

Supplementary tables to the paper: Computational identification of functional introns: high positional conservation of introns that harbor RNA genes

Michal Chorev, Liran Carmel

Species	Build	Source
<i>Acyrtosiphon pisum</i>	assembly2	AphidBase [1]
<i>Aedes aegypti</i>	AaegL1	Ensembl [2]
<i>Ailuropoda melanoleuca</i>	ailMel1	Ensembl [2]
<i>Anolis carolinensis</i>	anoCar1	Ensembl [2]
<i>Anopheles gambiae</i>	anoGam1	VectorBase [3]
<i>Apis mellifera</i>	apiMel3	Ensembl [2]
<i>Aplysia californica</i>	aplCal1	Ensembl [2]
<i>Arabidopsis thaliana</i>	TAIR10	Ensembl [2]
<i>Aspergillus nidulans</i>	ASM14920v1	Ensembl [2]
<i>Bombyx mori</i>	v2.0	SilkDB [4]
<i>Bos Taurus</i>	bosTau4	Ensembl [2]
<i>Branchiostoma floridae</i>	braFlo1	JGI [5]
<i>Caenorhabditis brenneri</i>	caePb2	UCSC [6]
<i>Caenorhabditis briggsae</i>	cb3	UCSC [6]
<i>Caenorhabditis elegans</i>	ce6	UCSC [6]
<i>Caenorhabditis japonica</i>	caeJap1	UCSC [6]
<i>Caenorhabditis remanei</i>	caeRem3	UCSC [6]
<i>Callithrix jacchus</i>	calJac3	Ensembl [2]

<i>Canis familiaris</i>	canFam2	Ensembl [2]
<i>Cavia porcellus</i>	cavPor3	Ensembl [2]
<i>Choloepus hoffmanni</i>	choHof1	Ensembl [2]
<i>Ciona intestinalis</i>	ci2	UCSC [6]
<i>Ciona savignyi</i>	CSAV2.0	Ensembl [2]
<i>Danio rerio</i>	danRer6	Ensembl [2]
<i>Daphnia pulex</i>	Dappu1	Ensembl [2]
<i>Dasypus novemcinctus</i>	dasNov2	Ensembl [2]
<i>Dictyostelium discoideum</i>	dictybase.01	Ensembl [2]
<i>Dipodomys ordii</i>	dipOrd1	Ensembl [2]
<i>Drosophila ananassae</i>	droAna3	FlyBase [7]
<i>Drosophila erecta</i>	droEre2	FlyBase [7]
<i>Drosophila grimshawi</i>	droGri2	FlyBase [7]
<i>Drosophila melanogaster</i>	dm3	UCSC [6]
<i>Drosophila mojavensis</i>	droMoj3	FlyBase [7]
<i>Drosophila persimilis</i>	droPer1	FlyBase [7]
<i>Drosophila pseudoobscura</i>	dp4	UCSC [6]
<i>Drosophila sechellia</i>	droSec1	FlyBase [7]
<i>Drosophila simulans</i>	droSim1	FlyBase [7]
<i>Drosophila virilis</i>	droVir3	FlyBase [7]
<i>Drosophila yakuba</i>	droYak2	FlyBase [7]
<i>Echinops telfairi</i>	TENREC	Ensembl [2]
<i>Equus caballus</i>	equCab2	Ensembl [2]
<i>Erinaceus europaeus</i>	eriEur1	Ensembl [2]

<i>Felis catus</i>	felCat3	Ensembl [2]
<i>Fugu rubripes</i>	fr2	UCSC [6]
<i>Gallus gallus</i>	galGal3	Ensembl [2]
<i>Gasterosteus aculeatus</i>	gasAcu1	Ensembl [2]
<i>Giardia lamblia</i>	1.1	Refseq [8]
<i>Gorilla gorilla</i>	gorGor3	Ensembl [2]
<i>Homo sapiens</i>	hg19,GRCh37	Ensembl [2]
<i>Leishmania infantum</i>	1.1	Ensembl [2]
<i>Leishmania major</i>	ASM272v2	Ensembl [2]
<i>Linepithema humile</i>	1.2	Fourmidable [9]
<i>Loxodonta Africana</i>	loxAfr3	Ensembl [2]
<i>Macropus eugenii</i>	Meug_1.0	Ensembl [2]
<i>Microcebus murinus</i>	micMur1	Ensembl [2]
<i>Monodelphis domestica</i>	monDom5	Ensembl [2]
<i>Mus musculus</i>	mm9	Ensembl [2]
<i>Myotis lucifugus</i>	myoLuc1	Ensembl [2]
<i>Nematostella vectensis</i>	Nemve1	Ensembl [2]
<i>Neurospora crassa</i>	ASM18292v1	Ensembl [2]
<i>Ochotona princeps</i>	OchPri2.0	Ensembl [2]
<i>Ornithorhynchus anatinus</i>	ornAna1	Ensembl [2]
<i>Oryctolagus cuniculus</i>	oryCun2.0	UCSC [6]
<i>Oryza sativa</i>	MSU6	Ensembl [2]
<i>Oryzias latipes</i>	oryLat2	Ensembl [2]
<i>Otolemur garnettii</i>	otoGar1	Ensembl [2]

