# Report

# Birth of New Spliceosomal Introns in Fungi by Multiplication of Introner-like Elements

Ate van der Burgt,[1,*] Edouard Severing,[2]
Pierre J.G.M. de Wit,[1,3] and Jérôme Collemare[1,*]
[1]Laboratory of Phytopathology
[2]Laboratory of Bioinformatics
Wageningen University, 6708PB Wageningen,
The Netherlands
[3]Department of Botany and Microbiology, King Saud
University, Riyadh, Saudi Arabia 11514

## Summary

Spliceosomal introns are noncoding sequences that separate exons in eukaryotic genes and are removed from pre-messenger RNAs by the splicing machinery. Their origin has remained a mystery in biology since their discovery [1, 2] because intron gains seem to be infrequent in many eukaryotic lineages [3, 4]. Although a few recent intron gains have been reported [5, 6], none of the proposed gain mechanisms [7] can convincingly explain the high number of introns in present-day eukaryotic genomes. Here we report on particular spliceosomal introns that share high sequence similarity and are reminiscent of introner elements [8]. These elements multiplied in unrelated genes of six fungal genomes and account for the vast majority of intron gains in these fungal species. Such introner-like elements (ILEs) contain all typical characteristics of regular spliceosomal introns (RSIs) [9, 10] but are longer and predicted to harbor more stable secondary structures. However, dating of multiplication events showed that they degenerate in sequence and length within 100,000 years to eventually become indistinguishable from RSIs. We suggest that ILEs not only account for intron gains in six fungi but also in ancestral eukaryotes to give rise to most RSIs by a yet unknown multiplication mechanism.

## Results and Discussion

### Characterization of a New Type of Spliceosomal Intron that Is Able to Multiply in Unrelated Genes

Spliceosomal introns are one of the key innovations of eukaryotes [1]. They are an important component of the eukaryotic gene structure mainly because they enlarge the proteome diversity by alternative splicing and regulate gene expression at the posttranscriptional level [11, 12]. Canonical regular spliceosomal introns (RSIs) share GT/AG donor and acceptor sites that are required for their recognition and removal by the spliceosome [10]. Although most RSIs contain branch point sequences and polypyrimidine tracts involved in the splicing mechanism, the sequences of RSIs are usually not conserved [9, 10, 13]. Since their discovery, the origin of introns has remained a mystery to molecular biologists. Evolution of the eukaryotic gene structure seems to have been predominated by intron loss [4] and the theoretically calculated

intron gain rates cannot explain the large number of introns in present-day eukaryotic genomes [14]. However, extensive recent intron gains have been reported in the microcrustacean *Daphnia pulex* [5], the fungus *Mycosphaerella graminicola* [6], and possibly in the green alga *Micromonas pusilla* [8] and the urochordate *Oikopleura dioica* [15]. These reports suggest that the number of introns in a given genome is not only subject to losses but also to substantial gains. In this study, we report on a particular type of spliceosomal introns that show a high level of sequence similarity and some reminiscence of introner elements found in *Micromonas* [8]. These so-called introner-like elements (ILEs) likely originate from multiplication of a discreet number of ancestral elements that are present in related fungal genomes. Although they are typical spliceosomal introns, ILEs are significantly longer and predicted to fold into more stable secondary structures than RSIs. Rigorous intron gain analyses in six fungal species revealed that the vast majority of gained introns are ILEs. By analyzing closely related fungi that diverged less than 100,000 years ago, we could show that the majority of newborn ILEs rapidly degenerate in length, sequence, and stability to become indistinguishable from RSIs. We propose that ILEs are the predecessors of RSIs in these six fungal species. This multiplication mechanism might also be involved in intron gains in other fungi and possibly in ancestral eukaryotic lineages.

### Identification of Near-Identical Introns in Six Fungal Genomes

Using a simple BlastN-based method, numerous introns with near-identical sequences could be identified in the intronome of the Dothideomycete fungus *Cladosporium fulvum*. Our analysis thoroughly excluded repeated sequences originating from repetitive elements, segmental duplications, or recombinant genes. This observation corroborates recent findings of near-identical intronic sequences in the marine picoeukaryote *Micromonas* [8], the tunicate *Oikopleura dioica* [15], and the Dothideomycete fungus *M. graminicola* [6]. In *Micromonas*, such repeat sequences are located within introns and were called introner elements [8]. The analysis was extended to the intronomes of 22 additional Ascomycete fungi and one Basidiomycete. The use of hidden Markov models (HMMs) corresponding to near-identical introns resulted in the identification of 45 to 538 near-identical introns in six different fungi: *C. fulvum*, *Dothistroma septosporum*, *M. graminicola*, *Mycosphaerella fijiensis*, *Hysterium pulicare*, and *Stagonospora nodorum*. In each fungus, these introns could be grouped into one to eight different clusters, each likely originating from a single ancestral element that had been multiplied. Each cluster contains between 10 and 180 members. Thereafter, they will be called ILEs to distinguish them from introner elements found in *Micromonas* [8]. According to expressed sequence tag (EST) support in the different species, ILEs are introns that are spliced out (see Table S1 available online). Although the available EST data for some fungal species is limited, EST support for ILEs is slightly higher than for the entire intronome (Table S1).

