.9</i>	3	357	TCAAAAAATgttaagttttatctttaaattgttaga
C6	NM_069297	<i>T05A1.7</i>	2	150	GTGGATATgtcaagtagagaagaataatgaactact
C7	NM_070325	<i>F02H6.3</i>	3	66	AAAAAGATgttaagtgaagagtttcaattctagta
C8	NM_061415	<i>F46F5.12</i>	2	44	TGGCATCAgttgagtcattttcctgcaacaattaat
C9	NM_061925	<i>NO NAME</i> <i>W02B12.4</i> ( <i>esterase</i> )	4	61	TAGAAAAATgttgagtcgaaaattgaaattcgcgcca
C10	NM_064042		3	578	TTCAAGAAgttgagttttgaactactaccggtac
C11	NM_064547	<i>K09E4.4</i>	3	53	TCTAAACAgttgagtcctcagggcaatcaggttgata
	Search	NYGGTYRAGT	4		
C12	NM_070973	<i>Y50D4B.7</i>	5	159	CTGGGGTGttaagttttattgtaatttaaacgaac
C13	NM_071110	<i>Y46H3A.5</i>	2	61	TGCGCATGgtcaggtttctcttttcacggccacgta
C14	NM_072502	<i>F18E3.10</i>	1	131	GGTTCCCGgttgagttgatgcatatcatcgtgtgta
C15	NM_066589	<i>F59B2.3</i>	2	50	ATTGAACGgttgagttgttcatttaacataattttc
	Search	NYAGTRRAGT	3		
C16	NM_059201	<i>T05E8.1</i>	5	185	TCGCGTCAgtgaagttgagtaagatctactctcca



C17	NM_074054	<i>srx-42</i>	7	1497	ACAAAGCAgtgaagtgaaggaaaaaaagaatgaataaa
C18	NM_061650	<i>F45C12.9</i>	2	75	TCTTGCAGtggagttatgaactggaacacctgaatt
	Search	NYAGTYYAGT	1		
C19	NM_171206	<i>NO NAME</i>	1	45	GAATCTTAgttcagtcattacgattttcagatga
	Search	NYAGTYRBGT	14		
C20	NM_058800	<i>D1037.5</i> ( <i>phospholipase</i> )	4	540	AACTAACAgttgggttgatgttttaactactgag
C21	NM_059316	<i>B0207.8</i>	3	48	AAAATATAgttaggtgaaacaattccggacatgtca
C22	NM_060727	<i>T09E11.5</i>	1	45	GTGGATCAgttaagttttcagtgataacctccagcat
C23	NM_060728	<i>T09E11.4</i>	1	44	GTTGATCAgttaggtttttcagtgataacctccagcat
C24	NM_071183	<i>srbc-21</i>	2	50	CCAAAACAgttgcgtttatccgcattttctaaatt
C25	NM_072749	<i>Y51B11A.1</i>	1	104	ATCAGACAgttatgtttcttaagaagtccgacc
C26	NM_074793	<i>nhr-271</i>	5	78	CGGAAATAgttagtaattgtaattgtaattgaaa
C27	NM_075237	<i>Y69H2.2</i>	10	310	TGAAACTAgttgggttgaaaaaatcgataattatc
C28	NM_067539	<i>asm-3</i>	7	53	TTGGAATAgttgggtttttgaagattggaaaaaa
C29	NM_068159	<i>R08C7.6</i>	4	1215	GATTGATAgttaggtttttcggctctcatagttt
C30	NM_070351	<i>LLC1.2</i>	4	3370	AATTAATAgttgggtttttttggaatttgaggg
C31	NM_062809	<i>R12C12.7</i>	1	362	TTTTGTAgttgggtgatattcgagataaatcgaa
C32	NM_075907	<i>tbx-31</i>	2	223	TTTCACAgttatgtaataagggaacatcaatggtt
C33	NM_065413	<i>unc-79</i>	2	82	GCCAGTCAgttgggtattttccaagggaattgataa
	Search	NYAGTYRAHT	23		
C34	NM_058984	<i>Y119C1B.9</i>	1	397	AATAATCAgttaatttttagttctatcaattgaaat
