

## Supplementary Information

Table S1. RNA-Seq data for human ultra-short introns.

Studies	BodyMap 2.0 project	Individual mRNA-Seq studies <sup>b</sup>			
Study ID	ERP000546	SRP002274	SRP010483	SRP012016	ERP000177
Tissues	16 tissues <sup>a</sup>	Brain	Pancreatic islet <sup>b</sup>	Stomach <sup>b</sup>	White blood cell
No. reads	–	46824452	204819206	426101413	247858317
Introns (SN)	Splicing observed tissues	Number of reads across ultra-short intron			
<b>1</b>	–	–	–	–	–
<b>2</b>	–	–	–	–	–
<b>3</b>	–	–	1	–	–
<b>4</b>	–	–	–	–	–
<b>5</b>	–	–	4	–	1
<b>6</b>	–	–	–	–	–
<b>7</b>	–	–	–	–	–
<b>8</b>	–	–	–	–	–
<b>9</b>	–	–	–	–	–
<b>10</b>	All 16 tissues	116	1826	12	1022
<b>11</b>	–	–	–	–	–
<b>12</b>	L	–	–	–	–
<b>13</b>	–	–	–	–	–
<b>14</b>	R, C, U, O, T, W	–	1	–	–
<b>15</b>	–	5	36	–	4
<b>16</b>	–	–	–	–	–
<b>17</b>	A, P, B, R, C, H, K, U, O, P, Y, T	26	28	–	2
<b>18</b>	W	–	2	–	3
<b>19</b>	A, R, K, M, P, Y, T	1	22	1	4
<b>20</b>	B, R	125	292	34	–
<b>21</b>	A, P, C, U, M, O, P, Y, W	5	1	–	9
<b>22</b>	B, H, O, P, T	43	–	–	1

See Table 1 for the introns and their serial numbers (SN; bold numbers are the confirmed introns).

<sup>a</sup> Abbreviations used; A: adrenal; P: adipose; B: brain; R: breast; C: colon; H: heart; K: kidney; L: liver; U: lung; M: lymph; O: ovary; P: prostate; S: skeletal muscle; T: testes; Y: thyroid; W: white blood cells; <sup>b</sup> “Pancreatic islet” and “stomach” are not included in the BodyMap 2.0 project.

**Table S2.** Results of RT–PCR experiments using cells and tissues.

<b>Introns (SN)</b>	<b>HIT ID number -intron no.</b>	<b>RT–PCR analysis <sup>a</sup></b>	<b>Tissues or cells tested <sup>b</sup></b>
<b>1</b>	HIT000059291-1	Expressed	4 Cultured cells, Pancreas
<b>2</b>	HIT000276161-4	Expressed	Placenta
<b>3</b>	HIT000008845-6	<b>Spliced</b>	4 Cultured cells
<b>4</b>	HIT000325704-2	No-Exp	
<b>5</b>	HIT000009363-12	<b>Spliced</b>	4 Cultured cells
<b>6</b>	HIT000084762-8	Expressed	Cerebellum
<b>7</b>	HIT000325704-3	Expressed	Testis
<b>8</b>	HIT000333308-1	Expressed	Lung
<b>9</b>	HIT000278575-1	Expressed	Brain, Lung, Placenta
<b>10</b>	HIT000192494-7	<b>Spliced</b>	4 Cultured cells, 8 Tissues
<b>11</b>	HIT000302202-1	Expressed	Testis
<b>12</b>	HIT000279220-1	<b>Spliced</b>	Placenta
<b>13</b>	HIT000333305-1	Expressed	Brain, Cerebellum
<b>14</b>	HIT000495960-1	<b>Spliced</b>	Peripheral blood leukocyte
<b>15</b>	HIT000191419-3	n/a	
<b>16</b>	HIT000091849-1	Expressed	Testis, Placenta
<b>17</b>	HIT000324311-10	<b>Spliced</b>	Brain
<b>18</b>	HIT000058074-1	No-PCR	
<b>19</b>	HIT000052133-11	No-PCR	
<b>20</b>	HIT000082518-3	<b>Spliced</b>	Cerebellum
<b>21</b>	HIT000252921-4	<b>Spliced</b>	Lung, Cerebellum
<b>22</b>	HIT000058190-7	<b>Spliced</b>	Brain

See Table 1 for the introns with their serial numbers (SN; bold numbers are the confirmed introns). <sup>a</sup> RT–PCR detection of the endogenous spliced product or transcript. See Table 1 (footnote ‘i’) for the definitions; <sup>b</sup> Eight tissues were tested: Cerebrum (total brain), cerebellum, hippocampus, medulla oblongata, lung, peripheral blood leukocyte, placenta, and testis. Four cells were tested: Human mammary epithelial cells (HMEC); MCF7, MDA-MB231, and HepG2 cell lines.

