

Expansion and transformation of the minor spliceosomal system in the slime mold *Physarum polycephalum*

Graham E. Larue¹, Marek Eliáš², Scott W. Roy^{1,3}

1. Department of Molecular and Cell Biology, University of California, Merced, Merced, CA 95343, USA

2. Department of Biology and Ecology Faculty of Science, University of Ostrava, Ostrava, Czech Republic.

3. Department of Biology, San Francisco State University, San Francisco, CA 94132, USA
scottwroy@gmail.com

Abstract

While the vast majority of spliceosomal introns are removed by the so-called U2 (major) spliceosome, diverse eukaryotes also contain a mysterious second form, the U12 (minor) form, and associated introns. In all characterized species, U12-type introns are distinguished by several features, including being rare in the genome, containing extended evolutionary-conserved splicing sites, being generally ancient as judged by conservation between distant species, and being inefficiently spliced. Here, we report a remarkable exception in the slime mold *Physarum polycephalum*. The *P. polycephalum* genome contains > 20,000 U12-type introns—25 times more than any other species—with transformed splicing signals that have co-evolved with the spliceosome due to massive gain of efficiently spliced U12-type introns. These results reveal an unappreciated dynamism of minor spliceosomal introns and spliceosomal introns in general.

Introduction

Spliceosomal introns interrupt nuclear genes and are removed from RNA transcripts by machinery called spliceosomes. The vast majority of spliceosomal introns are removed by the so-called U2 (or “major”) spliceosome. However, diverse eukaryotes also contain a mysterious second form of spliceosomal intron (U12-type or “minor”), which are removed by a dedicated splicing machinery. In all characterized species, U12-type introns are distinguished by several features. First, U12-type introns are either rare or absent, ranging from 671 (0.36% of all introns) in humans (1) to 19 (0.05%) in fruitflies (2) to complete absence in diverse lineages (3). Second, U12-type introns show distinct extended splicing motifs at the 5' splice site ([G/A]TATCCTT) and branchpoint sequence (TTCCTT[G/A]AC, ≤ 45 bases from the 3' splice site) which exactly basepair with complementary stretches of core non-coding RNAs in the splicing machinery (4, 5). Third, U12-type introns are typically ancient (e.g., 94% of human U12-type introns are conserved as U12-type in chicken (6)), implying low rates of

U12-type intron creation through evolution (3, 6–8). Finally, U12-type introns show slow rates of splicing, suggesting inherently low efficiency of the U12 spliceosomal reaction (9–11).

Results

During manual annotation of GTPase genes in the genome of the slime mold *Physarum polycephalum*, we observed several introns lacking typical GT/C-AG boundaries, including both AT-AC and non-canonical introns (i.e., neither G[T/C]-AG nor AT-AC; Figure 1A). Most of these atypical introns also contained extended U12-like 5' splice site (5'SS) motifs ([G/A]TATC[C/T]TTT), consistent with previous evidence of U12 splicing in this species (12, 13). However, genome-wide analysis of the current *P. polycephalum* genome annotation revealed that all annotated introns have GY-AG boundaries, a pattern suggesting non-GY-AG introns may have been discarded by the annotation pipeline (14, 15). Indeed, an RNA-seq based genome reannotation using standard pipelines (Materials and Methods) allowing for non-GY-AG introns improved overall annotation quality (73.3% versus 60.1% BUSCO (16) sets present), and revealed a large number of previously unannotated introns, including a substantial number of introns with AT-AC splice boundaries (1,830 AT-AC, 54,816 GY-AG).

Our updated *P. polycephalum* annotation contains 3648 introns with perfect matches to the canonical U12 5'SS motif (3021 with GTATCCTT, 627 with ATATCCTT). In contrast, far fewer introns exhibit a classic U12-type branchpoint sequence (BPS) motif (561 with CCTT[G/A]AC present in the last 45 bases out of all introns, and only 20 of the 3,648 introns with perfect U12-type 5'SS motifs), and standard position weight matrix (PWM) methods (following (1, 12, 17, 18) failed to clearly identify U12-type introns (Materials and Methods and Figure S1). Lack of classic U12-type branchpoints were confirmed for a subset of conserved U12-type introns (those with U12-like 5'SS motifs found at positions that match those of U12-type introns in other species (Materials and Methods)). Instead, we noted the motif TTTGA falling within a short region near the 3'SS (terminal A 9-12 bp upstream of splice site), a feature also common in the manually identified non-GY-AG introns (Figure 1A). Genome-wide analysis of the 5'SS and TTTGA motifs showed a clear correspondence: TTTGA motifs are present 9-12 bp upstream of the 3' splice site (3'SS) in 59% (41/70) of conserved U12-type introns, as well as 42% (3,107/7,462) of GTATCYTT-AG introns and 67% (417/625) of ATATCYTT-AC introns, but only 6% (10,313/167,111) of other introns (Figure 1B). Consistent with a role in splicing, among introns with U12-like 5' splice sites, introns containing the TTTGA motif had lower average retention than those without it (Figure S2).

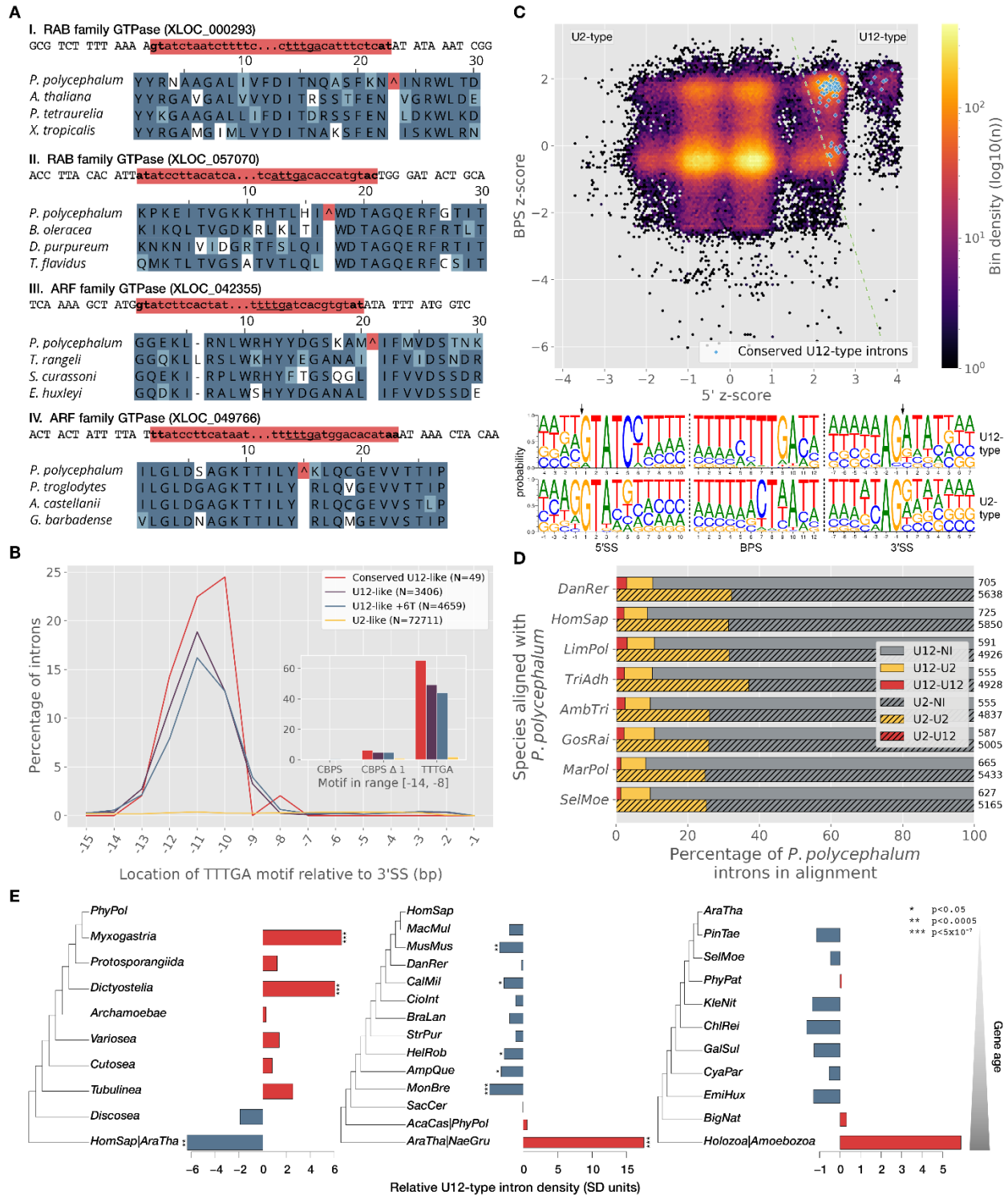


Figure 1. Evidence for massive gain of U12-type introns in *Physarum polycephalum*. **(A)** Canonical and non-canonical U12-like introns in conserved *P. polycephalum* GTPase genes. Δ indicates intron position. Lowercase red indicates intron sequence, with terminal dinucleotides in bold and putative BPS motifs underlined. **(B)** Presence of BPS motif in various groups of *P. polycephalum* introns. (Main) Occurrence of TTTGA motif as a function of number of nucleotides upstream of the 3'SS, for U12-like ([AG]TATCCTT-A[CG]) or [AG]TATCITT-A[CG] SS's for "U12-like" and "U12-like +6T", respectively), U2-like (GTNNG-AG), and conserved U12-like ([AG]TATC[CT]-NN and conserved as a U12-type intron in another species). (Inset) The same data as a cumulative bar plot for positions -14 through -8. **(C)** Conservation status of

P. polycephalum introns in other species, showing substantially lower U12- than U2-type conservation. For each species, the pair of bars shows the fractions of *P. polycephalum* introns of each intron type (U12-type, unhashed; U2-type, hashed) that are conserved as either U12-type (red) or U2-type (yellow) introns, or not conserved (gray). Total numbers of *P. polycephalum* introns assessed are given at right; for full species names, see Table S1. **(D)** Intron type classification and associated motifs. The main plot shows BPS-vs-5'SS PWM U12/U2 log-ratio z-scores for all *P. polycephalum* introns, with conserved U12-type introns highlighted in blue. The dashed green line indicates the approximate U2-U12 score boundary (Materials and Methods). Below the scatter plot are sequence logos showing motif differences between the two groups (20899 U12-type, top; 154299 U2-type, bottom). **(E)** Comparison of U12-type intron density (fraction of introns that are U12-type) in genes of different age categories for *P. polycephalum* (*PhyPol*), *Homo sapiens* (*HomSap*) and *Arabidopsis thaliana* (*AraTha*), relative to expectation (blue/red = below/above expectation) and scaled by standard deviation. U12-type intron densities in *P. polycephalum* are significantly overrepresented in newer genes, in contrast to the pattern seen in both human and *Arabidopsis*. Significance assessed via Fisher's exact test, with multiple-testing correction using the Holm step-down method. See Table S1 for species abbreviations.

Combining this position-specific atypical branchpoint motif with species-specific splice site motifs under a support-vector machine and PWM strategy (18) (Materials and Methods) led to a clearer separation of putative U12- and U2-type introns (Figures 1C, S3). Using a conservative criterion, we identified 20,899 putative U12-type introns in *P. polycephalum*, 25 times more than previously observed in any species. The true U12-type nature of these introns was further supported by two additional findings. First, comparisons of 8,267 pairs of *P. polycephalum* paralogs showed strong conservation of U12-type character: among intron positions shared between paralogs, an intron was 34-45 times more likely to be predicted to be U12-type if its paralogous intron was predicted to be U12-type (Materials and Methods, Table S2 and Figure S4). Second, putative *P. polycephalum* U12-type introns as a group are strongly biased away from phase 0 (26% compared with 39% for U2-type introns; phase is not a component of the scoring process), consistent with the phase bias observed in other species (Figure S5) (19, 20).