C35	NM_059145	<i>C32E12.4</i>	21	53	TGGCTATAgttaaattgaaattgtagattatgcaca
C36	NM_059245	<i>F46F11.7</i>	1	61	TCGGTACAgttaattctgataaagaatgatacaac
C37	NM_071952	<i>K11D12.8</i>	3	97	TTGGATCAgttgaataaaactaaaataatttgac
C38	NM_074745	<i>NO NAME</i>	5	123	TTGCAGCAgttgaatatgacgtaagtgaagtggagc
C39	NM_075053	<i>srz-45</i>	3	86	TTCTGATAgttgattgataacaatgaagatccgat
C40	NM_075147	<i>C14A6.5</i>	7	944	TTGTCCCAgttgattttcgaatcttcaaccgtttt
C41	NM_067921	<i>cyp-25A5</i>	1	293	GACCACCAgttaattggctttagggaaattgaaatt
C42	NM_068589	<i>C01G5.7</i>	3	46	ATTAGTTAgttactttaactttaaaactaaaatgct
C43	NM_070559	<i>Y51H4A.2</i>	1	82	AGCACTCAgtcaatttagttagtttggtgtgtgtt
C44	NM_182271	<i>K02E7.4</i>	2	246	TATGAACAgttaattgataataatattgatgttat
C45	NM_061573	<i>Y51H7C.12</i>	2	1095	ATTGTTCAgttaattttaactgtgataactataaa
C46	NM_061922	<i>sri-47</i>	2	49	CCAAAATAgttaattgaaattgaaattcagtcaaa
C47	NM_062134	<i>F18A12.7</i>	6	48	TCTGGTCAgttaatttttaacaaaactgactgaaat
C48	NM_063632	<i>T19H5.1</i>	3	794	ATTCAATAgttcaattgaaatgcaatcacttgacggat
C49	NM_064069	<i>NO NAME</i>	1	748	CAGTAACAgttaatttaaatattttttggtccgatt
C50	NM_076485	<i>H28G03.3</i>	2	44	TAAAAACAgttaatttttactacttaaacatctaaa
C51	NM_076877	<i>F46C8.8</i>	4	57	AGGTATCAgttaaatgttcgtagtatatttccctgtt
C52	NM_077618	<i>W03G11.4</i>	4	797	TTCCCCCAgtcaattatcactctttttttttgttt
C53	NM_077707	<i>R09A8.2</i>	6	105	AAAGGCTAgttgatttaagtagactcacttaatgat
C54	NM_078201	<i>Y7A5A.9</i>	5	614	TTGTTCCAgttaactttccgaaaaactaatatgitt
C55	NM_064757	<i>ver-2</i>	1	258	AAACGTTAgttaaactgtttgaatttgatgaattt
C56	NM_067102	<i>Y66A7A.5</i>	5	933	AAAATCCAgttgaattttctaaattttttattttt
	Search	NYAGTYRAGV	7		
C57	NM_071343	<i>srg-61</i>	3	56	G TTCACCAgttaagcaaaaggtaaaagccgctactgt
C58	NM_072056	<i>C24G6.3</i>	13	74	CCTGAACAgttgagatttttttaatttttaataaaa
C59	NM_073280	<i>Y32F6A.1</i>	6	43	GTCTATCAgttgagatatttatagaagttaatagt
C60	NM_069993	<i>sru-17</i> <i>R09D1.9</i>	5	52	CAGGAATAgttcaagagagcaggagacaattatgtgat
C61	NM_063629	( <i>chitinase</i> ) <i>R09D1.11</i>	3	44	AGAATACAgttgagggtcgcttattttgtaataat
C62	NM_063631	( <i>chitinase</i> )	3	44	AGAATACAgttaagggtacacttatagttgtaataat
C63	NM_063632	<i>T19H5.1</i>	1	44	AGGAGACAgttgagatattttgagttagacaatgac
PLANT		<i>Arabidopsis thaliana</i>			

Search	NNHGTYRAGT	55	Int.	Int.	Sequence
Accession	Gene Name	#	#	Length	
