

INFORMATION TO USERS

This manuscript has been reproduced from the microfilm master. UMI films the text directly from the original or copy submitted. Thus, some thesis and dissertation copies are in typewriter face, while others may be from any type of computer printer.

The quality of this reproduction is dependent upon the quality of the copy submitted. Broken or indistinct print, colored or poor quality illustrations and photographs, print bleedthrough, substandard margins, and improper alignment can adversely affect reproduction.

In the unlikely event that the author did not send UMI a complete manuscript and there are missing pages, these will be noted. Also, if unauthorized copyright material had to be removed, a note will indicate the deletion.

Oversize materials (e.g., maps, drawings, charts) are reproduced by sectioning the original, beginning at the upper left-hand corner and continuing from left to right in equal sections with small overlaps.

Photographs included in the original manuscript have been reproduced xerographically in this copy. Higher quality 6" x 9" black and white photographic prints are available for any photographs or illustrations appearing in this copy for an additional charge. Contact UMI directly to order.

Bell & Howell Information and Learning
300 North Zeeb Road, Ann Arbor, MI 48106-1346 USA
800-521-0600

UMI[®]



Université d'Ottawa • University of Ottawa

**THE EVOLUTION OF LAND PLANTS INFERRED USING DOMAINS D-F OF THE
LARGEST SUBUNIT OF RNA POLYMERASE II**

by

SHEHRE-BANOO MALIK

Thesis submitted to the
School of Graduate Studies and Research
University of Ottawa
in partial fulfilment of the
MSc degree
at the
Ottawa-Carleton Institute of Biology

Ottawa, Ontario, Canada

12th January 2000

© 2000 Shehre-Banoo Malik



National Library
of Canada

Acquisitions and
Bibliographic Services

395 Wellington Street
Ottawa ON K1A 0N4
Canada

Bibliothèque nationale
du Canada

Acquisitions et
services bibliographiques

395, rue Wellington
Ottawa ON K1A 0N4
Canada

Your file *Votre référence*

Our file *Notre référence*

The author has granted a non-exclusive licence allowing the National Library of Canada to reproduce, loan, distribute or sell copies of this thesis in microform, paper or electronic formats.

The author retains ownership of the copyright in this thesis. Neither the thesis nor substantial extracts from it may be printed or otherwise reproduced without the author's permission.

L'auteur a accordé une licence non exclusive permettant à la Bibliothèque nationale du Canada de reproduire, prêter, distribuer ou vendre des copies de cette thèse sous la forme de microfiche/film, de reproduction sur papier ou sur format électronique.

L'auteur conserve la propriété du droit d'auteur qui protège cette thèse. Ni la thèse ni des extraits substantiels de celle-ci ne doivent être imprimés ou autrement reproduits sans son autorisation.

0-612-48169-7

Canada

ACKNOWLEDGEMENTS:

I would like to thank Dr. Guy Drouin for giving me the opportunity and major guidance for this project. Guy, Dr. Donal Hickey, Dr. Lynn Gillespie and Dr. George Carmody -- members of my supervisory committee -- are thanked for all of their constructive advice throughout this project. Many thanks to Michael Ell for taking the time to guide me around. Dr. Andrew J. Roger, Dr. Sandie L. Baldauf, Dr. Hans-Peter Klenk, Dr. Ben D. Hall, Betty McConaughy, Dr. John W. Stiller, Tim Chipman, Dr. Peter G. Foster, Dr. Selvi Subramanian and Dr. John M. Logsdon Jr. are thanked for helpful discussions and practical advice regarding methodology. J. Stiller and J. Logsdon also kindly shared the RNA polymerase alignments from their latest publications. The *Giardia* RPB1 sequence is from H.-P. Klenk. Thanks to Marc Carrier for help with sequencing and for his dedicated work. Thanks to Mehrdad Babaei for providing two *Ephedra viridis* RPB1 (D-F) clones. Many thanks to Guy for the Southern analysis of *Magnolia* and *Nymphaea*. Cathy, Ed, Lia, Patrick, Gitanjali and Dave are thanked for their guidance and friendship. Special thanks also go to Louise Cossette, Lise Maisonneuve, Caroline Pharand, Annie Robillard and the department's technical staff for always being pleasantly helpful. Dr. David Currie, Dr. Jim Fenwick, Dr. Linda Bonen and Dr. Scott Findlay are thanked for their constructive remarks. I would like to thank each of my labmates for their patience, friendship, teamwork and for the pleasure of their company.

I am indebted to my mother and grandmother, for always stressing the importance of a good education and a decent life, and to my young brothers for filling us with hope and youth. I have not achieved anything without their endless love, perseverance and sense of humor.☺