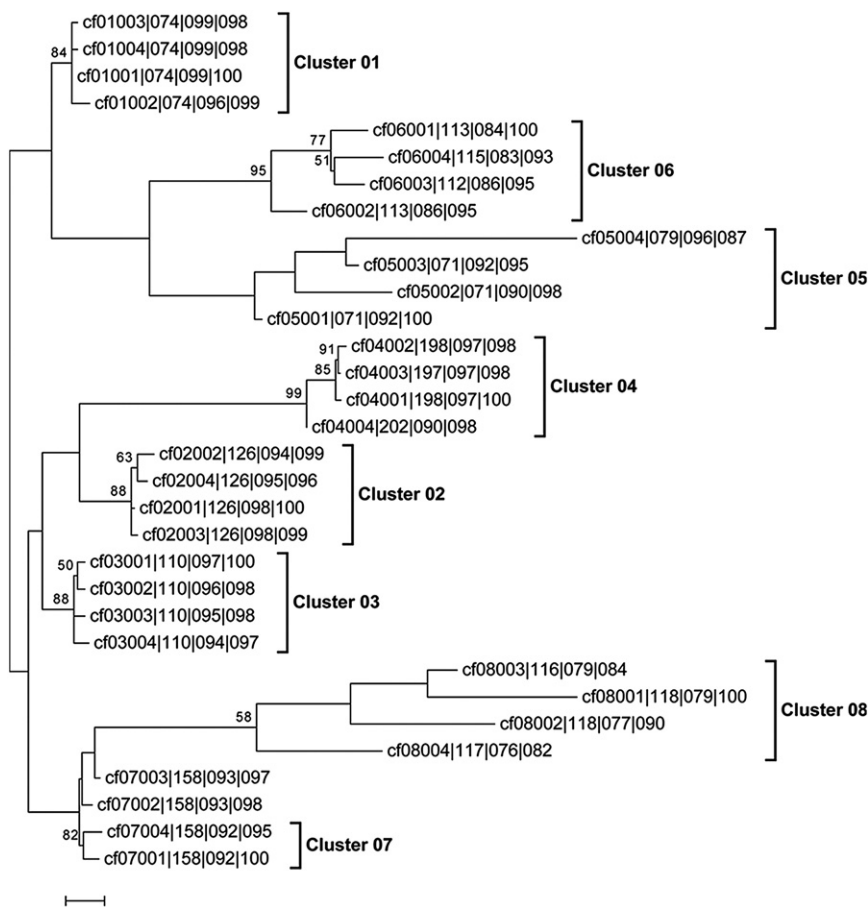*Correspondence: ate.vanderburgt@wur.nl (A.v.d.B.), jerome.collemare@wur.nl (J.C.)

**Figure 1. Introner-like Elements Originate from Common Ancestor Elements**

The four ILEs with highest HMM expect similarity of each *Cladosporium fulvum* cluster were aligned to construct a maximum likelihood phylogenetic tree. Identifiers contain ILE number|length (nt)|pairwise identity (%)|HMM expect similarity (%). ILE cluster number is indicated next to brackets. The midpoint rooting method was used to estimate the root of the tree. Only bootstrap values over 50 are shown. Scale indicates 0.1 substitutions per site. See also Figure S1 and Table S2.

is given for *C. fulvum* in Figure 2A). In addition, 96% of all ILEs contain a predicted CURAY branch point sequence, 95% contain a 5′ polypyrimidine tract, and 76% contain both 5′ and 3′ polypyrimidine tracts, frequencies that are higher than those reported for fungal RSIs [9]. Introns that lack an identifiable branch point sequence likely contain another noncanonical sequence that can be recruited as a branch point. Indeed, many introns lacking an identifiable branch point sequence have EST support for proper splicing (Table S1). Similar to RSIs [16], ILE clusters with the highest number of members show a preference for insertion in AG/GY sites and in phase 0 of coding sequences. They do, however, show a slight bias for being present in the center of genes in contrast to RSIs, which are more frequently found at the 5′ end of genes (Figure 2B; Figure S2). The biased location of RSIs was reportedly due to intron losses that primarily occur at the 3′ end of genes [17]. Altogether, these hallmark features suggest that ILEs are model RSIs and, more importantly, they can be perfectly spliced by the spliceosome immediately after multiplication and insertion in a new location.