Pan troglodytes	panTro2	Ensembl [2]
Petromyzon marinus	petMar1	Ensembl [2]
Phaeodactylum tricornutum	ASM15095v1	Ensembl [2]
Physcomitrella patens	ASM242v1	Ensembl [2]
Phytophthora infestans	ASM14294v1	Ensembl [2]
Plasmodium falciparum	ASM276v1	Ensembl [2]
Pogonomymex barbatus	1.2	Fourmidable [9]
Pongo abelii	ponAbe2	Ensembl [2]
Pristionchus pacificus	priPac1	Ensembl [2]
Procapra capensis	proCap1	Ensembl [2]
Pteropus vampyrus	pteVam1	Ensembl [2]
Puccinia graminis	ASM14992v1	Ensembl [2]
Puccinia triticina	ASM15152v1	Ensembl [2]
Rattus norvegicus	rn4	UCSC [6]
Rhesus macaque	rheMac2	Ensembl [2]
Saccharomyces cerevisiae	sacCer2	Ensembl [2]
Schistosoma mansoni	sma_v3.1	Ensembl [2]
Schizosaccharomyces pombe	1.1,ASM294v1	Refseq [8], Ensembl [2]
Sorex araneus	sorAra1	Ensembl [2]
Spermophilus tridecemlineatus	speTri1	Ensembl [2]
Strongylocentrotus purpuratus	strPur2	Ensembl [2]
Sus scrofa	susScr2	Ensembl [2]
Taeniopygia guttata	taeGut1	Ensembl [2]
Tetraodon nigroviridis	tetNig2	Ensembl [2]

Tribolium castaneum	Tcas3.0	BeetleBase [10]
Trichoplax adhaerens	ASM15027v1	Ensembl [2]
Tupaia belangeri	tupBel1	Ensembl [2]
Tursiops truncatus	turTru1	Ensembl [2]
Ustilago maydis	UM1	Ensembl [2]
Vicugna pacos	vicPac1	Ensembl [2]
Vitis vinifera	IGGP_12x	Ensembl [2]
Xenopus tropicalis	xenTro2	Ensembl [2]

Table S1. List of annotated genomes taken into our initial data set.

Species	# orthologous groups	Fraction
Physcomitrella patens	79	0.18
Plasmodium falciparum	108	0.24
Saccharomyces cerevisiae	121	0.27
Ustilago maydis	123	0.27
Schizosaccharomyces pombe	133	0.30
Phaeodactylum tricornutum	139	0.31
Neurospora crassa	140	0.31
Leishmania major	143	0.32
Arabidopsis thaliana	163	0.36
Oryza sativa	165	0.37
Vitis vinifera	167	0.37
Ciona savignyi	171	0.38
Dictyostelium discoideum	174	0.39
Phytophthora infestans	191	0.42
Caenorhabditis elegans	206	0.46
Schistosoma mansoni	213	0.47
Drosophila melanogaster	256	0.57
Xenopus tropicalis	336	0.75
Danio rerio	346	0.77
Ornithorhynchus anatinus	352	0.78
Gallus gallus	363	0.81
Anolis carolinensis	378	0.84

Mus musculus	394	0.88
Gasterosteus aculeatus	402	0.89
Monodelphis domestica	403	0.90
Pan troglodytes	409	0.91
Pongo abelii	427	0.95
Homo sapiens	450	1.00

Table S2. A list of the final 28 species used in the analysis. For each species, we provide the number of orthologous groups in which it has a representative gene, and the fraction it makes out of the 450 groups that were analyzed.

	miRNA-bearing	miRNA-mixed	miRNA-lacking	total
snoRNA-bearing	3	0	96	99
snoRNA-mixed	0	1	22	23
snoRNA-lacking	50	10	3,981	4,041
total	53	11	4,099	4,163

Table S3. Overlaps between miRNA/snoRNA-bearing/mixed/lacking unique patterns.