TABLE OF CONTENTS

| <u>Description</u> | <u>Page</u> |
|---|-------------|
| ACKNOWLEDGEMENTS | i |
| TABLE OF CONTENTS | ii |
| ABSTRACT | 1 |
| LIST OF FIGURES & APPENDICES | 2 |
| LIST OF TABLES | 4 |
| INTRODUCTION | 5 |
| MATERIALS AND METHODS | 14 |
| <i>Plant Material</i> | 14 |
| <i>Plant DNA and RNA Purification</i> | 14 |
| <i>Amplification of RPBI</i> | 16 |
| <i>Cloning</i> | 20 |
| <i>Sequencing</i> | 23 |
| <i>Northern and Southern Hybridizations</i> | 24 |
| <i>Sequence Analysis</i> | 25 |
| RESULTS | 31 |
| <i>Preliminary Analysis of the Region Between Domains F and H of RPBI</i> | 31 |
| <i>Plant RPBI</i> | 37 |
| <i>Phylogenetic Analysis of Plant RPBI D-F DNA Sequences</i> | 45 |
| <i>Phylogenetic Analysis of Plant RPBI D-F Amino Acid Sequences</i> | 50 |
| <i>RPBI in Other Eukaryotes</i> | 63 |
| <i>Compositional Bias</i> | 71 |
| <i>Substitution Analysis of Region D-F of RPBI in Land Plants</i> | 76 |
| <i>Relative Rate Tests</i> | 81 |
| <i>Copy Number of Plant RPBIs</i> | 83 |
| DISCUSSION | 87 |
| <i>Phylogenetic Analysis of DNA Sequences of Regions D-F of RPBI in Plants</i> | 87 |
| <i>Phylogenetic Analysis of Amino Acid Sequences of Regions D-F of RPBI in Plants</i> | 88 |
| <i>Substitution Analysis</i> | 93 |
| <i>Relative Rate Tests</i> | 93 |
| <i>Phylogenetic Analyses of Regions D-F of RPBI in Other Eukaryotes</i> | 94 |
| <i>Conclusions</i> | 100 |
| REFERENCES | 101 |

ABSTRACT

The aim of this project has been to develop the use of the nucleotide and amino acid sequences of the largest subunit of RNA polymerase II (RPB1) to study the evolutionary relationships of diverse land plants. 1.3 kb of this gene, which includes 900bp of coding sequence corresponding to regions D-F of the protein, have been cloned and sequenced from genomic DNA of seven plant species, one from each of seven divisions. 900 bp of cDNA sequence of regions D-F of RPB1 from *Zea mays* revealed identity with the genomic sequence. Phylogenetic analysis of these nucleotide and corresponding amino acid sequences by parsimony, neighbor-joining and maximum likelihood methods revealed inconsistent results. The maximum-likelihood method presented the best resolved trees. It is unlikely that the inconsistencies are the result of compositional biases, multiple substitutions, or unequal substitution rates. Analysis of amino acid sequences of regions D-F and F-H of RPB1 in animals, plants and fungi showed that animals are more closely related to plants than to fungi, contrary to recent literature. In comparing regions D-F of the largest subunits of RNA polymerases I, II and III (RPA1, RPB1 and RPC1) of all eukaryotes, though, the evolutionary relationships of the major eukaryotic kingdoms were not resolved but RPA1 was found to be more closely related to RPB1 than to RPC1. Since analyses in the literature of longer RPB1 sequences from fewer taxa did clearly resolve relationships of major eukaryotic lineages, this study shows that region D-F of this gene is too short to be informative. Due to the conserved nature of this partial sequence, it appears that a longer sequence from this gene with more synapomorphic sites is required to better assess the evolutionary relationships of plants.

LIST OF FIGURES AND APPENDICES

| <u>Description</u> | <u>Page</u> |
|--|-------------|
| Figure 1: Schematic diagram of domains of RPB1 and PCR primer locations | 17 |
| Figure 2: A 606-position amino acid alignment of 13 sequences from RPB1 F-H, with the slime mold <i>Dictyostelium discoideum</i> as the outgroup to Animals, Plants and Fungi | 32 |
| Figure 3: Maximum-likelihood distance tree of 13 amino acid sequences of RPB1 regions F-H from animals, plants and fungi | 36 |
| Figure 4: RPB1 D-F PCR products from genomic and maize cDNA templates | 38 |
| Figure 5: An 834 bp alignment of the RPB1 D-F coding sequence of plants | 40 |
| Figure 6: Two equally parsimonious phylogenetic trees of nucleotide sequences encoding plant RPB1 regions D-F. | 47 |
| Figure 7: Neighbor-joining phylogenetic tree of nucleotide sequences encoding RPB1 D-F in plants. | 48 |
| Figure 8: Maximum-likelihood phylogenetic tree of nucleotide sequences encoding RPB1 D-F in plants. | 49 |
| Figure 9: 557 position amino acid alignment of 38 D-F sequences from type A RNA polymerases. | 51 |
| Figure 10: Four equally parsimonious phylogenetic trees inferred from a 280-position alignment of plant RPB1 D-F amino acid sequences. | 57 |
| Figure 11: Neighbor-joining tree inferred from plant RPB1 D-F amino acid sequences. | 58 |
| Figure 12 (A, B) Maximum-likelihood analysis of plant RPB1 D-F amino acid sequences using the Blosum62 model of substitution. | 61 |
| Figure 13 (A, B) Maximum-likelihood analysis of plant RPB1 D-F amino acid sequences using the JTT model of substitution. | 62 |
| Figure 14: One of three equally parsimonious trees inferred from animal, plant and fungal amino acid sequences for regions D-F of RPB1. | 66 |
| Figure 15: Neighbor-joining phylogenetic tree inferred from animal, plant and fungal amino acid sequences for regions D-F of RPB1. | 67 |
| Figure 16: Maximum-likelihood phylogenetic tree inferred from animal, plant and fungal amino acid sequences for regions D-F of RPB1. | 68 |
| Figure 17: Maximum-likelihood phylogenetic tree inferred from amino acid sequences for domains D-F of eukaryotic type A RNA polymerases using the JTT model of substitution. | 69 |
| Figure 18: Maximum-likelihood phylogenetic tree inferred from amino acid sequences for domains D-F of eukaryotic type A RNA polymerases using the Blosum62 model of substitution. | 70 |
| Figure 19: The effects of G+C compositional bias on the three codon positions of regions D-F of RPB1 in (A) plants and in (B) diverse eukaryotes. | 72 |