A1	NM_099990	<i>AT1G01080 acetyl-CoA C-acyltransferase (AT1G04710)</i>	3	87	ATGGAACAgttaagtcgttatcataaaaaaatctt
A2	NM_100351		2	86	ATTGTCGCgttaagttctctctctttgatattgt
A3	NM_100846	<i>AT1G09740</i>	1	92	TGGTCCCAgttagctctgttactactaattttgtgat
A4	NM_101003	<i>AT1G11300</i>	8	91	AGAATTCAGttaagtattattccaacccctttcata
A5	NM_101933	<i>AT1G20810</i>	1	262	CACTGGCCgttaagtcctcctctctcttttagctate
A6	NM_101948	<i>AT1G20950</i>	16	99	CTATTACAgttgagttttctgtatttctatcaaaa
A7	NM_102114	<i>AT1G22670</i>	3	82	CAGTCACAgttgagtattgtgatagtttaaacatc
A8	NM_102414	<i>AT1G26520</i>	5	303	ATTCAACTgttaagtcctgaaactctctctctct
A9	NM_102778	<i>ATMRP13</i>	18	80	TGCTTCCAgttaagtggaatgaaggtttacacaactc
A10	NM_102779	<i>ATMRP12</i>	18	75	TGTTTCCAgttaagtaaaatgacagcctacacaaccc
A11	NM_103331	<i>AT1G36510</i>	3	492	ACCGTTAgttgagttctgcaaatataggattgggta
A12	NM_148569	<i>AT1G48635</i>	2	91	AAAACCCAgttaagttctgtctcataaaaactgatcc
A13	NM_180628	<i>SCD1</i>	1	506	TTCCTACTgttaagttttcacaacatctttgtaatt
A14	NM_202303	<i>AT1G54920</i>	6	99	TGTTATCAgttaagttctgtcggagttctcggccct
A15	NM_105876	<i>AT1G72175</i>	3	115	ATATCGCAgttaagttcgtgatcatctgtctctctg
A16	NM_105969	<i>AT1G73110</i>	1	112	AAAAGATTgttaagttctcttttcttttgaaaaa
A17	NM_106111	<i>AT1G74510</i>	1	133	AATACTCAGttgagtttctttttttttttgat
A18	NM_106305	<i>AT1G76550</i>	16	112	CTATCTCAGttgagtttctctctattataagccta
A19	NM_106594	<i>GA2</i>	13	95	TGAAAATTgttaagttatcaatccccacacaattata
A20	NM_126485	<i>AT2G04540</i>	3	89	AGAGAACAgttaagttctctcatcttagatcaaatg
A21	NM_127266	<i>AT2G17120</i>	2	106	AATTGGACgttaagtgattgataactctgtgtgttg
A22	NM_179653	<i>HPT1</i>	3	93	TTGGCACAgttaagttctcttttaaaatgtaactctt
A23	NM_127595	<i>AT2G20360</i>	1	130	TAATCACCgttaagttctcctcattgatttttctta
A24	NM_129047	<i>AT2G34940</i>	8	79	CTTGTTCAgttaagttactgctctttttctccactt
A25	NM_129519	<i>AT2G39630</i>	7	116	TTGCTACAgttgagttctgcaactctgtttctagaca
A26	NM_129584	<i>AT2G40260</i>	1	125	CCCAGATAgttgagttcaattatataatgaatcaat
A27	NM_111905	<i>AT3G10690</i>	1	95	AGAATTCTgtttgagtgactttcttctcctcaaaagta
A28	NM_112698	<i>AT3G18160</i>	1	92	AAAACCCAgttgagttttgtaacctctctctttatt
A29	NM_112821	<i>AT3G19340</i>	11	88	TACTTCACgttaagttcttttttttcccttcttg
A30	NM_113020	<i>ATMRP6</i>	3	121	CTCACCAAgttgagttcctcagaagttgatcaaat
A31	NM_113689	<i>MER3</i>	21	125	AGAATACAgttgagttcctcctcaataactcctctggtt
A32	NM_148769	<i>AT3G29786</i>	2	112	AGACGACTgttaagttttgtaattagaagtggaagcatt