## ILE Clusters Share Common Origins in Several Fungal Species

Phylogenetic analyses for each species showed that most ILE clusters are monophyletic clades, indicating that all elements of a given cluster share the same origin (Figure 1; Figure S1). A similarity matrix suggests that most ILE clusters in *M. graminicola* are related to each other and that clusters mf04 and mf05 in *M. fijiensis* are related to cluster cf01 in *C. fulvum* (Figure S1). In addition, the closely related fungi *C. fulvum* and *D. septosporum* share three ILE clusters (cf04/ds03, cf06/ds02, and cf08/ds04). These results indicate that ILE clusters present in different fungal species originate from the multiplication of a single ancestral element that was present before species divergence. Using HMMs to search the intronomes of closely related fungi not only showed that some of the identified ILE clusters contain additional members (mf04/mf05, cf05, cf06/ds02, mg05) but also revealed initially not identified shared clusters (Table S2). The intronomes of *D. septosporum* and *Septoria musiva* seem to contain less-conserved ILEs belonging to clusters cf07 and cf01/mf04/mf05, respectively. These results argue for the presence of ancestral ILEs that cannot easily be detected in other fungal intronomes.

## ILEs Harbor All Hallmark Features of Spliceosomal Introns

We characterized ILEs in more details to distinguish them from RSIs. Sequence analysis of all ILEs showed that they are genuine spliceosomal introns of which 99% contain canonical acceptor and donor sites (a representative example

## ILEs Are Predicted to be More Stable than RSIs but Are Also Prone to Degeneration

Although ILEs are model spliceosomal introns, they also have particular features that distinguish them from RSIs. Indeed, ILEs are longer than RSIs and show different length distributions with multiple peaks that correspond to the optimum lengths of different clusters (Figure 2C; Figure S2). In a given cluster, ILEs with the lowest identity are predominantly those showing sequence deletions (Figure S3). Further analyses confirmed a positive correlation between pairwise identity and length of all ILEs (Figure 3A). Substitutions and minor insertions are observed over ILEs' full length, but deletions occur less often around the conserved branch point sequence at the 3′ end (Figure 3B). This bias could be the result of selection pressure to retain splicing features. These observations suggest that ILEs may lose their ability to multiply due to degeneration in length and sequence. Although RSIs have some secondary structures that can facilitate splicing [12], the lower predicted Gibbs free energy ($\Delta G$) of ILEs suggests a significant greater molecular stability
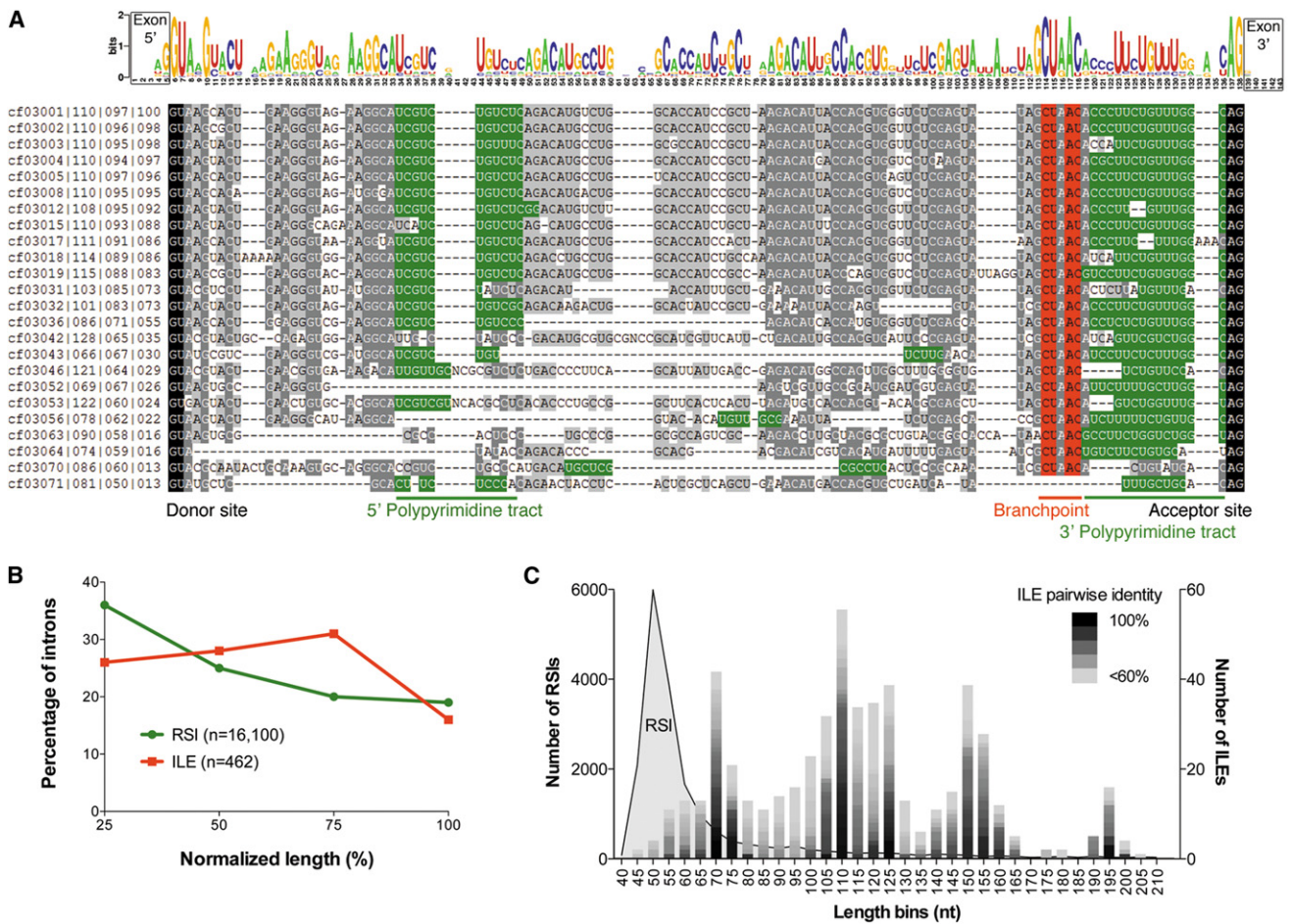
Figure 2. Characteristics of Introner-like Elements