| <u>Description</u> | <u>Page</u> |
|--|-------------|
| Figure 20: G+C content at the third base of codons reflects codon usage patterns in regions D-F of RPB1 in (A) plants and in (B) diverse eukaryotes. | 73 |
| Figure 21: Variation in G+C content among plant RPB1 sequences corresponding to regions D-F. | 74 |
| Figure 22: The effect of G+C content on the amino acid composition of regions D-F of RPB1 in (A) plants and in (B) diverse eukaryotes. | 75 |
| Figure 23: Comparison of mean nonsynonymous substitutions per site in regions D-F of RPB1 with divergence times of some plant groups. | 80 |
| Figure 24: Southern blots of (A) <i>Magnolia soulangeana</i> and (B) <i>Nymphaea odorata</i> total genomic DNA digests probed with homologous clones of regions D-F of RPB1. | 85 |
| Appendix 1: 556 bp alignment of 9 nucleotide sequences corresponding to the first and second bases of codons in the region of the RPB1 gene of plants which encodes domains D-F, which was used to infer the phylogenetic trees shown in figures 6-8. | 111 |
| Appendix 2: 278 position alignment of 9 amino acid sequences of regions D-F of RPB1 in land plants that was used to infer the phylogenetic trees shown in figures 10-13. | 113 |
| Appendix 3: 282 position alignment of 19 amino acid sequences of regions D-F of RPB1 in animals, plants and fungi that was used to infer the phylogenetic trees shown in figures 14-16. | 114 |
| Appendix 4: 252 position alignment of 38 amino acid sequences of regions D-F of RPA1, RPB1 and RPC1 from diverse eukaryotes that was used to infer the phylogenetic trees shown in figures 17 and 18. | 116 |
| Appendix 5: Results of this study of regions D-F of RPB1 in plants compared with the most recently published studies, which were based upon larger datasets and found that <i>Nymphaea</i> , a herbaceous dicot, was amongst the earliest extant angiosperms. | 119 |

LIST OF TABLES

| <u>Description</u> | <u>Page</u> |
|---|-------------|
| Table 1: Plants from which RPBI regions D-F were PCR amplified for this analysis. | 15 |
| Table 2: Alignments corresponding to domains D and F of RPBI, including chloroplast, α -proteobacterial, archaeal and RPA1 and RPC1 homologues. | 18 |
| Table 3: Codes and sources for sequences included in phylogenetic analyses of type A RNA polymerases. | 26 |
| Table 4: Identification of RPBI PCR products from various plants. | 39 |
| Table 5: The weighted average of nonsynonymous substitutions per site (K_a). | 77 |
| Table 6: The weighted average of synonymous substitutions per site (K_s). | 77 |
| Table 7: Estimate of the number of multiple substitutions at nonsynonymous sites in plant RPBI D-F sequences. | 78 |
| Table 8: Estimate of the number of multiple substitutions at synonymous sites in plant RPBI D-F sequences. | 79 |
| Table 9: Results of relative rate tests of substitution rates of amino acid sequences of RPBI regions D-F amongst plants. | 81 |
| Table 10: Results of relative rate tests of substitution rates of amino acid sequences of RPBI regions D-F between animals, plants and fungi. | 82 |
| Table 11: Results of relative rate tests of substitution rates of amino acid sequences of RPBI regions F-H between animals, plants and fungi. | 82 |
| Table 12: Results of relative rate tests of substitution rates of amino acid sequences of regions D-F of RPA1, RPBI and RPC1 in eukaryotes. | 82 |
| Table 13: Ploidy of representatives of some major plant phyla. | 86 |

INTRODUCTION

The phylogenetic relationships of the major groups of land plants, as well as of the basal clades of angiosperms, continually confound plant systematists. The rapid, ancient divergence of these taxa and the lack of sufficient extant intermediate taxa are major contributors to this systematic problem. While phylogenetic analyses of morphological features include data from the fossil record, these analyses lack the sensitivity afforded by the phylogenetic analysis of molecular data. Integration of both types of datasets is helpful, but they also conflict in some cases.

While both morphological and molecular analyses have their merits and pitfalls, molecules offer more characters with discretely different states, and increasingly stringent analyses of molecular data are being developed. Unlike morphological datasets, molecular datasets do not yet include data from fossils. Fossilization is rare, though, and even fossils do not yet provide enough intermediates to extant plant lineages (Ridley, 1996). Molecular datasets do not all evolve at the same rate, and sometimes they include enough compositional bias to misgroup taxa (Li and Graur, 1991; Foster and Hickey, 1999). For example, ribosomal RNA genes are affected more by compositional bias and heterogeneous rates of evolution than protein-coding genes are and such genes also evolve at different rates whether they are encoded in the nucleus, chloroplasts, or mitochondria (Hasegawa and Hashimoto, 1993; Li and Graur, 1991). In plants, each of these genomes is subject to different types or degrees of rearrangements. Even within these contexts, different genes also evolve at different rates, so it is important to find molecular phylogenetic markers which accurately represent the phylogeny of the taxa in question, with significant statistical support. Thus, it is useful to seek several independent datasets to corroborate with previous phylogenetic inferences and make our understanding of land plant phylogeny more precise.