A33	NM_114085	<i>AT3G42180</i>	2	2411	TTTCATCAgttaagtttatacttttgattacattttct
A34	NM_114662	<i>ATGLDH</i>	2	283	AATCACAAgttaagttatcctcaacttctctatattag
A35	NM_202688	<i>SPP2</i>	3	92	AACGTGGCgttaagtttactgttatctttttgtgg
A36	NM_115099	<i>AT3G52390</i>	3	88	GAATCATAgttgagttcctctattgattctttgt
A37	NM_116013	<i>AT3G61480</i>	6	89	AACTTCCAgttgagttatcatatccgttatccctt
A38	NM_116154	<i>AT3G62890</i>	1	118	GGATGGAAgtcaagttgatcagacacaaatggaggca
A39	NM_116199	<i>AT3G63340</i>	2	84	CAAATGCAgttgagttccttttattttctctt
A40	NM_116709	<i>CPK22</i>	10	163	AAACTGAAgttaagttatctcagcatcttgctctac
A41	NM_118234	<i>AT4G21150</i>	16	123	TAGTTCTAgttgagttcctcaactcattaaaattg
A42	NM_119061	<i>AT4G29170</i>	7	142	GGCAAAGAgttgagttctcataatcgttactttgt
A43	NM_202922	<i>AT4G31240</i>	1	444	TCAGGTATgttaagttttttattacttttcaaca
A44	NM_119844	<i>RCE1</i>	3	95	TATTACCAgttaagttagttcttttttcccaactgaaac
A45	NM_121001	<i>SNG2</i>	4	99	TTAGGAGTgttaagttgtttttgtctttgttacagtt
A46	NM_121132	<i>AT5G10940</i>	9	501	TATTTGCAgttgagttgttctgctgctaaagcttaaat
A47	NM_121713	<i>AT5G17070</i>	7	101	CGAAAACAgttgagtttctgattgatttcttctt
A48	NM_147869	<i>AT5G20165</i>	5	93	GGCATATAgttgagttctatctcctcacttttactgc
A49	NM_122394	<i>AT5G24850</i>	2	87	CTCGTACAgttgagttctttattgctgtgtgaag
A50	NM_122718	<i>AT5G28350</i>	6	90	AACTTCCAgttgagttatcatatccgttatccctt
A51	NM_180567	<i>NO NAME</i>	2	76	AGTTTGCTgttaagttttcttcagccctttatata
A52	NM_123979	<i>APE2</i>	8	100	CCATCATCgttaagttcactttcatttctcactagt
A53	NM_125498	<i>AT5G61050</i>	2	95	AGAATCACgttaagtttctctcctccttccagatt
A54	NM_125499	<i>AT5G61060</i>	10	105	CGTATCACgttaagttaccactctctcagcttcc
A55	NM_125658	<i>AT5G62630</i>	4	84	TATGCCAgttaagttgacgaatcgaatacaaaaagcc

Search	NYGGTYRAGT	17	Int.	Int.	Sequence
Accession	Gene Name	#	#	Length	
A56	NM_101086	<i>AT1G12140</i>	4	85	CCTGGATGgttaagttctctcaatcttactttgg

A57	NM_148464	<i>AT1G14205</i>	4	119	TTTAGCCGgttaagtaacatttagttggagttgtttt
A58	NM_104157	<i>AT1G52790</i>	2	83	CATTGATGgttaagttaaacataatacaaaaataca
A59	NM_104543	<i>AT1G56670</i>	4	82	CCAATACGgttaagtagagtactattaccctgacca
A60	NM_104587	<i>AT1G58050</i>	5	86	TGCTCTTGgttaagtttaacccttttttgtaact
A61	NM_104588	<i>AT1G58060</i>	5	84	TGATCTTGgttaagtttaactcttagtaagaatcagt
A62	NM_179988	<i>AT2G39670</i>	1	71	CAACCTCGgtcaagtaaatgatctgattttgcatata
A63	NM_180001	<i>AT2G40640</i>	7	73	AATGGATGgttaagctctgtgctctatctccata