Interestingly, very few U12-type intron positions in conserved coding regions are shared with distantly-related species (e.g., only 9% of *P. polycephalum* U12-type introns found as either U2- or U12-type introns in humans, far fewer than for U2-type intron positions (31%); Figure 1D, see Table S1 for species abbreviations), indicating either massive U12-type intron gain in *P. polycephalum* or equally massive loss in other species. There is, however, no evidence for widespread loss of U12-type introns in other species, and previous results have attested to significant U12-type intron conservation across long evolutionary distances (6, 18, 21). Indeed, among U12-type introns conserved between *P. polycephalum* and plants and/or animals (i.e. ancestral U12-type introns), 63% are retained as either U2- or U12-type in the variosean amoeba *Protostelium aurantium*, and 70% are similarly retained in the discosean *Acanthamoeba castellanii* (Figure S6). That *P. polycephalum* has recently gained many U12-type introns is also supported by the fact that putatively recently-evolved *P. polycephalum* genes (i.e., those lacking homology to genes outside of closely-related species)

average expression of U12 spliceosomal components, relative to U2 spliceosomal components, is significantly higher in *P. polycephalum* than other species (Materials and Methods).

We also scrutinized components of the U12 spliceosome in *P. polycephalum*. A genomic search revealed a single candidate for the U12 snRNA that basepairs with the branchpoint (as previously reported in (13)). Strikingly, this sequence exhibits two transition mutations relative to the core branchpoint binding motif (underlined): GCAAAGAA, which produce basepairing potential with the putative ITTGA branchpoint with a bulged A, comparable to the canonical structure (Figure 2B). This apparent coevolution of core U12 spliceosomal machinery and branchpoint sequence represents to our knowledge the first true instance of coevolution of core intronic splicing motifs and core spliceosomal snRNAs.

U12-type introns in other species have been reported to have lower splicing efficiency than U2-type introns (9, 11, 22–24), raising the question of how *P. polycephalum* copes with ubiquitous U12-type introns. To investigate, we used RNA-seq and IRFinder (25) to calculate intron retention, and estimated splicing efficiency by comparing fractions of spliced and unspliced junction support between U2- and U12-type introns (Materials and Methods). Surprisingly, U12-type introns show slightly lower average intron retention (and higher average splicing efficiency) when compared either *en masse* (Figures 2C and S8) or in matched pairwise comparisons with neighboring U2-type introns in the same gene (Figures S9–11). These data suggest that minor spliceosomal kinetics are not inherently less efficient, and that they have been optimized in *P. polycephalum* in concert with the spread of U12-type introns. Consistent with increased efficiency of the U12 machinery in this lineage, we also found that the difference in average expression between the U12 and U2 spliceosomal components was smaller in *P. polycephalum* than is the case in species with lower U12-type intron densities (Figure 2D).

The near absence of U12-type intron creation in most lineages has been argued to reflect the low likelihood of random appearance of the strict U12-type splicing motifs at a given locus (20, 26). How, then, did *P. polycephalum* acquire so many U12-type introns? Inspired by cases of U2-type intron creation by insertion of DNA transposable elements (27, 28), we scrutinized U12 splice sites in *P. polycephalum*. We noted that many *P. polycephalum* U12-type introns carry sequences that resemble the signature of DNA transposable elements, namely inverted repeats (rtatctt...aagATAT) flanked by a direct repeat of a TA motif. This suggests the possibility that *P. polycephalum* U12-type introns could have been created by a novel DNA transposable element with TCTT-AAGA termini and TA insertion site (Figure 3). It is of note that *P. polycephalum* U12-type introns differ at two sites from the corresponding classic motif (TCCT-YAGA), where both changes increase the repeat character. An ancestral

decrease in the length and stringency of the branchpoint motif could have increased the probability of *de novo* evolution of a DNA transposable element carrying sufficiently U12-like splice sites for new insertions to be recognized by the U12 spliceosome.

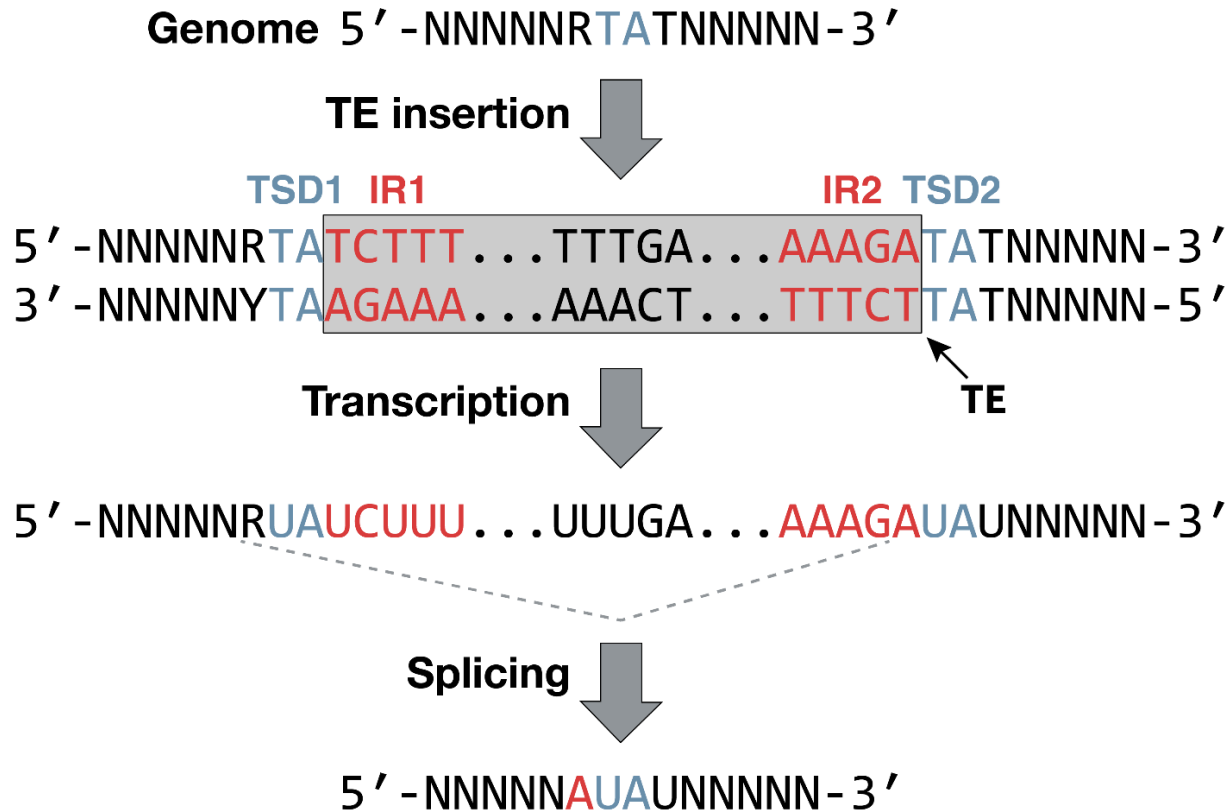


Figure 3. Proposed mechanism for transposon-driven creation of U12-type introns in *P. polycephalum*. Insertion of a transposable element (TE, gray box) carrying inverted repeats (IR1/IR2, red) leads to duplication of a TA target side (TSD1/TSD2, blue). Splicing at RT-AG boundaries leads to a spliced transcript with a sequence identical or nearly identical to the initial gene sequence with loss of an R (G/A) nucleotide and gain of the 3' A from the transposable element, maintaining the original reading frame.

Discussion

In contrast to the portrait of the U12 spliceosomal system as rare, ancient, static and suboptimal, these results expand our understanding of U12 diversity, by (i) increasing the upper bound of U12-type intron density per species by two orders of magnitude; (ii) showing that U12-type introns have been gained *en masse* through eukaryotic evolution; and (iii) showing that U12-type splicing is not necessarily less efficient than U2-type splicing. *P. polycephalum* provides promise for understanding of the flexibility of U12 splicing, a potentially important role given the increasing appreciation of U12 splicing errors in development and human disease. In addition, further study of *P. polycephalum* and related species may provide insights into the mechanisms and functional implications of the initial proliferation of introns in ancestral eukaryotes.

Acknowledgements

G.E.L. and S.W.R. were supported by the National Science Foundation (award no. 1616878 to S.W.R.).

M.E. was supported by the Czech Science Foundation project no. 18-18699S and the project “CePaViP” (CZ.02.1.01/0.0/0.0/16_019/0000759) provided by ERD Funds.

Materials and Methods

RNA-seq based reannotation of the Physarum polycephalum genome

We downloaded the *P. polycephalum* genome assembly and annotation from <http://www.physarum-blast.ovgu.de/>, and RNA-seq for *P. polycephalum* from NCBI's SRA database (accession numbers DRR047256, ERR089824-ERR089827, and ERR557103-ERR557120). To reannotate the genome, we combined *de novo* and reference-based approaches. First, we generated a *de novo* transcriptome from the aggregate RNA-seq data using Trinity (29). We also separately mapped the reads to the genome using HISAT2 (30), allowing for non-canonical splice sites (--pen-noncansplice 0), followed by StringTie (31) to incorporate the mapped reads with the existing annotations and generate additional putative transcript structures. Coding-sequence annotations for the assembled transcripts, informed by additional homology information from the SwissProt (32) protein database, were generated using TransDecoder (33), and further refined with the *de novo* transcriptome via PASA (34). In addition, an AUGUSTUS (15) annotation was generated from the mapped reads using BRAKER1 (35) explicitly allowing for AT-AC splice boundaries (--allow_hinted_splicesites=atac). Lastly, the AUGUSTUS- and StringTie-based gene predictions were merged using gffcompare (36), and updated again using PASA. To gauge the quality of our annotations versus those previously available, we performed a BUSCO (16) analysis against conserved eukaryotic genes; the previous annotations contained matches to 60.1% of eukaryotic BUSCO groups (54.5% single-copy; 27.1% fragmented; 12.8% missing); our annotation increased this percentage to 73.3% (64.4% single-copy; 18.5% fragmented; 8.2% missing).

Classification of intron types in P. polycephalum

All annotated intron sequences from our improved *P. polycephalum* genome annotation were collected and analyzed using a modified version of intronIC (18). Briefly, we first obtained high-confidence sets of U12- and U2-type *P. polycephalum* introns, as follows. High-confidence U2-type introns were defined as introns conserved as U2-type in at least three other species. Due to the low evolutionary conservation of putative *P. polycephalum* U12-type introns, the confident U12-type intron set was assembled by combining U12-type introns conserved as U12-type in one or more species, introns with perfect 5'SS motifs ([GA]TATCCTT) interrupting coding sequences in regions of good alignment to orthologs in one or more species, introns with near-perfect 5'SS motifs in addition to the TTTGA BPS motif 10-12 bp upstream of the 3'SS, and AT-AC introns (less likely to be false positives) with strong 5'SS consensus motifs in conserved eukaryotic genes (defined as representing a BUSCO match).

Sub-sequences of each intron corresponding to the 5'SS (from -3 to +8, where +1 is the first intronic base), 3'SS (from -5 to +4, where -1 is the last intronic base) and all 12mers within the branchpoint region (-45 to -5 where -1 is the last intronic base) were scored against position-weight matrices (PWMs) derived from the sets of high-confidence *P. polycephalum* U2- and U12-type introns to obtain U12/U2 log ratio scores for each motif. These log ratios were normalized to z-scores for each motif (5'SS and BPS), and were used to construct two-dimensional vector representations of each intron's score. In addition, to account for the narrow window of occurrence of the non-canonical TTTGA BPS, intronIC was modified to weight the branchpoint scores of introns whose BPS adenosines were found within the range [-12, -10] of the 3'SS, with the additional weight equal to the frequency of occurrence of the BPS adenosine at the same position within confident U12-type introns. Furthermore, unless explicitly stated otherwise, we used a more conservative U12-type probability score of 95% for classifying introns in *P. polycephalum*. The prominently-separated "cloud" in the upper-right of Figure 1D is composed mainly of AT-AC U12-type introns, whose 5'SS scores are more distinct than U12-type introns with other splice boundaries.