Feature	Depicting	Description
log-likelihood (LOGLIKE)	conservation	Given EREM's estimation of the evolutionary model parameters, this is the log-likelihood of observing the pattern, $\log l_p$.
log number of times observed	typicality	$\log n_p$. Combined with $\log l_p$, this gives, up to a multiplicative constant, the (log) ratio of the number of times we expect to observe that pattern in our data, divided by the actual observed number, $\log l_p - \log n_p = \log \left(\frac{l_p}{n_p} \right)$.
binomial test	typicality	p-value of the binomial test (with

		Bonferroni correction), measuring how likely it is to see the pattern n_p times, given its likelihood is l_p .
number of ones	conservation, antiquity	The number of 1's in the pattern.
fraction of ones (ONES_RATIO_KNOWN)	conservation, antiquity	The number of 1's divided by the total number of 1's and 0's in the pattern.
number of times observed	typicality	The number of occurrences of pattern p , n_p .
number of evolutionary events (gain weight 3, loss weight 1; SANKOFF_G3L1)	conservation	The minimum number of intron gain and loss events required to obtain the pattern, given that gains cost three times as much as losses (using the Sankoff algorithm).
number of evolutionary events (gain weight 1, loss weight 3; SANKOFF_G1L3)	conservation	The minimum number of intron gain and loss events required to obtain the pattern, given that losses cost three times as much as gains (using the Sankoff algorithm).
number of evolutionary events (gain weight 1, loss weight 1)	conservation	The minimum number of intron gain and loss events required to obtain the pattern, given that losses cost as

		much as gains (using the Fitch algorithm).
one in amphibians (IN_AMPHIBIAN)	conservation, antiquity	This feature is 1 if the pattern has a 1 in at least one amphibian (<i>A. carolinensis</i> or <i>X. tropicalis</i>), otherwise it is 0.
one in fish (IN_FISH)	conservation, antiquity	This feature is 1 if the pattern has a 1 in at least one fish (<i>D. rerio</i> or <i>G. aculeatus</i>), otherwise it is 0.
one in birds (IN_BIRD)	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>G. gallus</i> otherwise it is 0.
one in fungi (IN_FUNGI)	conservation, antiquity	This feature is 1 if the pattern has a 1 in at least one fungi (<i>U. maydis</i> , <i>S. pombe</i> , <i>S. cerevisiae</i> , or <i>N. crassa</i>), otherwise it is 0.
one in plants (IN_PLANT)	conservation, antiquity	This feature is 1 if the pattern has a 1 in at least one plant (<i>V. vinifera</i> , <i>A. thaliana</i> , <i>P. patens</i> , or <i>O. sativa</i>), otherwise it is 0.
one in protists (IN_PROTIST)	conservation, antiquity	This feature is 1 if the pattern has a 1 in at least one protist (<i>P. falciparum</i> , <i>D. discoideum</i> , <i>L. major</i> , <i>P. tricornutum</i> , or <i>P. infestans</i>),

		otherwise it is 0.
one in <i>U. maydis</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>U. maydis</i> , otherwise it is 0.
one in <i>S. pombe</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>S. pombe</i> , otherwise it is 0.
one in <i>S. cerevisiae</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>S. cerevisiae</i> , otherwise it is 0.
one in <i>N. crassa</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>N. crassa</i> , otherwise it is 0.
one in <i>V. vinifera</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>V. vinifera</i> , otherwise it is 0.
one in <i>A. thaliana</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>A. thaliana</i> , otherwise it is 0.
one in <i>P. patens</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>P. patens</i> , otherwise it is 0.
one in <i>O. sativa</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>O. sativa</i> , otherwise it is 0.
one in <i>P. falciparum</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>P. falciparum</i> , otherwise it is 0.
one in <i>D. discoideum</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>D. discoideum</i> , otherwise it is 0.
one in <i>L. major</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>L. major</i> , otherwise it is 0.