A64	NM_180313	<i>AT3G28270</i>	1	543	AATTTCTGgttaagttttctctctcaactctcttt
A65	NM_115438	<i>SBPASE</i>	2	92	AAAGCTTGgttgagtaacacaagaactcaactttatt
A66	NM_115577	<i>AT3G57170</i>	4	87	ATGAAATGgttgagtagtttccaaaatgctctctttg
A67	NM_117755	<i>AT4G16560</i>	12	255	TGCTCCTGgttaagctctctatccatttaactttgat
A68	NM_119341	<i>AT4G31900</i>	3	210	GCTTGCCGgttaagtgattcaagctgtctctttgtt
A69	NM_120931	<i>AT5G08460</i>	4	121	TAATTACGgttaagtagtctcaactatccctgttttac
A70	NM_121236	<i>AT5G11980</i>	7	92	AGTGCATGgttgagttctctgcatctgtgtattcct
A71	NM_124455	<i>AT5G50770</i>	3	85	CGGCTATGgtcaagttactccgaccaagaatctctcta
A72	NM_125539	<i>MIM</i>	22	85	AACTCATGgttgagctgcagatctgttatgttaaag
Search		NYAGTRRAGT	5		
A73	NM_100395	<i>AT1G05170</i>	3	88	GCATTACAgtaaagtttcaactctttgtctgtcattt
A74	NM_101251	<i>AT1G13830</i>	2	86	ACCTTCCAgtaaagtaataactcctaaaagctaaaa
A75	NM_101918	<i>AT1G20680</i>	5	1458	AACATCCAgtagagttacaagaacaattgcgagcat
A76	NM_114647	<i>ATATH7</i>	16	92	ACTCACCAgtagagctcttaagcatctatctcttt
A77	NM_124406	<i>AT5G50270</i>	1	148	ATCATCCAgtgaagttcctgagtttattactcagtc
Search		NYAGTYAGT	7		
A78	NM_100877	<i>AT1G10030</i>	1	1168	GACCACCAgttcagttctcccttttaetctatttt
A79	NM_103923	<i>AT1G50400</i>	1	285	AGCTACCAgttcagtataccctctctctactgctt
A80	NM_129435	<i>AT2G38780</i>	1	209	TAAACCCAgttcagtaaaaaatcatctttgacacaaa
A81	NM_129497	<i>AT2G39400</i>	2	98	CATGAACAgttcagtcagatattaaaaattactaac
A82	NM_112890	<i>TOM40</i>	1	498	AGCTCTCAgttcagtatacctaaaccttttccgat
A83	NM_118130	<i>AT4G20110</i>	7	92	CTTGCTCAgtttagtgattctctgtcataagccaaac
A84	NM_119538	<i>AT4G33800</i>	3	115	AATAACCAgttcagtaattttttattattctgtga
Search		NYAGTYRBGT	4		
A85	NM_180015	<i>AT2G41160</i>	1	87	CGGTTTCAgttacgtatcgaaattgaataaccatttc
A86	NM_180271	<i>AT3G17265</i>	1	86	CAAGAACAgttatgttttaggatatggaacacaacaa
A87	NM_118752	<i>AT4G26190</i>	3	71	TCTTCTTAgttagtgctctatacatagaaagttta
A88	NM_123512	<i>GLAI</i>	5	95	GTGCCCCAgttaggtggttttatgttttgcctatct
Search		NYAGTYRAHT	9		
A89	NM_100544	<i>NIH</i>	11	105	CTTATCCAgttaattcctctcttttatgtctctct
A90	NM_104690	<i>AT1G59980</i>	7	91	TTAATTCAgttaattatttctctacctcttaageta
A91	NM_105399	<i>AT1G67310</i>	1	141	AGTCCATAgtcaaaattatacctctttttgtctcta
A92	NM_106238	<i>AT1G75880</i>	2	87	TAGCGGTAgttaatttgagtaactatgcataatttaa
A93	NM_179823	<i>HVT1</i>	11	84	CTTATCCAgttaattctctttctctctatactcg