Phylogenetic analysis of morphological data in plants falls into several categories: classification according to shared biochemical and cellular properties of extant green algae and plants, comparisons of features of extant and fossilized pollen, and comparisons based on developmental features and ultrastructure of extant plants and fossils when possible (Bremer, 1985; Crane *et al.*, 1995; Kenrick and Crane, 1997; Mishler and Churchill, 1985). It was difficult to perform early such analyses due to a lack of clearly defined analytical criteria and raw data (Meacham, 1994; Takhtadzhian, 1997). In the past two decades, however, the accumulation of more fossil evidence and technological advances in the computation of phylogenetic relationships have clarified the relationships of land plants and angiosperms more than ever (Bateman *et al.*, 1998; Crane *et al.*, 1995; Donoghue, 1994; Kenrick and Crane, 1997). Based upon this information, it is clear that land plants arose from green algae, and that the earliest land plants were bryophytes, from which vascular plants arose as a monophyletic group (Bateman *et al.*, 1998 and references therein; Donoghue, 1994). This data also makes it clear that angiosperms are monophyletic and are composed of Eudicots, which have triaperturate pollen, and the Monocots and Magnoliid dicots, which have monocolpate pollen (Bateman *et al.*, 1998, Crane *et al.*, 1995; Taylor and Hickey, 1996; Loconte, 1996). However, the relationships amongst the early land plants are unclear, as are the relationships of the basal clades of angiosperms (Bateman *et al.*, 1998; Donoghue, 1994). The relationships are resolved differently depending upon the combinations of characters used, the equal or uneven distribution of taxa sampled, and whether maximum parsimony, maximum likelihood, or neighbour-joining methods of phylogenetic inference are used (Bateman *et al.*, 1998; Taylor and Hickey, 1996; Takhtadzhian, 1997; Meacham, 1994; Donoghue, 1994).

This data supports two opposing hypotheses for the origin of angiosperms. The first is that of a herbaceous origin of angiosperms, meaning that the Nymphaeales are the most basal angiosperm clade and monocots and eudicots such as *Arabidopsis* are more derived. The second hypothesis is that woody Magnoliids are most basal amongst angiosperms, and that they share some monocolpate pollen features with gymnosperms and others with monocots. Clearly, the abundance of homoplasy and the shortage of synapomorphies when comparing such diverse taxa as angiosperms and other land plants is a problem, as with any other anciently diversified group. This problem is hampered by the unequal rates of evolution of phenotypic characters in these groups, gaps in the fossil record and unbalanced and unequal sampling of extant taxa. A proper understanding of plant phylogeny requires that these gaps are filled.

Similar problems of homoplasy and unequal rates of evolution are faced by molecular systematists. The Siluro-Devonian radiation of land plants, approximately 480-450 million years ago (MYA), is estimated to have occurred over 35 - 100 million years, which is very rapid in comparison to the amount of time that has passed between then and now (Bateman *et al.* 1998; Kenrick and Crane, 1997; Stewart and Rothwell, 1993). Such relationships are difficult to measure since molecules (or characters) evolving fast enough to bear a record of the rapid radiation have been saturated with homoplasy by now, and molecules or characters evolving slowly enough to be less homoplastic likely evolved too slowly to be informative of the rapid radiation. An assemblage of molecular data could overcome problems of unequal rates of evolution and homoplasy, since different molecular markers and corresponding models of evolution could be applied to reconstruct different parts of the plant phylogenetic tree (Ritland and Eckenwalder, 1992). Donoghue (1994) pointed out that from 1989-1991, only about 0.6% of phylogenetic analyses considered molecular

data for analysis of the relations of the major groups of vascular plants, land plants or seed plants. Thus, it is imperative to assemble an integrated approach to the molecular systematics of plants, incorporating data from multiple loci and from large-scale genome characteristics such as inversions, insertions/deletions, and introns (Bateman *et al.* 1998; Qiu and Palmer, 1999).

Several molecular datasets from each of the three genomes (nuclear, chloroplast and mitochondrial) have been examined to try to solve the relationships of plants, but conflicts remain. The bulk of molecular phylogenetic analysis of plants lies in the examination of nuclear 18S ribosomal RNA genes from 223 plants (Soltis *et al.*, 1997) and of the chloroplast-encoded gene for the largest subunit of ribulose-1,5-bisphosphate carboxylase (*rbcL*) in 499 plants (Chase *et al.*, 1993), almost entirely from angiosperm species. Chaw *et al.* (1997) examined 18S rRNA in 65 gymnosperms and angiosperms. Nuclear 26S rRNA was examined in 15 angiosperms by Kuzoff *et al.* (1998). Hasebe *et al.* (1992, 1993) examined the phylogeny of 10-15 seed plants and ferns using *rbcL*. Mitochondrial introns and RNA editing have been examined in three genes from 28 plant species spanning the range of all plants (Malek *et al.*, 1996; Qiu *et al.*, 1998; Qiu and Palmer, 1999; Cho *et al.*, 1998). Multigene families of nuclear genes encoding actin, phosphoenol pyruvate carboxylase (PEPC), floral homeotic MADS-box genes and phytochromes have also been examined in up to 20 diverse land plant taxa (Moniz de Sá and Drouin, 1996 and unpublished data; An *et al.*, 1999; Gehrig *et al.*, 1998; Purugganan, 1997; Kolukisaoglu *et al.*, 1995; Donoghue and Matthews, 1998). The second-largest subunit of RNA polymerase II (RPB2), a single-copy ubiquitous nuclear-encoded protein-coding gene, was also recently used to infer the phylogeny of 12 green plants by Denton *et al.* (1998). These are some of the most notable recent examples of molecular markers