one in <i>P. tricornutum</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>P. tricornutum</i> , otherwise it is 0.
one in <i>P. infestans</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>P. infestans</i> , otherwise it is 0.
one in <i>C. elegans</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>C. elegans</i> , otherwise it is 0.
one in <i>C. savignyi</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>C. savignyi</i> , otherwise it is 0.
one in <i>S. mansoni</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>S. mansoni</i> , otherwise it is 0.
one in <i>G. aculeatus</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>G. aculeatus</i> , otherwise it is 0.
one in <i>H. sapiens</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>H. sapiens</i> , otherwise it is 0.
one in <i>A. carolinensis</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>A. carolinensis</i> , otherwise it is 0.
one in <i>P. troglodytes</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>P. troglodytes</i> , otherwise it is 0.
one in <i>P. abelii</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>P. abelii</i> , otherwise it is 0.
one in <i>M. domestica</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>M. domestica</i> , otherwise it is 0.
one in <i>M. musculus</i>	conservation,	This feature is 1 if the pattern has a

	antiquity	1 in <i>M. musculus</i> , otherwise it is 0.
one in <i>D. melanogaster</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>D. melanogaster</i> , otherwise it is 0.
one in <i>D. rerio</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>D. rerio</i> , otherwise it is 0.
one in <i>X. tropicalis</i>	conservation, antiquity	This feature is 1 if the pattern has a 1 in <i>X. tropicalis</i> , otherwise it is 0.
intron density under LCA	conservation, antiquity	The last common ancestor (LCA) of all the intron-bearing species is assumed to be the species in which the intron was originated. LCA_LEAVES_RATIO is the ratio of intron-bearing (1's) to intron-lacking (0's) species, from among all the descendants of LCA.
taxonomic level of intron origin	antiquity	The number of ancestor nodes separating the LCA above from the tree root.
age of intron origin (LCA_AGE)	antiquity	The age of LCA above [MYA].
mean relative intron position	position	The mean distance of the exon-exon junction from the beginning of the

		coding sequence (CDS) divided by the CDS length.
median relative intron position (MED_REL_POSITION)	position	The median distance of the exon-exon junction from the beginning of the CDS divided by CDS length.
mean intron position	position	The mean distance of the exon-exon junction from the beginning of the CDS [nucleotides].
median intron position (MED_POSITION)	position	The median distance of the exon-exon junction from the beginning of the CDS [nucleotides].

Table S4. The initial set of 48 pattern-characterizing features. Abbreviated names are provided for the 13 features that were used in the final analysis. For each pattern, we note what aspect of the pattern it describes.

RNA gene	Mean			median		
	-bearing pattern	-lacking pattern	p-value (t-test)	-bearing pattern	-lacking pattern	p-value (U-test)
miRNA	690.99	400.37	$7.6 \cdot 10^{-4}$	454.9	184	$8.8 \cdot 10^{-9}$
snoRNA	840.18	392.52	$4.7 \cdot 10^{-16}$	744	184	$9.4 \cdot 10^{-22}$

Table S5. The results when LCA_AGE was calculated by EREM instead of by using the Dollo parsimony. Mean and median for both miRNA and snoRNA -bearing and -lacking unique patterns. P-values are Bonferroni-corrected.

References

1. Legeai, F., et al., *AphidBase: a centralized bioinformatic resource for annotation of the pea aphid genome*. *Insect Mol Biol*, 2010. **19 Suppl 2**: p. 5-12.
2. Flicek, P., et al., *Ensembl 2012*. *Nucleic Acids Res*, 2012. **40**(Database issue): p. D84-90.
3. Lawson, D., et al., *VectorBase: a data resource for invertebrate vector genomics*. *Nucleic Acids Res*, 2009. **37**(Database issue): p. D583-7.
4. Duan, J., et al., *SilkDB v2.0: a platform for silkworm (Bombyx mori) genome biology*. *Nucleic Acids Res*, 2010. **38**(Database issue): p. D453-6.
5. Grigoriev, I.V., et al., *The genome portal of the Department of Energy Joint Genome Institute*. *Nucleic Acids Res*, 2012. **40**(Database issue): p. D26-32.
6. Rhead, B., et al., *The UCSC Genome Browser database: update 2010*. *Nucleic Acids Res*, 2010. **38**(Database issue): p. D613-9.
7. McQuilton, P., S.E. St Pierre, and J. Thurmond, *FlyBase 101--the basics of navigating FlyBase*. *Nucleic Acids Res*, 2012. **40**(Database issue): p. D706-14.
8. Pruitt, K.D., T. Tatusova, and D.R. Maglott, *NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins*. *Nucleic Acids Res*, 2007. **35**(Database issue): p. D61-5.
9. Wurm, Y., et al., *Fourmidable: a database for ant genomics*. *BMC Genomics*, 2009. **10**: p. 5.
10. Kim, H.S., et al., *BeetleBase in 2010: revisions to provide comprehensive genomic information for Tribolium castaneum*. *Nucleic Acids Res*, 2010. **38**(Database issue): p. D437-42.