Identification of homologous sequences and conserved intron positions

Genomes and annotations for all additional species were downloaded from various online resources (Table S1), and in cases where sufficient RNA-seq was available and we suspected that U12-type introns had been systematically suppressed (e.g. zero or very few AT-AC introns annotated), we performed RNA-seq based annotation updates using Trinity and PASA (29, 33). For each genome, annotated coding sequences were extracted and translated via a custom Python script (<https://github.com/glarue/cdseq>). Annotated intron sequences were collected and scored using intronIC (18) with default settings. Under these settings, only introns defined by CDS features from the longest isoform of each gene were included, and introns with U12-type probability scores > 90% were classified as U12-type. Furthermore, introns shorter than 30 nt and/or introns with ambiguous ("N") characters within scored motif regions were excluded.

Between *P. polycephalum* and each other species (or, in the case of paralogs, itself), we performed pairwise reciprocal BLASTP v2.9.0+ (37, 38) searches (E-value cutoff of 10^{-10}), and parsed the results to retrieve reciprocal best-hit pairs (defined by bitscore) using a custom Python script (<https://github.com/glarue/reciprologs>). Pairs of homologous sequences were globally aligned at the protein level using ClustalW v2.1 (39), and introns occurring at the same position in regions of good local alignment ($\geq 4/10$ shared amino acid residues on both sides of the intron) were considered to be conserved (based on the approach in (40)).

Calculation of dS values between paralogs

We identified 8267 pairs of paralogs in *P. polycephalum* using the same approach as for other homologs. Each pair sharing at least one intron position was globally aligned at the protein level using Clustal Omega v1.2.4 (41), and then back-translated to the original nucleotide sequence using a custom Python script. Maximum likelihood dS values for each aligned sequence pair were computed using PAML v4.9e (42) (runmode = -2, seqtype = 1, model = 0), with dS values greater than 3 treated as equal to 3 in subsequent analyses (as dS values > 3 are not meaningfully differentiable in this context) (Figure S3).

Relative gene ages

For the three focal species (FS) *P. polycephalum*, human and *Arabidopsis thaliana*, sets of node-defining species (NDS) were selected to represent a range of evolutionary distances from the FS based on established phylogenetic relationships. In the case of *P. polycephalum*, we used data from Kang et al. (43) and their amoebozoan phylogeny; for the other two FS, we downloaded the NDS genomes and annotation files from the publicly-available resources Ensembl, JGI and NCBI (Table S1). We then performed one-way BLASTP (v2.9.0+, E-value cutoff 10^{-10}) searches of each FS transcriptome against the transcriptomes of its NDS set to establish an oldest node for each gene, defined as the ancestral node of the FS and the most-distantly-related NDS where one or more BLASTP hits to the gene were found. For example, a human gene would be assigned to the human-*Danio rerio* ancestral node if a BLASTP hit to the gene was found in *Danio rerio* (and optionally any more closely-related NDS) but not in any other more distantly-related NDS.

Once gene ages were assigned, for each FS we examined the difference of the observed and expected number of U12-type introns at each node using an expected value based on the aggregate density of U12-type introns in all other nodes, and scaled the observed-minus-expected value by dividing by the node's expected standard deviation (SD). For a given node, if n is the total number of introns in the node and p is the expected frequency of U12-type introns at the node based on combined frequencies from the other nodes, then per the binomial theorem $SD = \sqrt{np(1-p)}$. The significance of the observed numbers of U12-type introns at each node was calculated with a Fisher's exact test (scipy v1.5.2 (44)), and p -values were corrected for multiple-testing using the Holm step-down method as implemented in the Python library statsmodels (45) (v0.11.1).

Estimation of intron splicing efficiency and retention

For each annotated intron defined by CDS features from the longest isoform of each gene, splice junctions for the spliced (5' exon + 3' exon) and retained (5' exon + intron, intron + 3' exon) structures were created in silico using a custom Python script. RNA-seq reads

(accession numbers listed in the reannotation section) were then mapped in single-end mode to the junction constructs using Bowtie v1.2.2 (46) with parameter -m 1 to exclude multiply-mapped reads. Reads overlapping a junction by ≥ 5 nt were counted and corrected by the number of mappable positions on the associated junction construct. For each RNA-seq dataset, introns with no reads supporting the spliced form were excluded from the analysis, as were introns with no junctions supported by at least 10 reads. For all other introns, efficiency was calculated as the ratio of splice-supporting read coverage (C_s) over the total read coverage, which is just C_s plus the average of the retention-supporting read coverage (C_r), expressed as a percentage, i.e. $\frac{C_s}{(C_r/2) + C_s} \cdot 100\%$. Each intron's splicing efficiency was then computed as either the average across all RNA-seq samples containing mapped reads (for the bulk analysis), or using per-sample read counts (for the comparisons of neighboring introns from the same transcript).

To help validate our splicing efficiency results, we also employed an established method to evaluate intron retention using the same RNA-seq data. We obtained intron retention values for all annotated *P. polycephalum* introns with IRFinder (25) (v1.3.0), which produced an equivalent (inverted) pattern to our splicing efficiency metric (Figures S7-10). Introns with IRFinder warnings of "LowSplicing" and "LowCover" were excluded.

U12-type introns in *P. polycephalum* paralogs and non-canonical U12-type introns

Introns conserved across *P. polycephalum* paralogs were identified as described for homologous introns. We then examined all intron positions conserved between paralogs, and tabulated the intron types at each position. To determine the relative likelihood of a given U12-type intron being conserved as U12-type across paralogs, we calculated the relative probability of an intron *B* being U12-type conditioned on its paralogous intron *A* being 1) U12-type versus 2) U2-type as $\frac{P(B_{U12}|A_{U12})}{P(B_{U12}|A_{U2})}$, which results in a likelihood fold-increase of $\frac{(866/1074)}{(208/8695)} \approx \frac{0.806}{0.024} \approx 34$. This value is most likely conservative, as decreasing the stringency of U12-type classification results in a further increase in the relative likelihood (Table S2).

To avoid inclusion of spurious intron annotations representing artifacts of the RT-PCR process ("RTfacts", (47)) in our non-canonical intron analysis, we used a fairly simple heuristic to detect unexpectedly high similarities between extended sequences around the 5' and 3' splice sites. For each intron, we considered regions of 24 bp centered around the 5' and 3' splice sites (12 bp from the exon and 12 bp from the intron in each case) and used a 12 bp sliding window to compare every 5'SS 12mer against every 3'SS 12mer. For each 12mer pair, we defined their pairwise similarity *s* as $s = 1 - \frac{d}{l}$, where *d* is the Hamming distance between the two strings and *l* is their length in bp (i.e. 12), and treated the highest value found as the overall similarity score. Introns with similarity scores ≥ 0.916 (corresponding to one

mismatch between the pair of splice site 12mers) were considered possible RTfacts and were excluded (n = 1,624, 0.93% of 175,198 total introns).

In our survey of non-canonical introns in *P. polycephalum*, we took advantage of the greater number of conserved U12-type intron positions within paralogs (versus with other species) to gauge support for non-canonical U12-type intron boundaries present in our annotations. Of the non-canonical U12-type introns found in regions of good alignment between paralogs, 66% (42/63) contained the U12-type BPS motif 9-12 bp upstream of the 3'SS; in the smaller set conserved as introns between paralogs, the same motif was present in 73% (22/30). The BPS motif enrichment within these introns supports their identity as genuine non-canonical U12-type, and the distribution of the most common boundaries found within paralogs is consistent with the broader set of non-canonical U12-type introns (Figure S6B).

snRNP relative expression

Orthologs for components of the major and minor spliceosome (major: SF3a120/SAP114, SF3a60/SAP61, U1-70K, U1 A, U2 A'; minor: U11/U12 20K, 25K, 25K, 31K, 35K, and 65K) were identified via reciprocal BLASTP searches (as described in the section on ortholog identification) using the components' annotated human transcripts as queries (Table S3). For each species, a series of RNA-seq samples (curated by size and wild-type status) were aligned to the coding sequences of all available components using HISAT2, and the output processed with StringTie using the "-A" option to obtain per-transcript TPM values. Mean per-species TPM values across all RNA-seq samples for the U12- and U2-type components were then compared to calculate the U12/U2 expression ratios. An ANOVA test was performed on the group of ratios (p -value 8.8×10^{-13}), followed by pairwise two-tailed Mann-Whitney U tests between all combinations of ratios. The difference between *P. polycephalum* and every other species was significant at $p < 0.05$ after multiple-testing correction using the Holm step-down method as implemented in the Python library statsmodels v0.11.1 (45).

Supplemental Figures/Tables

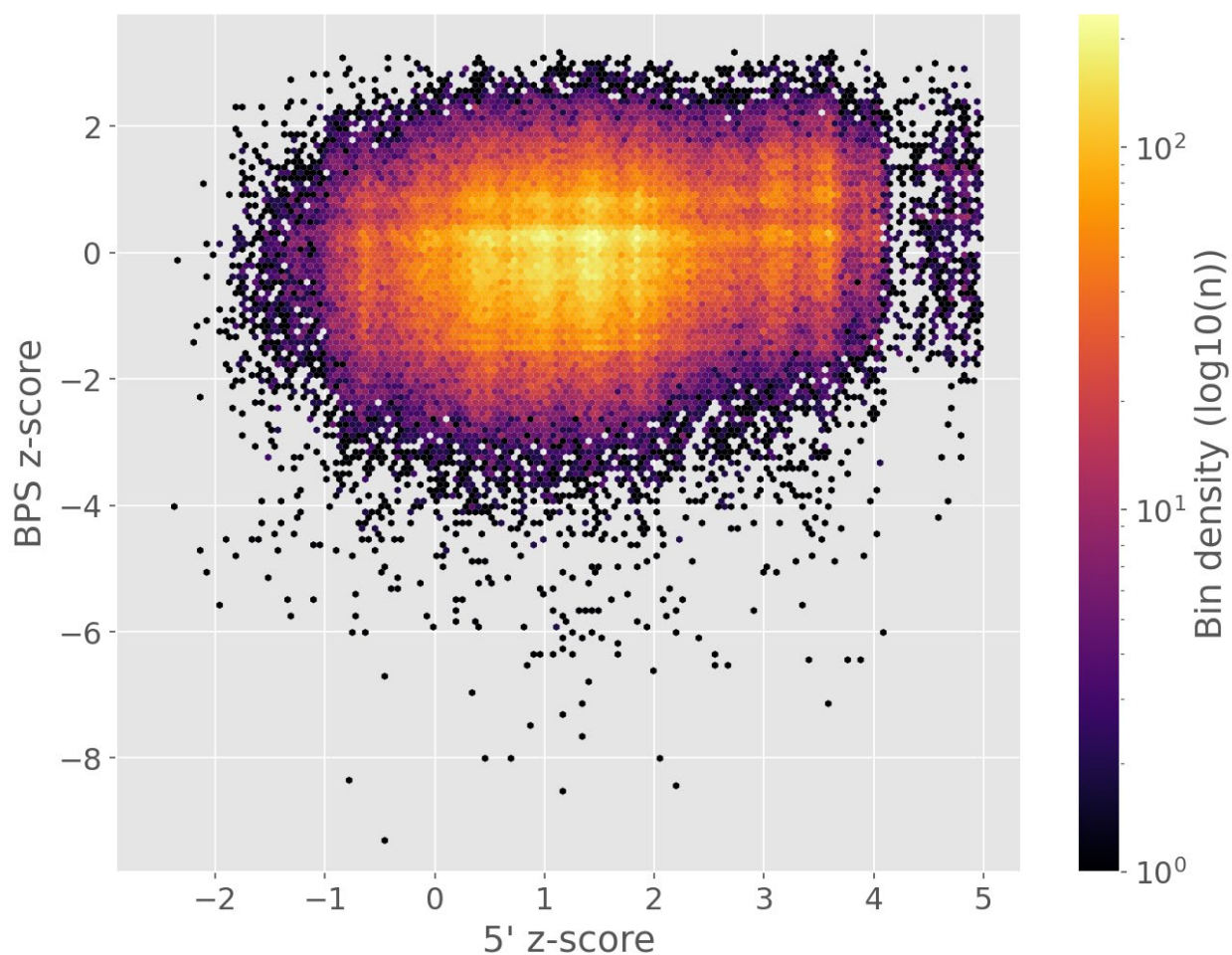


Figure S1. BPS-vs-5'SS score plot for all *P. polycephalum* introns under default settings via intronIC (18). The default PWMs used by intronIC are derived from human introns, and for divergent motifs like those present in *P. polycephalum* (especially the BPS motif) they fail to produce clear differentiation (i.e. separation of U12-type introns into a distinct cloud in the first quadrant). Curation of species-specific PWMs for *P. polycephalum* resulted in clearer differentiation along both axes (as in Figure 1D).