applied to resolve the deep phylogenetic relationships amongst green plants and amongst the basal clades of angiosperms.

While molecular data support the origin of land plants from Charophycean algae and the monophyly of land plants, of vascular plants, of Gnetales and of angiosperms, some ambiguities still lie in the early relations of land plants, amongst gymnosperm taxa, and amongst the basal angiosperm clades. The analyses which included the broadest range of taxa were made with molecular markers which evolve slowly, such that they are more suited to studying more recent evolutionary events amongst plants. Other molecular analyses (actin, PEPC, MADS-box genes, phytochrome genes, RPB2) lack the breadth of taxa included in the rRNA and *rbcL* analyses. It is possible that orthologous members of the multigene families may clarify the phylogeny of subsets of plants, by rooting such genes from the ingroup with their latest paralogue in the outgroup (Donoghue and Matthews, 1998; Brown and Doolittle, 1995). In comparison to angiosperms, gymnosperms are underrepresented by many of these analyses, but not as underrepresented as early nonvascular and vascular plants are. While the absence of certain mitochondrial introns supports that liverworts such as *Marchantia* are the earliest land plants (Qiu *et al.*, 1998) the authors also emphasize the need to combine data from multiple genes subject to different evolutionary pressures to reduce the influence of homoplasy (Qiu and Palmer, 1999). Also, it is important that all of these datasets be analysed in a manner that minimizes the effects of GC content bias, saturation of synonymous substitutions, and heterogeneous rates of evolution amongst the ingroups, which can result in artefactual phylogenies.

The largest datasets amongst these analyses are those of *rbcL* and 18S rRNA. The *rbcL* dataset supports monophyly of Gnetales, of angiosperms and of monocots and a herbaceous origin

of angiosperms, with monocots either basal to or amongst the paleoherbs (such as *Nymphaea*), depending on the methods of phylogenetic tree inference used (Chase *et al.*, 1993). The 18S rRNA dataset also supports the monophyly of Gnetales, of angiosperms and of monocots, but does not clearly resolve any of the relationships of the major angiosperm clades (Soltis *et al.*, 1997). In parsimony analyses of rRNA data, a Magnoliid origin of angiosperms is only one step less likely than a herbaceous origin is (Nickrent and Soltis, 1995).

Only two of the datasets briefly described above include large protein-coding genes. The RPB2 and phytochrome datasets are 3564 bp and 3264 bp in length respectively, including the third base of codons (Denion *et al.*, 1998; Donoghue and Matthews, 1998). Protein-coding genes are phylogenetically useful since most of the GC content bias, if any, is absorbed in the synonymous sites of codons, which can be removed from the datasets used for phylogenetic inference. Also, alignments made at the amino acid level and then fitted to corresponding nucleotide sequences can be made with less ambiguity than with nucleotides alone when studying distantly related taxa (Sidow and Wilson, 1990). Large protein-coding genes are also phylogenetically useful since they provide more informative sites for analysis. To date, analyses of full-length datasets of these genes which include diverse gymnosperm, vascular and nonvascular plant taxa have not been shown. The largest subunit of RNA polymerase II, RPB1, has approximately 5 kb of protein-coding sequence in plants (Nawrath *et al.*, 1990). The only full-length RPB1 gene available from plants is from *Arabidopsis*. Since RPB1 is a key component of the RNA polymerase II holoenzyme in eukaryotes which transcribes messenger RNA and small nuclear RNAs, it would be useful not only from the phylogenetic point of view but also from a biochemist's point of view to further characterize this gene and protein in land plants.

RPB1 has recently been used to infer phylogenetic relationships of deeply diverging groups of eukaryotes with better statistical support than with data from rRNA, elongation factor (EF-1 α , EF-2) and heat shock protein (Hsp70) sequences (Hirt *et al.*, 1999; Stiller *et al.*, 1998; Stiller and Hall, 1997 and 1998). The Archaeal homologue to the largest subunit of eukaryotic RNA polymerases I, II and III, called the A' subunit, also provided the data for a better-resolved inference of Archaeal phylogenetic relationships than the rRNA dataset had (Puhler *et al.*, 1989a; Iwabe *et al.*, 1991; Klenk and Zillig, 1994). Partial RPB1 sequences have also been used recently to infer the phylogeny of invertebrates, arthropods, nematodes, and *Leishmania* species (Sidow and Thomas, 1994; Regier and Schultz, 1997; Baldwin *et al.*, 1997; Croan and Ellis, 1996). These data show that RPB1 is useful for studying deep-level and later levels of eukaryotic phylogeny. RPB1 is ubiquitous and usually found as a single-copy gene amongst eukaryotes. An RPB1 dataset can thus offer a phylogeny of orthologous molecules, encoded in the nuclear genome, with more informative sites than previous plant molecular datasets and less compositional bias than incurred by non-protein-coding sequences such as rRNA (Sidow and Wilson, 1990).