Each panel presents representative results obtained for ILEs from *Cladosporium fulvum*. Similar results were obtained for the five other fungal species.

(A) Alignment and consensus sequence of selected ILEs. Identifiers are as shown in Figure 1. Colors show splicing elements typical of RSIs.

(B) Distribution of RSIs and ILEs within *C. fulvum* genes. Genes were divided in four sections expressed as percentages. Number of RSIs and ILEs in each section was counted and expressed as a percentage of all RSIs and ILEs, respectively.

(C) Length distribution of RSIs and ILEs from *C. fulvum*. Scale numbers on the x axis indicate the shortest length of 5 nt bins. See also Figure S2.

(Figure 3C). Remarkably, the increase in *ΔG* values of ILEs correlates with their degree of degeneration, until *ΔG* values of RSIs are eventually reached (Figures 3A and 3C). The low *ΔG* values of ILEs are explained by their predicted alignment-based secondary structures that often consist of three stem-loops containing many G-U pairs (Figures 4A and 4B). Consistent with the similarity matrix analysis, structure predictions for all *M. graminicola* ILE clusters suggest that they all have a common structure due to stretches of identical nucleotides (Figure 4B; Figure S4). Moreover, alignments of related ILEs revealed many compensatory mutations that conserve hairpin structures. Together, these observations argue for evolutionary constraints that preserve ILE secondary structures. We propose that they are an important feature because such predicted stable secondary structures are known for noncoding RNAs with specific functions [18]. Overall, our results strongly suggest that highly structured ILEs are likely mobile, but they are prone to degenerate mainly through deletions. They seem to gradually evolve to become RSIs that lost the ability to multiply and lack conserved predicted secondary structures. Thus, we hypothesize that RSIs might originate from ILEs.

## ILEs Account for the Vast Majority of Intron Gains in Six Fungi

The proposed hypothesis for the origin of RSIs implies that recent intron gains in fungi are mainly the result of ILE multiplication. Previous studies suggested that intron losses prevail over intron gains in most eukaryotic lineages [4], although several extensive gains have been reported as well [5, 6]. The balanced rates estimated in fungi are consistent with the presence of active ILEs [4, 14, 17]. By using up to six nodes between *C. fulvum/D. septosporum* and the most distant outgroup, the Basidiomycete *Cryptococcus neoformans*, single intron gains in *Dothideomycetes* could confidently be assigned as intron presence in only one of all species included in this study (Figure 5A). Thirteen single gains were identified in both *C. fulvum* and *D. septosporum* when using the maximum number of outgroups, which allowed inspection of 951 orthologs only. As outgroups are removed and more orthologs become available for inspection, the number of single gains increased to 199 (Table S3). Strikingly, on average, 50% of single gains originate from ILEs, irrespective of the number of outgroups used in the analysis (Figure 5B). It is noteworthy that 75% to 90% of ILEs present in a set of orthologs are
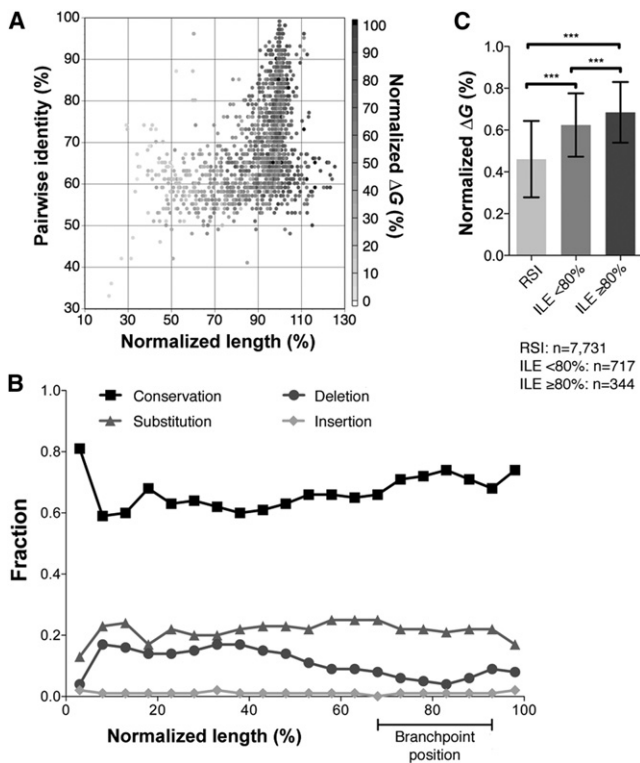
Figure 3. Degeneration of Introner-like Elements