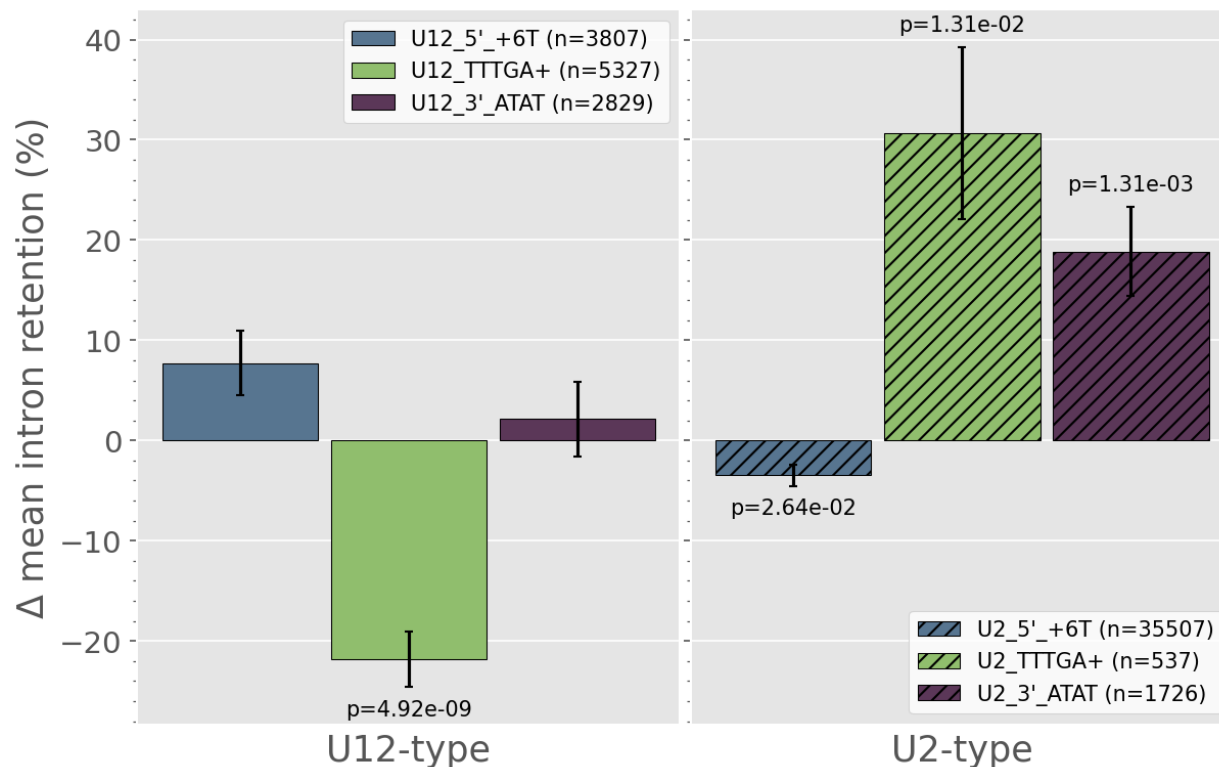


Figure S2. Relative intron retention for U12- (left) and U2-type (right) introns based on sequence features. Differences from the mean for each category are relative to all other introns of the same type. A negative/positive value indicates that introns with the given feature exhibit more/less efficient splicing relative to other introns of the same type. Features shown are "5'+6T", introns with a T at position +6 in the intron; "TTTGA+", introns with the TTTGA motif within the last 55 bases of the intron; "3'_ATAT", introns with the motif ATAT immediately downstream of the 3'SS.

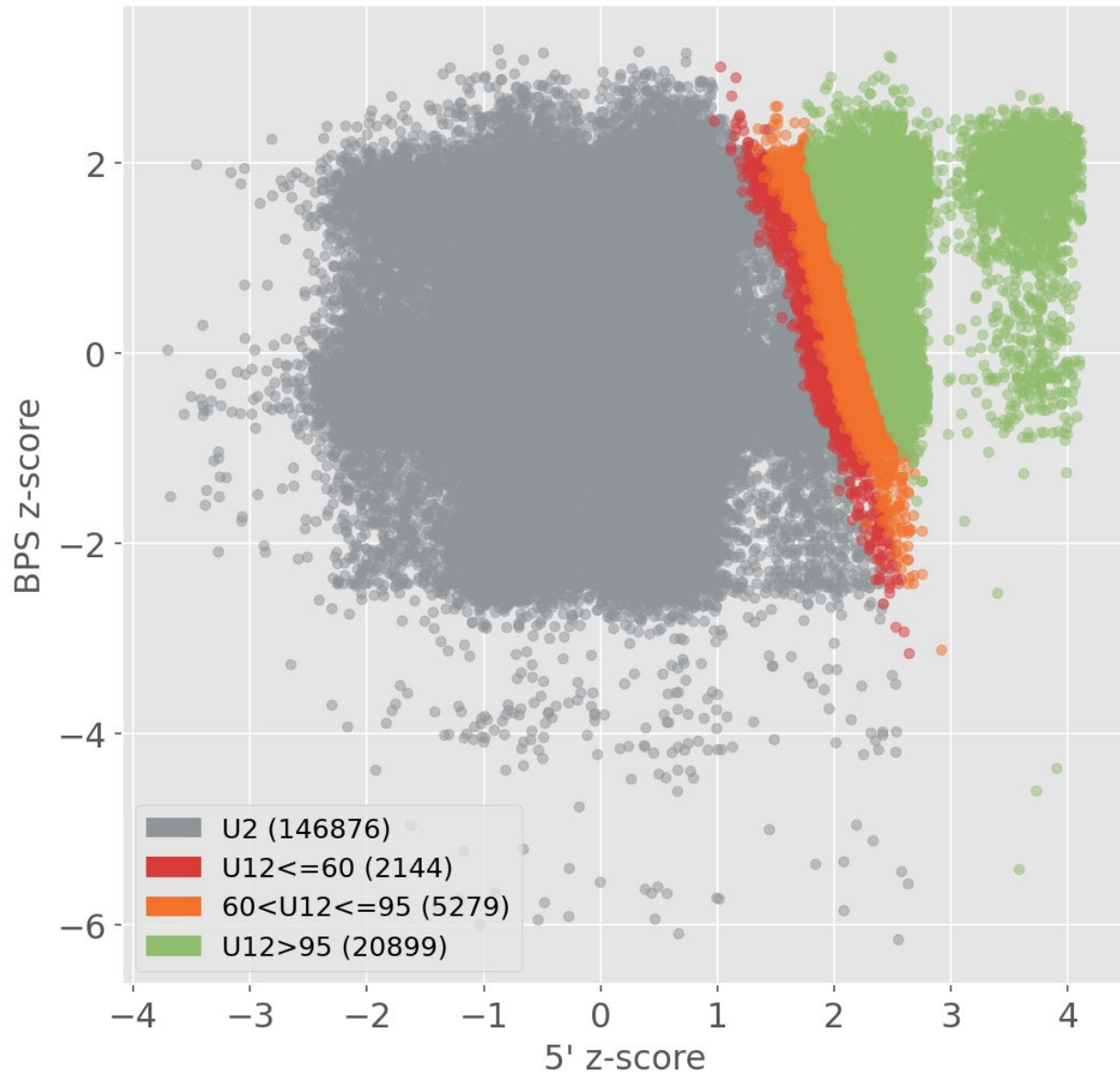


Figure S3. BPS-vs-5'SS score plot with assigned classifications for all *P. polycephalum* introns. The same underlying data as Figure 1C, where each point represents an intron, and the color indicates the U12-type probability classification (gray, U2-type; red, U12-type with probability ≤60%; orange, U12-type with probability 60-95%; green, U12-type with probability >95%). Throughout our analyses, only the >95% category were considered U12-type (unless explicitly stated otherwise).

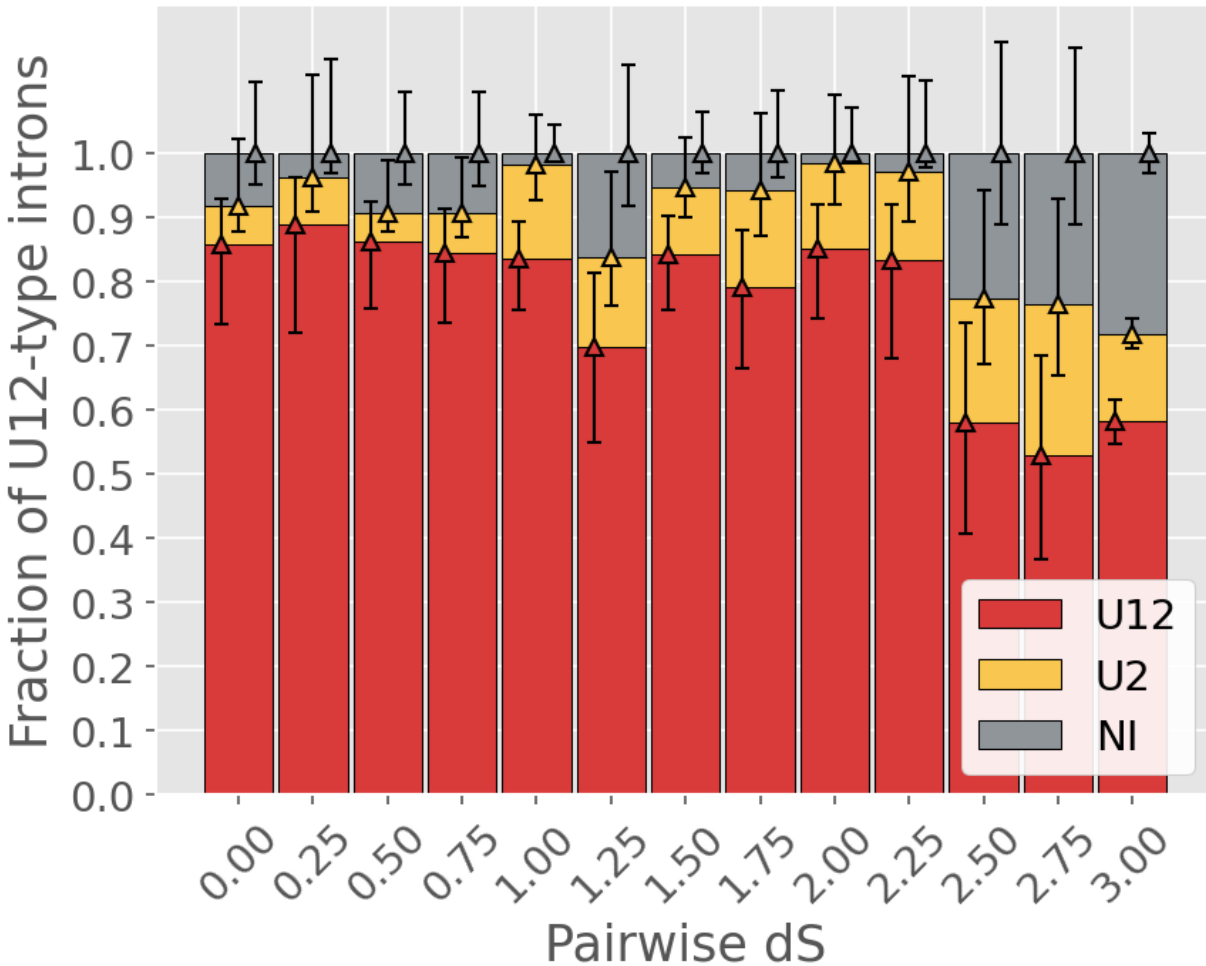


Figure S4. Between-paralog comparison provides little evidence for ongoing U12-type intron gain in *P. polycephalum*. For U12-type intron-containing paralog pairs sharing at least one intron of either type (to exclude recent retrogenes), pairwise dS values were used to bin all pairs into the range [0, 3]; dS values ≥ 3 were binned together. Within a given bin, each U12-type intron has one of three possible conservation states in its corresponding paralog: U12-type (red), U2-type (yellow) or no intron present (“no intron”, gray). These data suggest that there have not been major U12-type intron gains in *P. polycephalum* since a time corresponding to at least dS ~ 2.5 .