While the paucity of available plant RPB1 sequences and the homoplasy accumulated in other such ancient characters may discourage the use of RPB1 for inferring land plant phylogeny, the results of other phylogenetic analyses of RPB1 and its homologues are very encouraging. The largest subunits of eukaryotic RNA polymerases I, II and III (RPA1, RPB1 and RPC1), archaeal A' and A'' RNA polymerases, and bacterial β and β' RNA polymerases are collectively known as type A RNA polymerases due to their shared homology (Young, 1991). Their protein sequences share eight conserved domains, termed A-H (Young, 1991; Palenik, 1992). RPB1 is also characterized by carboxy-terminal repeats (YSPTSPX), considered as conserved domain "I", which are correlated

with the presence of spliceosomal introns in eukaryotic genes (Stiller and Hall, 1998). These conserved domains simplify the alignment of RNA polymerase genes from highly divergent taxa. The conserved domains are also very suitable for designing degenerate PCR primers for amplifying fragments of RPB1 from all eukaryotes. Sequences interspersed between domains A-I which are less conserved and are difficult to unambiguously align amongst deeply diverging protist lineages can, however, provide informative phylogenetic data when examining more closely related species.

Domains D and F of RPB1 provide a useful starting point for the examination of RPB1 in land plants. The conserved motif "PYNADFDGDEM^N" in domain D is involved with domain G in elongation and contains a DNA binding element of RPB1 which is found in all type A RNA polymerase homologues. The aspartic acid residues bind Mg²⁺. The nucleotide sequence for this motif in *Arabidopsis* RPB1 also lacks any introns, repeats or sixfold degenerate codons, is rich in twofold degenerate codons, and is long enough to provide an RNA polymerase-specific PCR primer with a stringent annealing temperature (i.e., high enough to avoid non-specific priming). Domain F, which is involved in elongation, contains a motif "VGQQNVEG" which represents the alpha-amanitin sensitive region which is a characteristic of RPB1. This motif, as stated, is not as highly conserved amongst all eukaryotes as is the domain D motif, but is highly conserved amongst crown eukaryotic lineages, and is the consensus of *Arabidopsis thaliana* and *Glycine max* plant RPB1 sequence for that domain. In these plants, the nucleotide sequence for this motif also lacks introns or sixfold degenerate codons and contains a few twofold degenerate codons, making it a stringent and useful candidate for a PCR primer as well. While these motifs only span 900 bp of coding sequence for RPB1 in *Arabidopsis*, that region should be easily obtained by PCR and provide a valuable starting point for further characterization of RPB1 in land plants.

On the basis of these characteristics of RPB1 and the successful use of complete and partial sequences of this gene in addressing a variety of phylogenetic questions, I hypothesize that it will be a useful tool for studying land plant phylogeny. In addition to this, characterization of RPB1 from more plant taxa will eventually lead to a better understanding of its fundamental biochemical properties in transcription.

MATERIALS AND METHODS

Plant Material:

Plant material for this study was collected from greenhouses at Carleton University and the University of Ottawa, the Rideau River and from a park in Cumberland, with the kind and generous assistance of Hughette Allard, Ed Bruggink and Elizabeth Thompson. Species names are as indicated in Table 1.

Plant DNA and RNA purification:

Plant total genomic DNA was prepared using the CTAB method of Doyle and Doyle (1990) and using the Qiagen Plant DNeasy Maxi and Mini Kits exactly as directed by the manufacturer. Two modifications to the CTAB method were that the extraction buffer included 2% polyvinylpyrrolidone and that if a thick, creamy interface persisted after the chloroform extractions, then a phenol-chloroform extraction was performed. Genomic DNA preparations from *Ephedra tweediana*, *Gnetum gnemon*, and *Welwitschia mirabilis* were kindly provided by Dr. Lynn Gillespie. CsCl-purified genomic DNA for *Zea mays* was obtained from Dr. Guy Drouin. These gifts of DNA as well as DNA prepared using the CTAB method from all of the other plant species included in this project were used as templates for PCR amplification of RPB1. While suitable as PCR templates, some of these DNA preparations still contained polysaccharides or secondary metabolites which interfered with restriction enzyme digestion, hindering Southern blot analysis. In such cases, either existing total DNA preparations were further purified using the Qiagen Plant DNeasy kit, or the same kit was used to prepare more DNA from fresh tissue.

Table 1: Plants from which RPBI D-F were PCR amplified for this analysis.