(A) Correlation between pairwise identity and normalized length of all identified fungal ILEs. The gray scale indicates normalized Gibbs free energy ($\Delta G$).

(B) Conserved, substituted, deleted, and inserted nucleotide positions within ILEs. ILE length was expressed as percentage and arranged in 5% bins. Numbers of conserved positions, substitutions, deletions, and insertions were counted in each bin and expressed as a fraction. Distribution is shown for all ILEs from *Cladosporium fulvum*, *Dothistroma septosporum*, *Mycosphaerella graminicola*, and *Mycosphaerella fijiensis*.

(C) Mean and SD of normalized $\Delta G$ values of all RSIs, ILEs with < 80%, and ILEs with $\geq$ 80% pairwise identity from all six fungi. A non-parametric Kruskall-Wallis test was carried out (***p < 0.0001), followed by a Dunn's pairwise comparison test at $\alpha$ = 0.05 significance level. See also Figure S3.

associated with gains in both *C. fulvum* and *D. septosporum* (Table S3). Up to 351 single gains were inferred in their common ancestor, of which ILEs only represent 10%. Similar numbers of single gains were not only found in *M. graminicola* and *M. fijiensis* but also in *S. musiva* that only contains highly degenerated ILEs (Table S3). From these results, we conclude that due to ILE degeneration, it is essential to compare species with shorter evolutionary distances in order to confidently estimate the contribution of ILEs to intron gains.

### Newborn ILEs Become Undistinguishable from RSIs within 100,000 Years

The genomes of *M. graminicola* (three isolates), two sister species S1 and S2 (five and four isolates, respectively), and *Septoria passerinii* (one isolate) were used because divergence between *M. graminicola*, S1 and S2 was dated to 11 and 22.3 thousand years, respectively [19]. *C. fulvum* and *D. septosporum* were included as an outgroup and, in contrast to the recent report on extensive intron gains in *M. graminicola* [6], polymorphic positions that mainly represent segregating introns within populations were excluded from the analysis. This new data set confirmed that 50% of the single gains in

*C. fulvum* and *D. septosporum* originate from ILEs (Figure 5C). ILE contribution to intron gains reaches 90% in *M. graminicola* and S2, which diverged more recently than *C. fulvum* and *D. septosporum*, but only 40% of the single gains in *S. passerinii* and the common ancestor of *M. graminicola*, S1, and S2, could be ascribed to distinguishable ILEs (Figure 5C; Table S4). ILE contribution to single gains even drops to 6%–10% in older ancestors. These observations support the hypothesis that many or even potentially all RSIs are degenerated ILEs. Rough dating of the species divergences suggests that extensive intron gains have occurred in all these species during the last 100 thousand years. The dating also shows that all six *M. graminicola* ILE clusters must have multiplied more than 22 thousand years ago, but members of clusters mg05 and mg06 are no longer active (Figure 5D). Certainly, degenerated ILE clusters identified in *D. septosporum* and *S. musiva* lost their ability to multiply a long time ago. Additionally, over this short time frame, average pairwise identity and average length of ILEs have decreased (Figure 5E), showing that ILE clusters can emerge and successfully multiply, but multiplication may also stop due to rapid degeneration. As such, ILE identification becomes difficult in organisms with short generation times in which the last multiplication event has occurred more than 100 thousand years (according to the dated species tree).

### Conclusions

Based on this study, we conclude that ILEs account for most of the intron gains in at least six fungal species. Several mechanisms have been proposed for intron gains, from intron transposition to intronization of exons [7]. Most, however, are supported by few experimental data [7] and none can convincingly explain ILE multiplication. The contribution of these mechanisms to intron gains appears limited in comparison to the ILE multiplication reported here. A recently proposed mechanism involves reverse splicing of introns directly into the genome followed by reverse transcription [3]. This hypothetical mechanism requires fewer steps than the hitherto accepted mechanism of intron transposition but requires RNA-DNA hybridization. This mechanism could apply to ILEs because new ILEs perfectly insert into genes (no sequence deletion or duplication), and it would not involve homologous recombination. Such a step is required in the other proposed mechanisms although it occurs at low frequency in filamentous fungi [20]. ILEs multiplication might be linked to transcription because they are always found on the coding strand of genes as reported for introner elements in *Micromonas* [8].