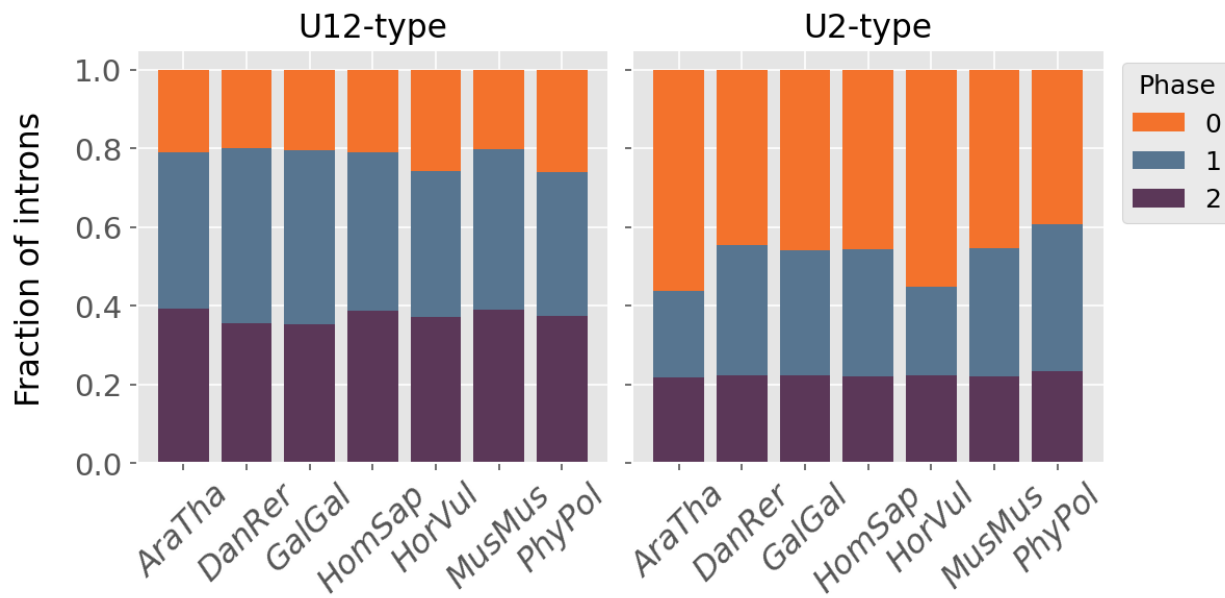


Figure S5. Phase distribution of U12- (left) and U2-type (right) introns across different species. U12-type introns in *P. polycephalum* (*PhyPol*) display a bias away from phase 0, as in other species, whereas U2-type introns show a bias against phase 2. For each species, only introns interrupting coding sequence from the longest isoform of each gene were included. See Table S1 for species abbreviations.

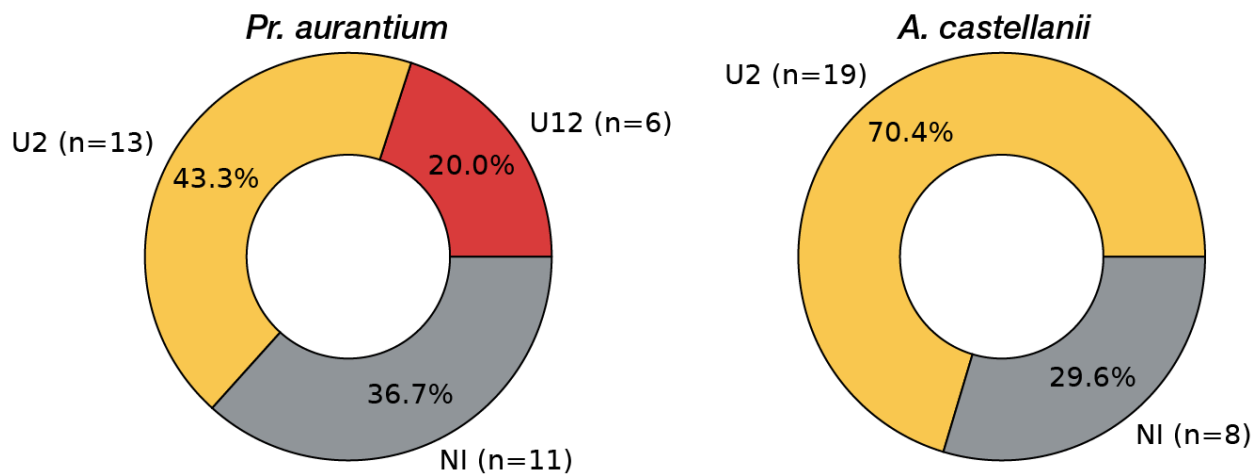


Figure S6. Ancestral U12-type introns in *P. polycephalum* are conserved as introns in other amoebozoans. Each pie chart shows the conservation status (red, U12-type; yellow, U2-type; gray, no intron) of the same ancestral set of *P. polycephalum* U12-type introns (introns conserved as U12-type with one or more non-amoebozoans) in the variosean amoeba *Protostelium aurantium* (left) and the discosean amoeba *Acanthamoeba castellanii* (right). In each case, a significant majority of the U12-type introns are conserved as introns. These

data suggest that these species have not undergone massive loss of U12-type introns, and thus that the unprecedented numbers of U12-type introns in *P. polycephalum* represents U12-type intron creation in *P. polycephalum*, not loss in related species.

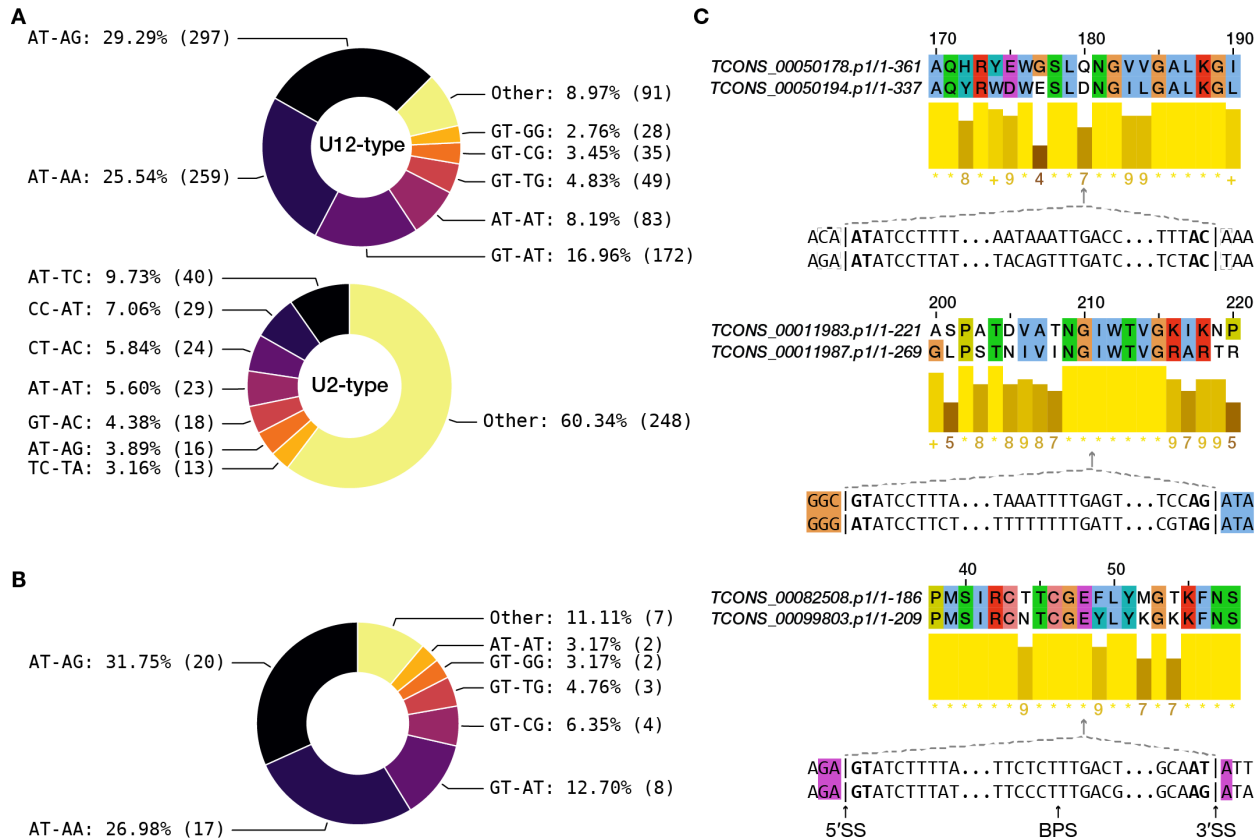


Figure S7. Non-canonical introns in *P. polycephalum*. **(A)** U12- (top) and U2-type (bottom) non-canonical intron subtypes (using a 60% probability threshold for the U12/U2-type classification instead of the 95% threshold used elsewhere e.g. Figure 2A, thereby including “likely” U12-type introns), highlighting the degree to which non-canonical U12-type introns are greatly enriched for a subset of boundary pairs. By contrast, the U2-type non-canonical subtype distribution is much more diffuse. **(B)** Distribution of subset non-canonical U12-type introns which are found in regions of good alignment between pairs of *P. polycephalum* paralogs (but not necessarily conserved as introns between pairs), increasing confidence that they are truly introns, showing general consistency with part A. **(C)** Example alignments of *P. polycephalum* paralogs, showing conserved U12-type introns (canonical and non-canonical). Coloring is based on chemical properties of the amino acids, and bars underneath each alignment represent chemical similarities of the aligned amino acids. Colored nucleotides before and after the intron splice sites correspond to the colors of the amino acid(s) in the alignment that are interrupted by the shared intron position. Transcript names appear in italics.