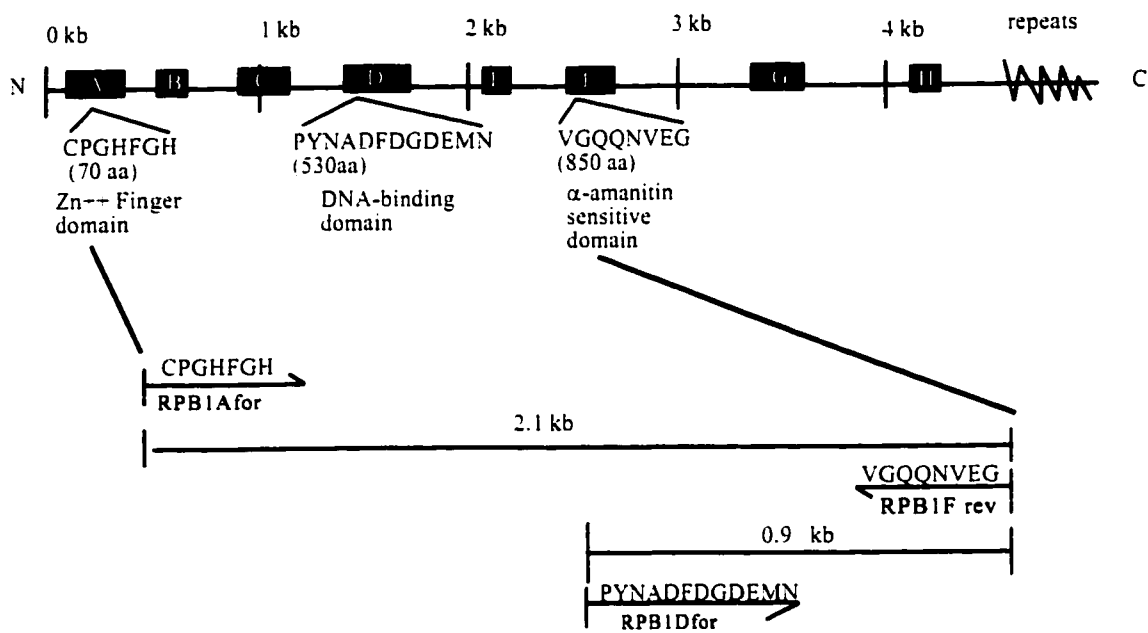
| <i>Taxa</i> | <i>Species Name</i> | <i>Common Name</i> |
|---------------------------|-------------------------------|---------------------|
| Coniferophyta | <i>Araucaria heterophylla</i> | Norfolk Island pine |
| Cycadophyta | <i>Cycas revoluta</i> | cycad |
| Gnetales, Ephedraceae | <i>Ephedra viridis</i> | Mormon tea |
| Gnetales, Ephedraceae | <i>Ephedra tweediana</i> | Mormon tea |
| Equisetophyta | <i>Equisetum hyemale</i> | horsetail |
| Ginkgophyta | <i>Ginkgo biloba</i> | gingko |
| Gnetales, Gnetaceae | <i>Gnetum gnemon</i> | gnetum |
| Magnoliales | <i>Magnolia soulangeana</i> | magnolia |
| Marchantiopsida | <i>Marchantia polymorpha</i> | liverwort |
| Nymphaeaceae | <i>Nymphaea odorata</i> | water lily |
| Filicophyta | <i>Osmunda claytonia</i> | fern |
| Psilotophyta | <i>Psilotum nudum</i> | whisk fern |
| Gnetales, Welwitschiaceae | <i>Welwitschia mirabilis</i> | welwitschia |
| Liliopsida (monocots) | <i>Zea mays</i> | maize/corn |

Plant total RNA was prepared according to the CTAB method of Chang *et al* (1993) and using the Qiagen Plant RNeasy Kit as directed by the manufacturer. RNA prepared using the CTAB method was successfully used as a template for RT-PCR of RPBI from maize. RT-PCR was unsuccessful using RNA prepared using the CTAB method from *Cycas revoluta*, *Equisetum hyemale*s, *Marchantia polymorpha* and *Psilotum nudum*. This, too, was inhibited by secondary metabolites that could not be separated from the RNA. The Qiagen RNeasy kit was later used to successfully further purify total RNA from *E. hyemale*s that had been prepared three years previously by the CTAB method. The Qiagen RNeasy kit was then used to successfully prepare total RNA from fresh tissue of *Magnolia soulangeana*, *Nymphaea odorata*, *Osmunda claytonia*, *Ephedra viridis* and *Zea mays* (Pioneer variety 3953). RNA from *Araucaria heterophylla* and *Ginkgo biloba* could not be obtained using this method, though, so the CTAB method was used to obtain RNA from these taxa.

Amplification of RPBI:

By aligning RPBI sequences from animals, fungi, *Arabidopsis* and soybean (*Glycine max*) using ClustalV (Higgins, 1994), conserved regions (A-I) were identified (Figure 1), and degenerate oligonucleotide primers were designed, with reference to conserved amino acid motifs identified by Sidow and Thomas (1994). The primers correspond to the ECPGHFGH motif in domain A (5' GAI TGY CCI GGI CAY TTY GG 3' forward), the PYNADFDGDEM N motif in domain D (5' CCI TAY AAY GCI GAY TTY GAY GGI GAY GAR ATG AA 3' forward) and the VGQQNVEG motif in domain F (5' CC YTC IAC RTT YTG YTG ICC IAC 3' reverse). Table 2 represents sections of the RPBI amino acid alignment from which D and F PCR primers were designed, including parts of

Conserved Domains in RNA Polymerase II Largest Subunit



Anticipated PCR Products

Figure 1: Schematic diagram of the coding region of RPB1, with the blocks of conserved amino acids highlighted to identify the positions of corresponding PCR primers. Zigzag lines indicate the carboxy terminal repeats, while vertical lines represent the size of the protein-coding nucleotide sequence.