Our findings show that introns in fungi are highly dynamic because only 51% of introns are conserved from *S. passerinii* to *C. fulvum* (Figure 5C). Because up to 90% of intron gains in the fungal species included in this study originate from ILEs, we propose that ILEs, too degenerated for detection, represent the species-specific introns in other fungi. Introner elements in *Micromonas* extensively multiplied in thousands of copies, which suggests that intron multiplication could also occur outside the fungal kingdom [8]. However, introner elements differ from ILEs because they lack predicted stable secondary structures. We could not find ILEs in intronomes of other eukaryotic lineages that contain near-identical introns such as *O. dioica* [15]. However, we speculate that active ILEs might have appeared very early in eukaryotes evolution but have become indistinguishable from RSIs. It is also possible that similar ILEs are currently spreading in other not yet studied eukaryotes. Further studies on ILEs are required to increase our understanding of eukaryotic gene structure evolution.
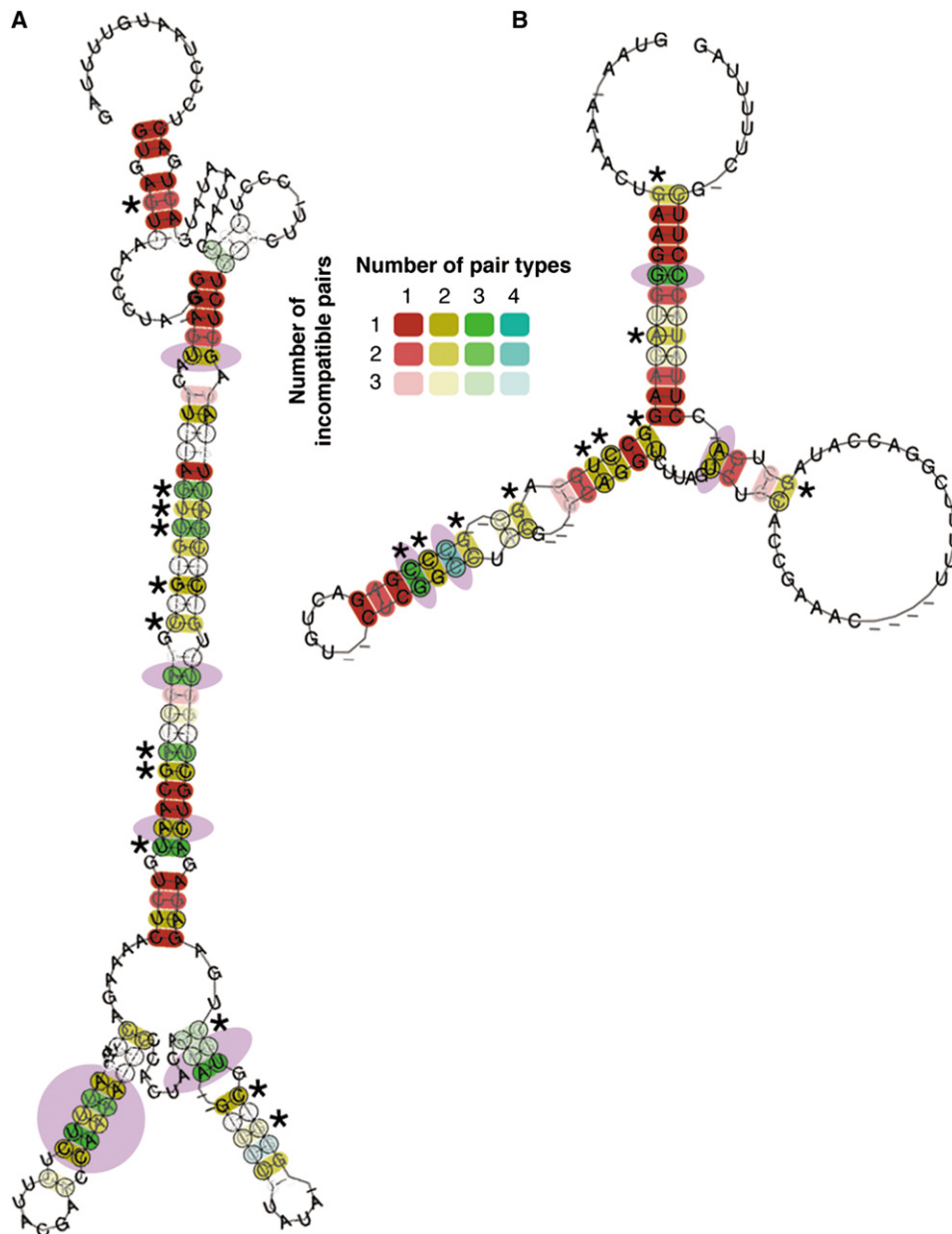
**Figure 4. Predicted Secondary Structure of Introner-like Elements**

(A) Alignment-based predicted secondary structure of ILEs from clusters cf03 and ds04 of *Cladosporium fulvum* and *Dothistroma septosporum*, respectively. Stars indicate pairs that involve a G-U pair in at least one sequence of the alignment. Pink circles highlight pairs that contain compensatory mutations. The color scale indicates the number of compatible pair types (C-G, G-C, A-U, U-A, G-U, or U-G). The saturation decreases with the number of incompatible base pairs.