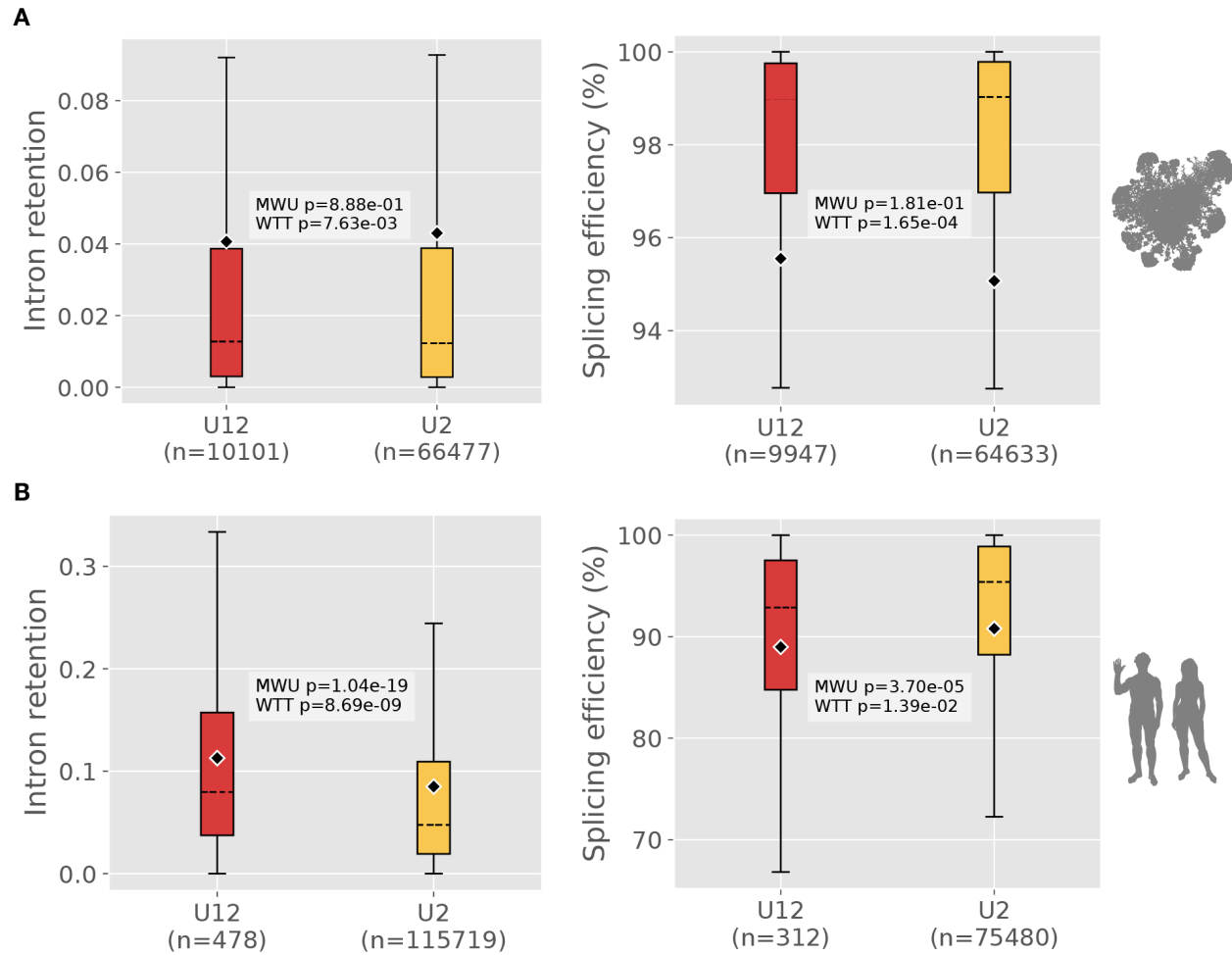


Figure S8. Comparison of intron retention (left) and splicing efficiency (right) in *P. polycephalum* and human. **(A)** Box plot of average intron retention and splicing efficiency data for *P. polycephalum* introns, showing that U12-type introns are neither more retained nor less efficiently-spliced than U2-type introns. Note that although the differences in means between U12- and U2-type introns are significant, this difference is inverted relative to data from other species. The left panel is the same as Figure 2C. **(B)** As in (A), but for *Homo sapiens*. Here, by both statistical measures there are significant differences between the two types of introns, with U12-type introns being more retained and less-efficiently spliced as has been reported elsewhere. MWU: two-tailed Mann-Whitney U test; WTT: two-tailed Welch's t-test. Note that the y-axis scales differ between plots.

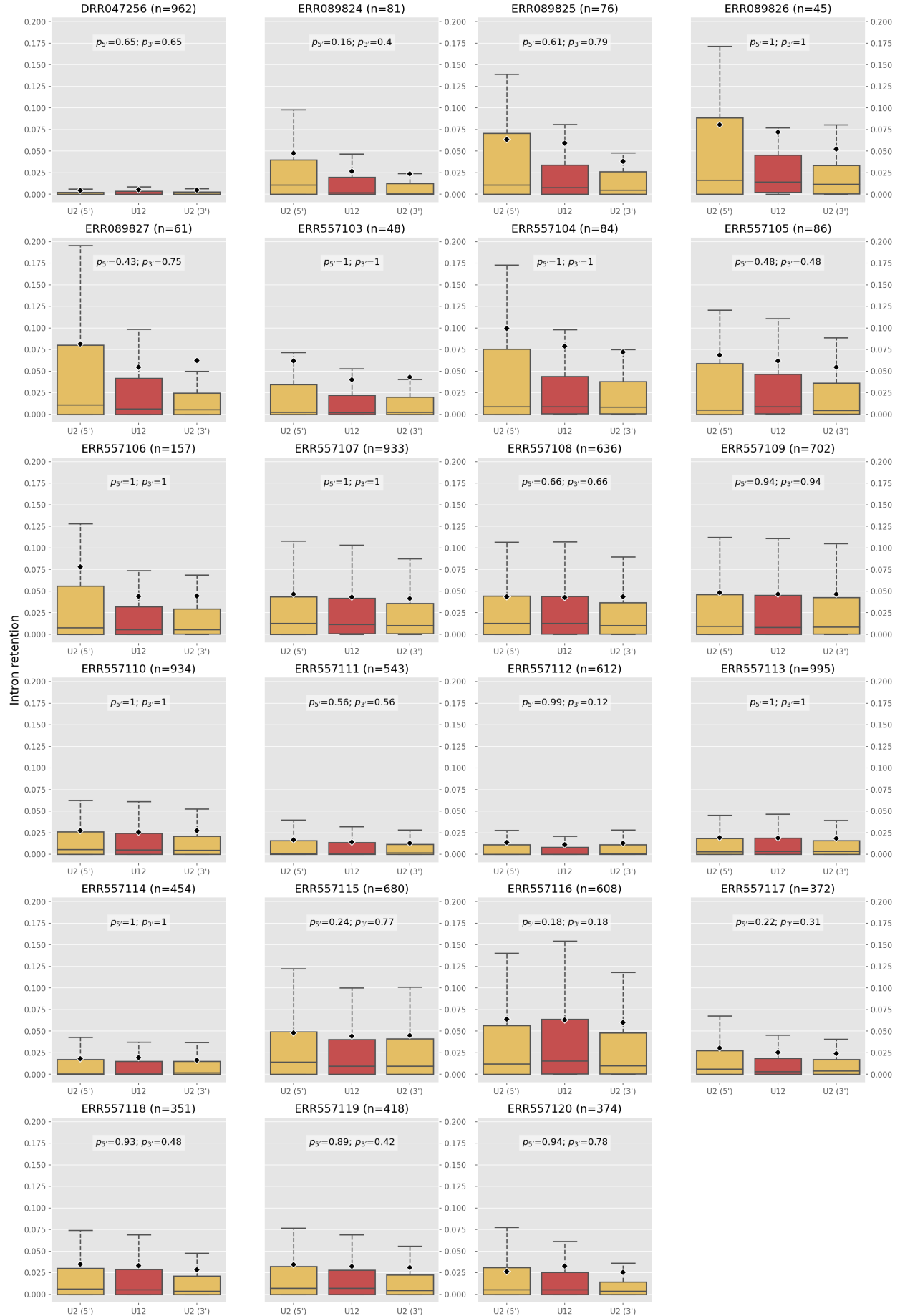


Figure S9. U12-type intron retention is not significantly different from that of neighboring U2-type introns in *P. polycephalum*. Each subplot represents data from a different RNA-seq sample (accession number in subplot title, along with number of introns represented), showing the distribution of intron retention values for U12-type (red) and neighboring U2-type (yellow) introns on both sides (left: 5', right: 3'). For each neighboring U2-type (defined as introns with U12-type probability scores < 5%) dataset, p -values vs. the corresponding U12-type data were obtained via Mann-Whitney U tests, and corrected for multiple testing using the Holm step-down method (reported as $p_{5'}$ and $p_{3'}$ for the 5' and 3' U2-type introns, respectively).

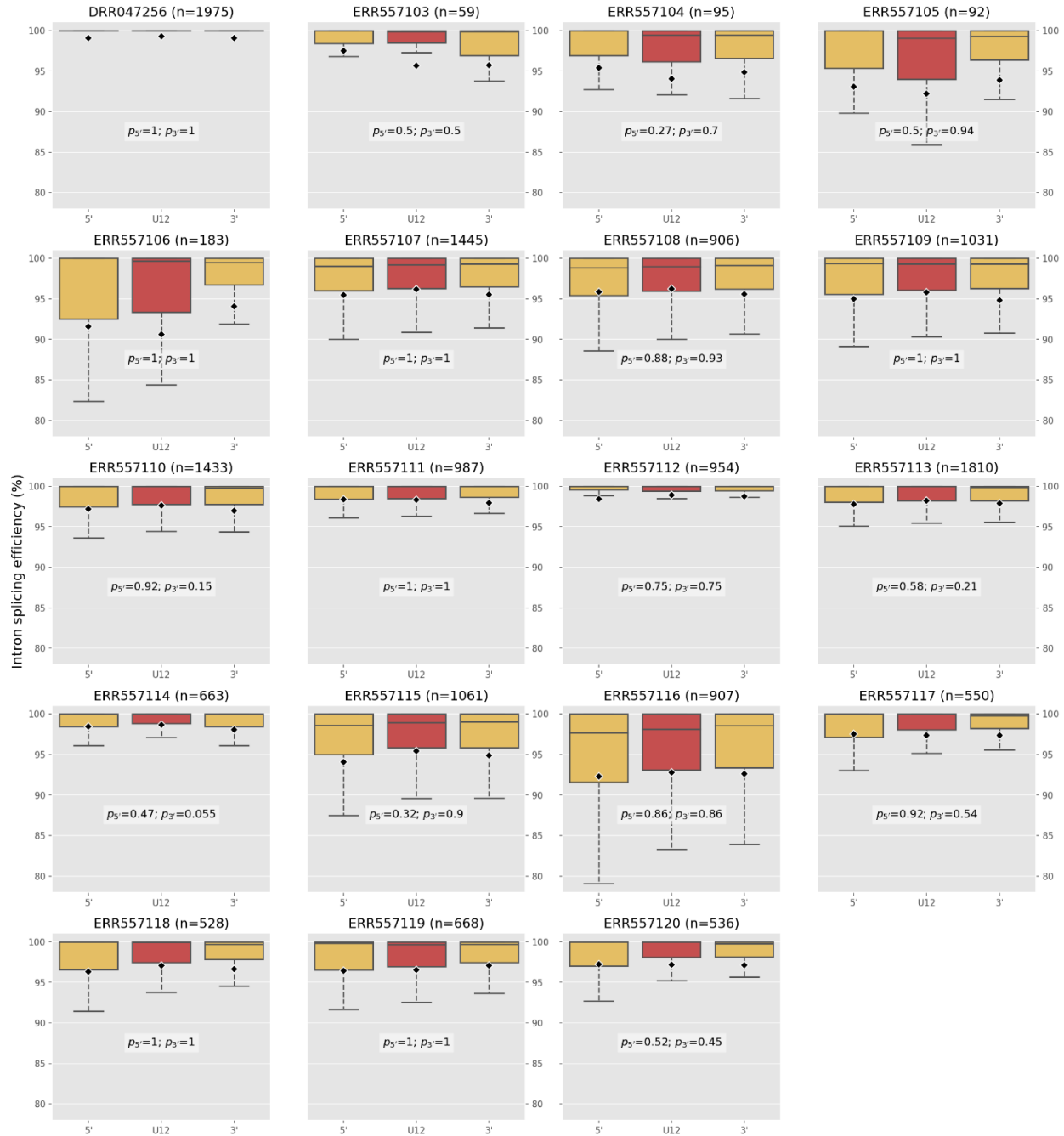


Figure S10. Following the same procedure as in Figure S8, but with splicing efficiency instead of intron retention in *P. polycephalum*, again showing no statistically significant differentiation between intron types.