A94	NM_179859	<i>ATCSLB03</i>	5	94	TACAATCAgttaattctctctctcttattttttt
A95	NM_114201	<i>AT3G43330</i>	7	61	AAGAATCAgttaactacatgtttctactagctttgc
A96	NM_119567	<i>AT4G34060</i>	2	41	TGTAATCAgttaactcattttatttttaactttt
A97	NM_122884	<i>AT5G34930</i>	1	117	CAGAGACAgttgaatccatcatcatcatcatcagtc
Search		NYAGTYRAGV	18		
A98	NM_100923	<i>AT1G10490</i>	20	94	AAATACCAgttgagaatccctcttccctctctctcg
A99	NM_127387	<i>AT2G18290</i>	1	101	TATTGGCAgttgagaaacgaaaaaaaactaaagctc
A100	NM_127544	<i>ATHXK2</i>	3	105	ACGATACAgttgagctgttctctctctctctctcc
A101	NM_128895	<i>AT2G33350</i>	1	239	TTAACATAgttaagccagacaaaattacttaagtata
A102	NM_179985	<i>AT2G39410</i>	3	194	CATGAATAgttgagatcaaccacttcacaaattctat

A103	NM_129500	<i>AT2G39420</i>	3	97	CATGAACAgttgagatactcagccaactcatgatt
A104	NM_113609	<i>AT3G26950</i>	6	92	CAGTCTCAgtaagatcgttttttagtttcctcaacac
A105	NM_117800	<i>AT4G16970</i>	11	89	CCTGAACAgtaagatcgaactgaaatgagtagacacat
A106	NM_119057	<i>ATHXK1</i>	3	86	AAGAAGCAgtaagctcgatttctcaactattca
A107	NM_119554	<i>AT4G33940</i>	3	210	TCTTTGCAgtaagaacatctaccctgaaattacat
A108	NM_120026	<i>AT4G38650</i>	3	87	CTTCTCCAgttaagcatttcattctctctagaa
A109	NM_120494	<i>AT5G04120</i>	3	108	CTGTTGCAgtaagatccaatcccttaactccctgt
A110	NM_120572	<i>AT5G04900</i>	12	114	TATTCTCAgttaagacgattaatagctcgtcgtagaag
A111	NM_122051	<i>AT5G20440</i>	3	93	CATGAACAgtaagacacacattatcatccatagtt
A112	NM_122301	<i>AT5G23960</i>	5	78	AGAAATCAgttaagacaaaagaatcaacttaagcat
A113	NM_148058	<i>AT5G38386</i>	1	105	AGAAACTAgtaagctcattatctgattgatattc
A114	NM_125499	<i>AT5G61060</i>	8	217	CTACTCCAgttgagaccctctttttgtaccattt
A115	NM_125736	<i>AT5G63410</i>	5	192	ATAACACAgttaagctcgtttctttaatgtctgtttc

**Supplementary Table 2 Complete list of predicted atypical 5' ss in five species.** Each atypical 5' ss is labeled with the initial of the species where it was found, and a number. For each 5' ss, we provide the accession number, the gene name, the intron number and length, and a sequence that includes the atypical 5' ss (8 nucleotides in the exon, 30 in the intron). The query sequence used to obtain each hit is also shown. For the human and mouse collections, the orthologous 5' ss in the other species is indicated, whether it is conserved or not. Cases for which the orthologous intron or gene could not be found are also indicated.