(B) Alignment-based predicted secondary structure of HMM consensus sequences of clusters mg01, mg02, mg03, mg04, and mg06 from *Mycosphaerella graminicola*. See also Figure S4.

**References**

1. Koonin, E.V. (2006). The origin of introns and their role in eukaryogenesis: a compromise solution to the introns-early versus introns-late debate? Biol. Direct *1*, 22.
2. Rodríguez-Trelles, F., Tarrío, R., and Ayala, F.J. (2006). Origins and evolution of spliceosomal introns. Annu. Rev. Genet. *40*, 47–76.
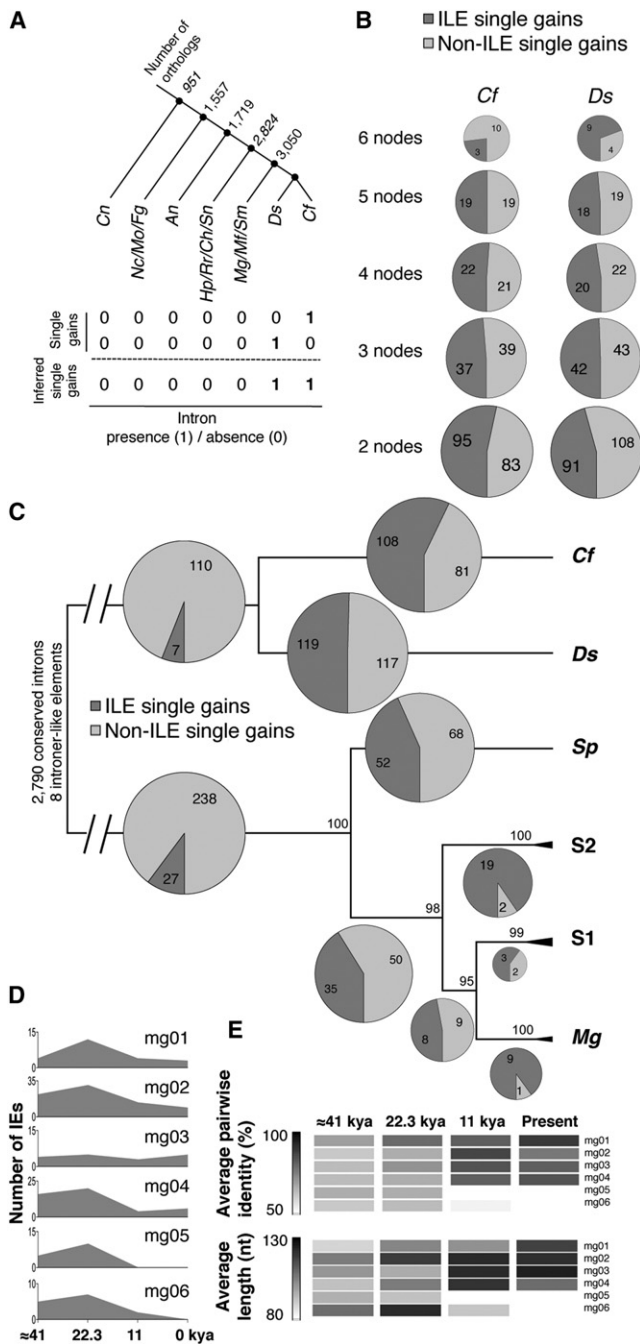
Figure 5. Contribution of Introner-like Elements to Single Intron Gains

(A) Species tree [21] and pattern used to assign single intron gains in Dothideomycetes.

(B) Contribution of ILEs and non-ILEs to single gains in *Cf* and *Ds* when including outgroups of 2 to 6 nodes. Size of the pies is proportional to the total number of single gains.

(C) Contribution of ILEs to single gains in *Mg*, S1, S2, *Sp*, *Cf*, and *Ds*. Inferred gains in the oldest ancestors were determined using other Dothideomycetes as outgroup. Size of the pies is proportional to the total number of single gains. The species tree was constructed using four genes (*TUB1*, *EIF3b*, *PAP1*, and *RPS9*). Numbers at the nodes indicate bootstrap values of 500 replicates. The tree was rooted according to Wang et al. [21]. The molecular clock was calibrated using the dating of S2 speciation [19].

(D) Number of ILEs per *Mg* cluster (mg01 to mg06) that are conserved in all species at each dated node. The dating on the x axis corresponds to the age of the nodes determined in the phylogenetic tree. Kya, thousand years ago.