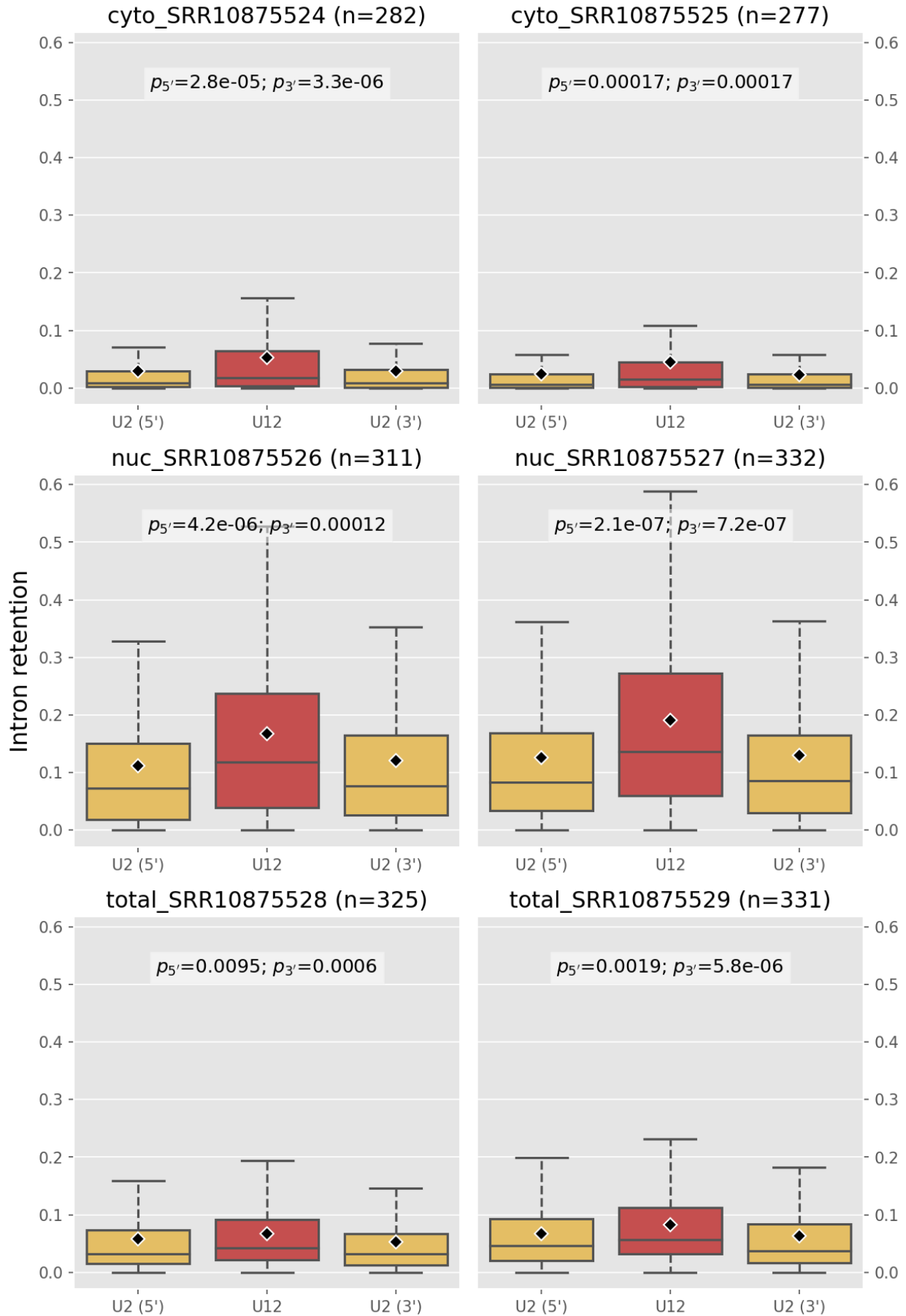


Figure S11. U12-type intron retention is significantly higher than that of neighboring U2-type introns in *Homo sapiens*. RNA-seq sets from three different protocols (cytosolic (“cyto”), nuclear (“nuc”) and total RNA, indicated in the titles of each subplot) were used to compare U12-type intron retention to that of neighboring U2-type introns. In each case, U12-type introns were more retained than their neighboring U2-type introns, in contrast to our results for *P. polycephalum* (Figure S9).

species	abbreviation	source	genome version	annotation version	updated with RNA-seq (Y/N)	notes
<i>Acanthamoeba castellanii</i>	AcaCas	Ensembl	strNEFF_v1	45	N	
<i>Amborella trichopoda</i>	AmbTri	JGI	291_v1.0	291_v1.0	Y	
<i>Amphimedon queenslandica</i>	AmpQue	Ensembl	Aqu1	45	N	
<i>Arabidopsis thaliana</i>	AraTha	Ensembl	TAIR10	40	N	
<i>Bigeloviella natans</i>	BigNat	Ensembl	Bigna1	45	N	
<i>Branchiostoma lanceolatum</i>	BraLan	Ensembl	BraLan2	45	N	
<i>Callorhynchus milii</i>	CallMil	Ensembl	6.1.3	98	N	
<i>Chlamydomonas reinhardtii</i>	ChiRei	Ensembl	v5.5	45	N	
<i>Ciona intestinalis</i>	CioInt	Ensembl	KH	98	N	
<i>Cyanophora paradoxa</i>	CyaPar	http://cyanophora.rutgers.edu	N/A	N/A	N	http://cyanophora.rutgers.edu/cyanophora:Cyanophora_paradoxa_MAKER_gene_predictions-022111-aa.fasta
<i>Danio rerio</i>	DanRer	Ensembl	GRCz11	95	N	
<i>Emiliania huxleyi</i>	EmiHux	Ensembl	CCMP1516_main_genome_assem_bly_v1.0	45	N	
<i>Galdieria sulphuraria</i>	GalSul	Ensembl	ASM34128v1	45	N	
<i>Gallus gallus</i>	GalGal	Ensembl	GRCg6a	98	N	
<i>Gossypium raimondii</i>	GosRai	JGI	221_v2.0	221_v2.1	N	
<i>Helobdella robusta</i>	HelRob	Ensembl	Helro1	45	N	
<i>Homo sapiens</i>	HomSap	Ensembl	GRCh38	95	N	
<i>Hordeum vulgare</i>	HorVul	Ensembl	Hv_IBSC_PGSSB_v2	40	N	
<i>Klebsormidium nitens</i>	KleNit	http://www.plantmorphogenesis.bio.titech.ac.jp	120824_klebsormidium_Scaffolds_v1.0	171026_klebsormidium_v1.1	N	http://www.plantmorphogenesis.bio.titech.ac.jp/~algae_genome_project/klebsormidium/v1_download.htm
<i>Limulus polyphemus</i>	LimPol	NCBI/GenBank	GCF_000517525.1 v2.1.2	2.1.2	N	
<i>Macaca mulatta</i>	MacMul	Ensembl	Mmul_10	98	N	
<i>Marchantia polymorpha</i>	MarPol	Ensembl	v1	45	N	
<i>Monosiga brevicollis</i>	MonBre	Ensembl	mx1_gca_000002865.V1.0	45	N	
<i>Mus musculus</i>	MusMus	Ensembl	GRCm38	95	N	
<i>Naegleria gruberi</i>	NaeGru	Ensembl	gca_000004985.V1	45	N	
<i>Physcomitrella patens</i>	PhyPat	Ensembl	Phypa_V3	45	N	
<i>Pinus taeda</i>	PinTae	https://treegenesdb.org/	v2.01	v2.01	Y	https://treegenesdb.org/FTP/Genomes/Pita/v2.01/ Indexed on GenBank as <i>Planoprotestelium fungivorum</i> ; originally described here https://doi.org/10.1093/jbe/evy011 , but under the corrected name <i>Protestelium aurantium</i> var. <i>fungivorm</i> , with <i>P. fungivorum</i> being used only in the supplemental material. Supporting evidence for the correct classification of the organism as <i>Pr. aurantium</i> can be found here https://doi.org/10.1111/jev.12475
<i>Protestelium aurantium</i>	ProAur	NCBI/GenBank	GCA_003024175.1_ASM302417v1	GCA_003024175.1_ASM302417v1	Y	
<i>Saccharomyces cerevisiae</i>	SacCer	Ensembl	R64-1-1	98	N	
<i>Salaginella moellendorffii</i>	SelMoe	JGI	v1	v1	Y	
<i>Salaginella moellendorffii</i>	SelMoe	Ensembl	v1.0	45	N	This version was used in the gene age analysis
<i>Strongylocentrotus purpuratus</i>	StrPur	Ensembl	Spur_3.1	45	N	
<i>Trichoplax adhaerens</i>	TriAdh	Ensembl	ASM15027v1	26	Y	

Table S1. Genome and annotation information for additional species.

U12/U2-type score threshold	Total U12:intron alignments	Total U2:intron alignments	U12:U12	U12:U2	P(U12 U12) / P(U12 U2)
95	1072	8683	864	208	33.64537887
60	1261	8494	1098	163	45.37450558

Table S2. Paralogous intron data for *P. polycephalum* under different U12-type probability thresholds. For each intron type, the number of introns of that type with corresponding paralogous introns (of either type) are shown (e.g. “Total U12: intron alignments” is the total number of U12-type introns found in regions of good alignment between paralogs where there is an intron in the same position in the other paralog). “U12:X” are the total numbers of U12-type introns whose paralogous introns are of type X, and the final column lists the relative likelihood of an intron being U12-type given that its paralogous intron is U12-type.

species	RNA-seq samples						
<i>Arabidopsis thaliana</i>	SRR5197911, SRR5197985, SRR5197986, SRR5820083, SRR6874228, SRR7663610, SRR7726615, SRR9019675, SRR9265357, SRR934391, SRR995072						
<i>Homo sapiens</i>	SRR1617450, SRR1617451, SRR1617452, SRR1617454, SRR1617461, SRR5442314, SRR5442315, SRR5442317						
<i>Drosophila melanogaster</i>	ERR1145740, ERR1145741, ERR1145742, ERR1145743, ERR1145746, ERR1145750, SRR10005788, SRR6652839, SRR6652840, SRR6665463, SRR7450965, SRR8949126						
<i>Danio rerio</i>	ERR3365998, SRR4017373, SRR5274769, SRR6384886, SRR8441448, SRR8922974, SRR8944755, SRR9966394, SRR9966395, SRR9966396						
spliceosomal system		gene ID	gene	species-specific transcript ID			
major	NP_005868	SF3a120, SAP114	AT1G14650.1	ENS DART00000061378	F Btr0083374	ENST00000215793	TCONS_00075169.p1
major	NP_006793	SF3a60, SAP61	AT5G06160.1	ENS DART00000160433	F Btr0078751	ENST00000373019	TCONS_00069287.p1
major	NP_004587	U1 A	-	ENS DART00000012018	-	ENST00000243563	TCONS_00026432.p1
major	NP_003080	U1-70K	AT3G50670.1	-	F Btr0079355	ENST00000598441	TCONS_00063511.p1
minor	NP_061976	U11/U12-20K	AT5G26749.2	ENS DART00000144510	F Btr0077965	ENST00000344318	TCONS_00023854.p1
minor	NP_078847	U11/U12-25K	AT3G07860.1	ENS DART0000006783	-	ENST00000383018	TCONS_00025616.p1
minor	NP_149105	U11/U12-31K (MADP1)	AT3G10400.1	ENS DART00000067172	-	ENST00000266529	TCONS_00052334.p1
minor	NP_851030	U11/U12-35K	AT2G43370.1	ENS DART00000172191	-	ENST00000412157	TCONS_00095609.p1
minor	NP_060089	U11/U12-65K	AT1G09230.1	ENS DART0000017427	F Btr0339103	ENST00000423855	TCONS_00087173.p1
major	NP_003081	U2 A'	AT1G09760.1	ENS DART00000128530	F Btr0088941	ENST00000254193	TCONS_00018540.p1

Table S3. RNA-seq accession numbers and spliceosomal component information used in the snRNP expression analysis (Figure 2D).