**Table S3.** Potential NMD targets following retention of the ultra-short introns.

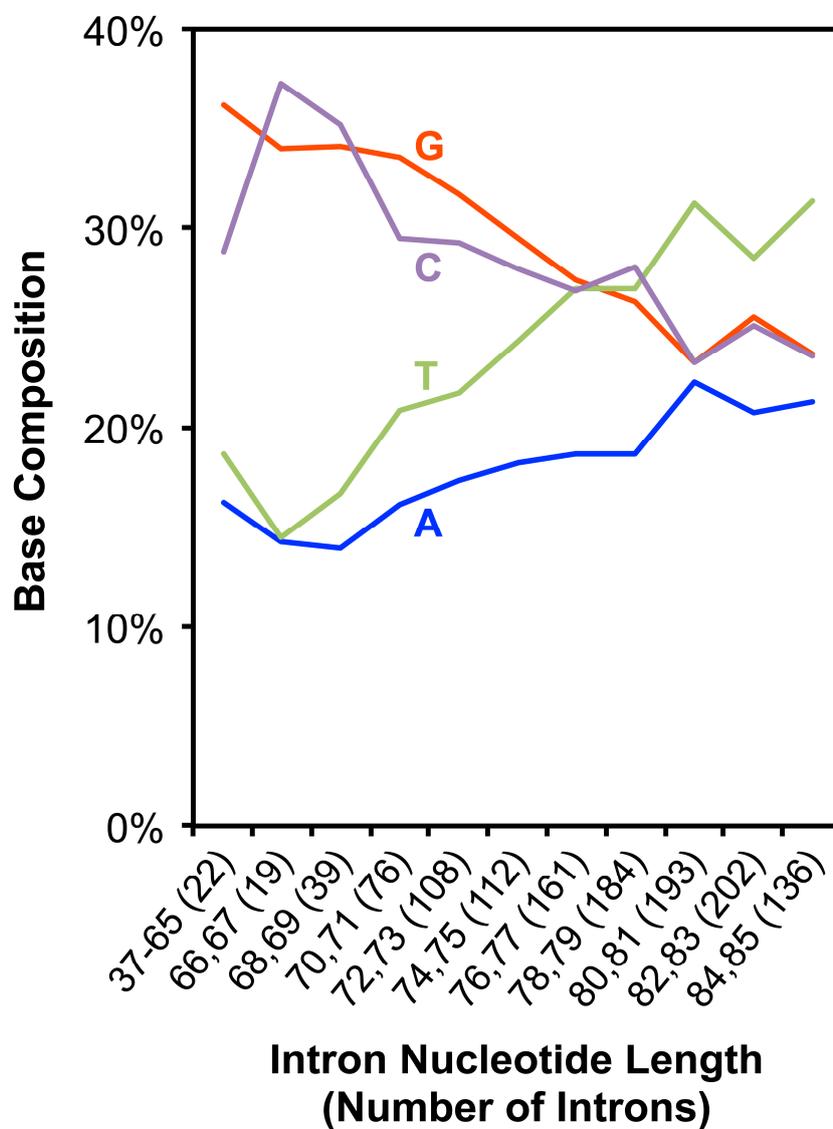
Introns SN <sup>a</sup>	Intron Length (nt) <sup>a</sup>	Intron position/ no. of introns <sup>a</sup>	Site of intron <sup>a</sup>	Effects of the retained intron			
				Frame- shift	Position of PTC generated		NMD target <sup>b</sup>
					From 3' end of intron (nt downstream)	From 3'-most exon- exon junction (nt upstream) <sup>b</sup>	
<b>1</b>	37	1/3	CDS	Yes	51	515	<b>Yes</b>
2	41	4/4	CDS	Yes	34 (in 3'-most exon)	-34 (downstream)	No
<b>3</b>	43	6/14	CDS	Yes	(within intron)	1424	<b>Yes</b>
<b>4</b>	47	2/15	CDS	Yes	184	1805	<b>Yes</b>
<b>5</b>	49	12/13	CDS	Yes	426 (close to 3'-end)	-303 (downstream)	No
6	50	8/10	CDS	Yes	(within intron)	311	<b>Yes</b>
7	54	3/15	CDS	No	(no PTC)	(no PTC)	No
8	54	1/2	CDS	No	(within intron)	478	<b>Yes</b>
9	55	1/5	CDS	Yes	93	332	<b>Yes</b>
<b>10</b>	56	7/13	CDS	Yes	(within intron)	589	<b>Yes</b>
11	61	1/13	5'UTR	-	-	-	-
<b>12</b>	62	1/7	CDS	Yes	(within intron)	1120	<b>Yes</b>
13	62	1/2	CDS	Yes	(within intron)	92	<b>Yes</b>
<b>14</b>	62	1/6	CDS	Yes	(within intron)	935	<b>Yes</b>
<b>15</b>	63	3/4	CDS	No	(no PTC)	(no PTC)	No
16	63	1/2	5'UTR	-	-	-	-
<b>17</b>	65	10/28	CDS	Yes	(within intron)	2838	<b>Yes</b>
<b>18</b>	65	1/20	CDS	Yes	(within intron)	3453	<b>Yes</b>
<b>19</b>	65	11/13	CDS	Yes	224	173	<b>Yes</b>
<b>20</b>	65	3/11	CDS	Yes	97	1030	<b>Yes</b>
<b>21</b>	65	4/4	CDS	Yes	(no PTC)	(no PTC)	No
<b>22</b>	65	7/26	CDS	Yes	55	1629	<b>Yes</b>

<sup>a</sup> See Table 1 for the explanations with their serial numbers (SN; bold numbers are the confirmed introns); <sup>b</sup> If this number is >50–55 nt, it fulfills the criterion for NMD target (see text).