(E) The average pairwise identity and average length of conserved ILEs is indicated per cluster for each dated node of the phylogenetic tree. *An*, *Aspergillus nidulans*; *Cf*, *Cladosporium fulvum*; *Ch*, *Cochliobolus heterostrophus*; *Cn*, *Cryptococcus neoformans*; *Ds*, *Dothistroma septosporum*; *Fg*, *Fusarium graminearum*; *Hp*, *Hysterium pulicare*; *Mf*, *Mycosphaerella fijiensis*; *Mg*, *Mycosphaerella graminicola*; *Mo*, *Magnaporthe oryzae*; *Nc*, *Neurospora crassa*; *Rr*, *Rhytidhysteron rufulum*; *Sm*, *Septoria musiva*; *Sn*, *Stagonospora nodorum*; *Sp*, *Septoria passerini*. See also Tables S3 and S4.

3. Roy, S.W., and Irimia, M. (2009). Mystery of intron gain: new data and new models. Trends Genet. *25*, 67–73.

4. Csuros, M., Rogozin, I.B., and Koonin, E.V. (2011). A detailed history of intron-rich eukaryotic ancestors inferred from a global survey of 100 complete genomes. PLoS Comput. Biol. *7*, e1002150.

5. Li, W., Tucker, A.E., Sung, W., Thomas, W.K., and Lynch, M. (2009). Extensive, recent intron gains in *Daphnia* populations. Science *326*, 1260–1262.

6. Torriani, S.F., Stukenbrock, E.H., Brunner, P.C., McDonald, B.A., and Croll, D. (2011). Evidence for extensive recent intron transposition in closely related fungi. Curr. Biol. *21*, 2017–2022.

7. Yenerall, P., Krupa, B., and Zhou, L. (2011). Mechanisms of intron gain and loss in *Drosophila*. BMC Evol. Biol. *11*, 364.

8. Worden, A.Z., Lee, J.H., Mock, T., Rouzé, P., Simmons, M.P., Aerts, A.L., Allen, A.E., Cuvelier, M.L., Derelle, E., Everett, M.V., et al. (2009). Green evolution and dynamic adaptations revealed by genomes of the marine picoeukaryotes Micromonas. Science *324*, 268–272.

9. Kupfer, D.M., Drabenstot, S.D., Buchanan, K.L., Lai, H., Zhu, H., Dyer, D.W., Roe, B.A., and Murphy, J.W. (2004). Introns and splicing elements of five diverse fungi. Eukaryot. Cell *3*, 1088–1100.

10. Schwartz, S.H., Silva, J., Burstein, D., Pupko, T., Eyras, E., and Ast, G. (2008). Large-scale comparative analysis of splicing signals and their corresponding splicing factors in eukaryotes. Genome Res. *18*, 88–103.

11. Le Hir, H., Nott, A., and Moore, M.J. (2003). How introns influence and enhance eukaryotic gene expression. Trends Biochem. Sci. *28*, 215–220.

12. Warf, M.B., and Berglund, J.A. (2010). Role of RNA structure in regulating pre-mRNA splicing. Trends Biochem. Sci. *35*, 169–178.

13. Roy, S.W., and Irimia, M. (2009). Splicing in the eukaryotic ancestor: form, function and dysfunction. Trends Ecol. Evol. (Amst.) *24*, 447–455.

14. Roy, S.W., and Gilbert, W. (2005). Rates of intron loss and gain: implications for early eukaryotic evolution. Proc. Natl. Acad. Sci. USA *102*, 5773–5778.

15. Denoeud, F., Henriet, S., Mungpakdee, S., Aury, J.M., Da Silva, C., Brinkmann, H., Mikhaleva, J., Olsen, L.C., Jubin, C., Cañestro, C., et al. (2010). Plasticity of animal genome architecture unmasked by rapid evolution of a pelagic tunicate. Science *330*, 1381–1385.

16. Qiu, W.G., Schisler, N., and Stoltzfus, A. (2004). The evolutionary gain of spliceosomal introns: sequence and phase preferences. Mol. Biol. Evol. *21*, 1252–1263.

17. Nielsen, C.B., Friedman, B., Birren, B., Burge, C.B., and Galagan, J.E. (2004). Patterns of intron gain and loss in fungi. PLoS Biol. *2*, e422.

18. Mathews, D.H., Moss, W.N., and Turner, D.H. (2010). Folding and finding RNA secondary structure. Cold Spring Harb Perspect Biol *2*, a003665.

19. Stukenbrock, E.H., Bataillon, T., Dutheil, J.Y., Hansen, T.T., Li, R., Zala, M., McDonald, B.A., Wang, J., and Schierup, M.H. (2011). The making of a new pathogen: insights from comparative population genomics of the domesticated wheat pathogen *Mycosphaerella graminicola* and its wild sister species. Genome Res. *21*, 2157–2166.

20. Weld, R.J., Plummer, K.M., Carpenter, M.A., and Ridgway, H.J. (2006). Approaches to functional genomics in filamentous fungi. Cell Res. *16*, 31–44.

21. Wang, H., Xu, Z., Gao, L., and Hao, B. (2009). A fungal phylogeny based on 82 complete genomes using the composition vector method. BMC Evol. Biol. *9*, 195.