References

1. N. Sheth, X. Roca, M. L. Hastings, T. Roeder, A. R. Krainer, R. Sachidanandam, Comprehensive splice-site analysis using comparative genomics. *Nucleic Acids Res.* **34**, 3955–3967 (2006).
2. J. Janice, A. Pande, J. Weiner, C. F. Lin, W. Makalowski, U12-type Spliceosomal Introns of Insecta. *Int. J. Biol. Sci.* **8**, 344–352 (2012).
3. J. J. Turunen, E. H. Niemelä, B. Verma, M. J. Frilander, The significant other: Splicing by the minor spliceosome. *Wiley Interdiscip. Rev. RNA.* **4**, 61–76 (2013).
4. W. Y. Tarn, J. A. Steitz, Highly diverged U4 and U6 small nuclear RNAs required for splicing rare AT-AC introns. *Science.* **273**, 1824–1832 (1996).
5. S. L. Hall, R. A. Padgett, Requirement of U12 snRNA for in vivo splicing of a minor class of eukaryotic nuclear pre-mRNA introns. *Science.* **271**, 1716–1718 (1996).
6. C.-F. Lin, S. M. Mount, A. Jarmołowski, W. Makalowski, Evolutionary dynamics of U12-type spliceosomal introns. *BMC Evol. Biol.* **10**, 47 (2010).
7. W. Zhu, V. Brendel, Identification, characterization and molecular phylogeny of U12-dependent introns in the *Arabidopsis thaliana* genome. *Nucleic Acids Res.* **31**, 4561–4572 (2003).
8. A. G. Russell, J. M. Charette, D. F. Spencer, M. W. Gray, An early evolutionary origin for the minor spliceosome. *Nature.* **443**, 863–866 (2006).
9. E. H. Niemelä, A. Oghabian, R. H. J. Staals, D. Greco, G. J. M. Pruijn, M. J. Frilander, Global analysis of the nuclear processing of transcripts with unspliced U12-type introns by the exosome. *Nucleic Acids Res.* **42**, 7358–7369 (2014).
10. J. Singh, R. A. Padgett, Rates of in situ transcription and splicing in large human genes. *Nat. Struct. Mol. Biol.* **16**, 1128–1133 (2009).
11. A. A. Patel, M. McCarthy, J. A. Steitz, The splicing of U12-type introns can be a rate-limiting step in gene expression. *EMBO J.* **21**, 3804–3815 (2002).
12. S. Bartschat, T. Samuelsson, U12 type introns were lost at multiple occasions during evolution. *BMC Genomics.* **11**, 106 (2010).
13. M. D. Lopez, M. Alm Rosenblad, T. Samuelsson, Computational screen for spliceosomal RNA genes aids in defining the phylogenetic distribution of major and minor spliceosomal components. *Nucleic Acids Res.* **36**, 3001–3010 (2008).
14. P. Schaap, I. Barrantes, P. Minx, N. Sasaki, R. W. Anderson, M. Bénard, K. K. Biggar, N. E.

- Buchler, R. Bundschuh, X. Chen, C. Fronick, L. Fulton, G. Golderer, N. Jahn, V. Knoop, L. F. Landweber, C. Maric, D. Miller, A. A. Noegel, R. Peace, G. Pierron, T. Sasaki, M. Schallenberg-Rüdinger, M. Schleicher, R. Singh, T. Spaller, K. B. Storey, T. Suzuki, C. Tomlinson, J. J. Tyson, W. C. Warren, E. R. Werner, G. Werner-Felmayer, R. K. Wilson, T. Winckler, J. M. Gott, G. Glöckner, W. Marwan, The Physarum polycephalum Genome Reveals Extensive Use of Prokaryotic Two-Component and Metazoan-Type Tyrosine Kinase Signaling. *Genome Biol. Evol.* **8**, 109–125 (2015).
15. M. Stanke, M. Diekhans, R. Baertsch, D. Haussler, Using native and syntenically mapped cDNA alignments to improve de novo gene finding. *Bioinformatics*. **24**, 637–644 (2008).
 16. F. A. Simão, R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva, E. M. Zdobnov, BUSCO: Assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* (2015), doi:10.1093/bioinformatics/btv351.
 17. C. Burge, P. A. Sharp, Classification of introns: U2-type or U12-type. *Cell*. **91**, 875–879 (1997).
 18. D. C. Moyer, G. E. Larue, C. E. Hershberger, S. W. Roy, R. A. Padgett, Comprehensive database and evolutionary dynamics of U12-type introns. *Nucleic Acids Res.* (2020), doi:10.1093/nar/gkaa464.
 19. A. Levine, R. Durbin, A computational scan for U12-dependent introns in the human genome sequence. *Nucleic Acids Res.* **29**, 4006–4013 (2001).
 20. C. B. Burge, R. a. Padgett, P. a. Sharp, Evolutionary fates and origins of U12-type introns. *Mol. Cell*. **2**, 773–785 (1998).
 21. M. K. Basu, W. Makalowski, I. B. Rogozin, E. V. Koonin, U12 intron positions are more strongly conserved between animals and plants than U2 intron positions. *Biol. Direct*. **3**, 19 (2008).
 22. I. Younis, K. Dittmar, W. Wang, S. W. Foley, M. G. Berg, K. Y. Hu, Z. Wei, L. Wan, G. Dreyfuss, Minor introns are embedded molecular switches regulated by highly unstable U6atac snRNA. *Elife*. **2**, e00780 (2013).
 23. H. Pessa, A. Ruokolainen, M. J. Frilander, The abundance of the spliceosomal snRNPs is not limiting the splicing of U12-type introns. *RNA*. **12**, 1883–1892 (2006).
 24. W. Y. Tarn, J. a. Steitz, A novel spliceosome containing U11, U12, and U5 snRNPs excises a minor class (AT-AC) intron in vitro. *Cell*. **84**, 801–811 (1996).
 25. R. Middleton, D. Gao, A. Thomas, B. Singh, A. Au, J. J.-L. Wong, A. Bomane, B. Cosson, E. Eyraas, J. E. J. Rasko, W. Ritchie, IRFinder: assessing the impact of intron retention on mammalian gene expression. *Genome Biol.* **18**, 51 (2017).

26. R. C. Dietrich, R. Incorvaia, R. A. Padgett, Terminal Intron Dinucleotide Sequences Do Not Distinguish between U2- and U12-Dependent Introns. *Mol. Cell.* **1**, 151–160 (1997).
27. J. T. Huff, D. Zilberman, S. W. Roy, Mechanism for DNA transposons to generate introns on genomic scales (2016), doi:10.1038/nature20110.
28. S. Henriët, B. Colom Sanmartí, S. Sumic, D. Chourrout, Evolution of the U2 Spliceosome for Processing Numerous and Highly Diverse Non-canonical Introns in the Chordate *Fritillaria borealis*. *Curr. Biol.* **29**, 3193–3199.e4 (2019).
29. M. G. Grabherr, B. J. Haas, M. Yassour, J. Z. Levin, D. a. Thompson, I. Amit, X. Adiconis, L. Fan, R. Raychowdhury, Q. Zeng, Z. Chen, E. Mauceli, N. Hacohen, A. Gnirke, N. Rhind, F. di Palma, B. W. Birren, C. Nusbaum, K. Lindblad-Toh, N. Friedman, A. Regev, Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat. Biotechnol.* **29**, 644–652 (2011).
30. D. Kim, J. M. Paggi, C. Park, C. Bennett, S. L. Salzberg, Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* **37**, 907–915 (2019).
31. M. Pertea, D. Kim, G. M. Pertea, J. T. Leek, S. L. Salzberg, Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat. Protoc.* **11**, 1650–1667 (2016).
32. UniProt Consortium, The universal protein resource (UniProt). *Nucleic Acids Res.* **36**, D190–5 (2008).
33. B. J. Haas, A. Papanicolaou, M. Yassour, M. Grabherr, P. D. Blood, J. Bowden, M. B. Couger, D. Eccles, B. Li, M. Lieber, M. D. Macmanes, M. Ott, J. Orvis, N. Pochet, F. Strozzi, N. Weeks, R. Westerman, T. William, C. N. Dewey, R. Henschel, R. D. Leduc, N. Friedman, A. Regev, De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nat. Protoc.* **8**, 1494–1512 (2013).
34. B. J. Haas, Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. *Nucleic Acids Res.* **31**, 5654–5666 (2003).
35. K. J. Hoff, S. Lange, A. Lomsadze, M. Borodovsky, M. Stanke, BRAKER1: Unsupervised RNA-Seq-based genome annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics* (2015), doi:10.1093/bioinformatics/btv661.
36. G. Pertea, M. Pertea, GFF Utilities: GffRead and GffCompare. *F1000Res.* **9**, 304 (2020).
37. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, D. J. Lipman, Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
38. C. Camacho, G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer, T. L. Madden,

- BLAST+: architecture and applications. *BMC Bioinformatics*. **10**, 421 (2009).
39. M. A. Larkin, G. Blackshields, N. P. Brown, R. Chenna, P. A. McGettigan, H. McWilliam, F. Valentin, I. M. Wallace, A. Wilm, R. Lopez, J. D. Thompson, T. J. Gibson, D. G. Higgins, Clustal W and Clustal X version 2.0. *Bioinformatics*. **23**, 2947–2948 (2007).
 40. S. W. Roy, A. Fedorov, W. Gilbert, Large-scale comparison of intron positions in mammalian genes shows intron loss but no gain. *Proc. Natl. Acad. Sci. U. S. A.* **100**, 7158–7162 (2003).
 41. F. Sievers, D. G. Higgins, Clustal Omega. *Curr. Protoc. Bioinformatics*. **2014**, 3.13.1–3.13.16 (2014).
 42. Z. Yang, PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**, 1586–1591 (2007).
 43. S. Kang, A. K. Tice, F. W. Spiegel, J. D. Silberman, T. Pánek, I. Cepicka, M. Kostka, A. Kosakyan, D. M. C. Alcântara, A. J. Roger, L. L. Shadwick, A. Smirnov, A. Kudryavtsev, D. J. G. Lahr, M. W. Brown, Between a Pod and a Hard Test: The Deep Evolution of Amoebae. *Mol. Biol. Evol.* **34**, 2258–2270 (2017).
 44. P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, SciPy 1.0 Contributors, SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods*. **17**, 261–272 (2020).
 45. S. Seabold, J. Perktold, in *Proceedings of the 9th Python in Science Conference* (Austin, TX, 2010), vol. 57, p. 61.
 46. B. Langmead, Aligning short sequencing reads with Bowtie. *Curr. Protoc. Bioinformatics*. **32**, 11–17 (2010).
 47. S. W. Roy, M. Irimia, When good transcripts go bad: Artifactual RT-PCR “splicing” and genome analysis. *Bioessays*. **30**, 601–605 (2008).