**Table S4.** Primer DNA sequences used in RT–PCR.

Target introns (SN)	Primer names (Forward & Reverse)	Primer sequences (5' to 3')
2	ENSA-F	CAAAAGTACTTTGACTCAGGAGACTACA
	ENSA-R	GACTAGGTGCCCATTTTCATCCCTGCCTA
3	ESRP2-kozak-F	<b>gccaccatgg</b> GTC AATTGTTTTCGAAGCCCGAGGTG
	ESRP2-R1	GCCACGGCAACCCACGAGCCCGTACC
5	NDOR1-kozak3-F	<b>gccaccatgg</b> GCTGCCGCTGGCGGGACCAAGACTTC
	NDOR1-SalI-R	<b>ggtcgacg</b> AGGGCAGAGAACAGGGCTGGGGGC
6	SAMD14-F	GCCCAGCTCCTGGCTCATGCACAGGAGG
	SAMD14-R	TCTGAAGACTGAGACAGCGTGTGGTAGG
7	IFRD2-47nt-F	CCAAGTGCCTGAGAAGCCAGTCCTGCC
	IFRD2-47nt-R	ATCCACGAGGCTAAAATACTAATCCCA
	IFRD2-54nt-F	CTAGCCGCCGTCGCCTGCGCGCTCAGGC
	IFRD2-54nt-R	GGGGGGCGCGGGTCAGGGACCCGGTGG
8	HSP90B2P-F	TCCAGTCTTTCATCATCCAGCTGACAT
	HSP90B2P-R	CCTTCCTTGGCAACATTCTGGAACCTCT
9	AKIRIN2-F2	GCCACTCTGAAAAGGACTCTGGATTTTCG
	AKIRIN2-R2	GAAGGGGGATGGCTCCATTCGGAGATAC
10	hnRNPH1-kozak-F	<b>gccaccatgg</b> GCTATGGAGGCTATGATGATTACAATG
	hnRNPH1-SalI-R	<b>ggtcga</b> CAGTAGCTCTGTAAGGTAATCCCCGC
11	MSTO1-F	ATTTTGCCGGTTTCGTGGGCGCGCACTG
	MSTO1-R1	GAGTCGCGCGTGTAGGTCTCTGGCCG
	MSTO1-R2	CTGGGGAAGCAGAGAGGACGAGGCAGGC
12	SIGLEC6-F	CCTCCACTCCCTTCCCCTTCTCTGCTCA
	SIGLEC6-R	GGCAATCTGCAGGGTACGAGGACGCACA
13	HSP90AB4P-F	CTGATGACATCACCCAGGAAGAGTGTGG
	HSP90AB4P-R	ACGGACAAAGCCGACATACTCTGGCATC
14	SIGLECP3-F	CACACAGGAAGCCAGGAAGCCTCTGCC
	SIGLECP3-R	TTCCAGCCATACTCGAGGTAGAAGACCG
16	HIT000091849-F	GGAGAATAGAAATGCATTTTCTCTAAAC
	HIT000091849-R	CATTGGTGGTCACTCTGAGTTTGTAGTTC
17	PLXNA1-F	TACTCCTACGAGGGGAACGATGTCAGCG
	PLXNA1-R	GCCACGCACCATCCGCACTCGAAGCGCG
18	RECQL4-F	ATGGAGCGGCTGCGGGACGTGCGGGAGC
	RECQL4-R	CGCGGGTCTCCTCCGGCGCCGCTCCAC
19	C11orf35-F	AGAAGCCGCTGCGCGCTCCTCGAGCCG
	C11orf35-R	GAGCGGGAAGCGGTCGATGGACAAGTCG
20	ARHGDIG-F	GCGCGAGCTGGCTGAGGAGTTTGGTGTG
	ARHGDIG-R	GCGCCTGGGCGGCCGCTCGTCTCCAG
21	TNFRSF18-F	ACTGCACATCTGGCAGCTGAGGAGTCAGTG
	TNFRSF18-R	CACCCACAGGTCTCCCAGCCGCCCTTC
22	ADAM11-F	GGACCCCTTCCCACCTCATTTACCGGA
	ADAM11-R	CTGCCTTTTCTTCTCAGCCTCGGCCGG

See Table 1 for the introns and their corresponding serial numbers (SN; bold numbers are the confirmed introns). The inserted sequences (Kozak sequence and SalI site sequence) are indicated in lower-case letters.



**Figure S1.** The G-base frequency shows a gradual increase in human short introns from 85 nt toward  $\leq 65$  nt. The average base compositions of the screened conserved human introns (total of 1252 introns; see section 3.2) were calculated and plotted in the ultra-short range (37–65 nt) and in the short range (66–85 nt) for each 2-nucleotide length. The number of short introns in each group is shown in parentheses.