Table 2: Alignments corresponding to domains D and F of RPB1, including chloroplast (cp), *E. coli*, archaeal and RPA1 and RPC1 homologues. Corresponding species names are given in Table 3. Domain D is specific to type A RNA polymerases. Domain F is characterized by an α -amanitin sensitive domain only in RPB1, and not in other homologues. Amino acid residues identical to *Arabidopsis* sequence are indicated by a dot.

| <u>Species</u> | <u>Domain D</u> | <u>Domain F</u> |
|-----------------------------------|-------------------------|------------------------|
| <i>Arab.th-II</i> | PYNADFDGDEM | VGQQNVEG |
| <i>Sac.cer-II</i> | |S.... |
| <i>Sch.pom-II</i> | |I.... |
| <i>Try.bru-II</i> | |A. |
| <i>Cae.ele-II</i> | | |
| <i>Dro.mel-II</i> | | |
| <i>Mus.mus-II</i> | | |
| <i>Sac.cer-I</i> | A..... | L...AL.. |
| <i>Try.bru-I</i> | SF...T..... | L...LFD. |
| <i>Sac.cer-III</i> | |IIS. |
| <i>Try.bru-III</i> | |T.S. |
| <i>H.halobium A'</i> | |A.R. |
| <i>E.coli β'</i> | A.....Q.A | R.LMAKPD |
| <i>Marchantia cp</i> | GF.....Q.A (β') | R.LMSDPQ (β'') |

alignments by Puhler *et al* (1989a). This alignment shows that the amino acid sequence encoded by primer D is specific to RNA polymerases, and that the amino acid sequence encoded by primer F is specific to RPB1, and that they would not amplify RNA polymerase genes from chloroplast and mitochondrial genomes.

All primers were synthesized by Cortec DNA Service Laboratories or by Life Technologies. PCR reactions were performed following cycling conditions of: 2 minutes at 94°C, 35 cycles of 40 seconds at 94°C, 1 minute at 55°C and 3 minutes at 74°C, followed by 10 minutes at 74°C. 200 ng of plant genomic DNA was used in PCR amplifications consisting of 2.5 - 5.5 mM MgCl₂, 0.75 μM each primer, 200 μM each dNTP, 0.025 U/μl *Taq* DNA polymerase, 1 X PCR buffer (supplied by the manufacturer of the *Taq* polymerase), and an additional buffer described by Ponce and Michol (1992) consisting of 30 μM tricine pH 8.4, 5 μM 2-mercaptoethanol, 10 μg/ml BSA and 0.05% Tween-20. PCR reagents were obtained from Promega, Pharmacia, Boehringer-Mannheim, and BDH.

Positive, negative and single primer controls were included in each batch of PCR amplifications. A clone of *Arabidopsis* RPB1 (pOSA2, obtained from R. Larkin and T. Guilfoyle) was used as the template for the positive control to show that the amplification conditions worked. No DNA template was included in the negative control, intended to indicate if any contaminants were present in the reagents. Single primer controls were made to indicate if any of the PCR products were artifacts which resulted from nonspecific binding of a single primer.

Southern blots of D-F PCR amplifications from genomic DNA of 13 plants were probed with the *Arabidopsis* RPB1 clone. This revealed a 1.3 kb putative RPB1 PCR product from each of the

13 plants. This corresponded to the 900 bp of coding sequence in addition to three introns found between my domain D and F RPBI primers in *Arabidopsis*.

First-strand cDNA was prepared from total maize RNA using the Amersham First-Strand cDNA Synthesis Kit, according to the manufacturer's instructions, priming with random hexanucleotides. One-tenth of the cDNA reaction was then used directly as the template for an RPBI PCR amplification, as described above.

Cloning:

PCR products of the expected size were isolated from agarose gels and purified using the GeneClean II Kit (Bio 101). Purified products from maize genomic templates were blunt-ended with Klenow fragment and T4 polynucleotide kinase purified, and ligated into the dephosphorylated *EcoRV* site of pBluescript II SK+ (Stratagene) with T4 DNA ligase at 10°C overnight. This was inefficient, and cloning attempts with the Stratagene pCRScript SK+ blunt-ended cloning kit were not fruitful, so all other RPBI PCR products were cloned using the TOPO TA Cloning Kit (Invitrogen) into pCR2.1 according to the manufacturer's instructions. Ligations were transformed into *E. coli* strain DH11S by heat shock into RbCl-competent cells or using the Transformaid Kit according to the manufacturer's instructions (MBI Fermentas), or heat-shock transformed into OneShot *E. coli* cells as directed (Invitrogen). Transformants were grown using Luria-Bertani agar and broth under ampicillin selection, with the addition of IPTG and X-gal to the agar to facilitate blue/white colour selection of transformants (Sambrook *et al* 1989). Enzymes were from New England Biolabs, Pharmacia, Stratagene and Boehringer-Mannheim. LB agar and broth components were from Difco and BDH.

Putative clones were identified by colony hybridization or according to the size of their inserts. Insert size was found either by restriction analysis of the plasmids or by PCR amplification of the vector's multiple cloning site using M13 forward and reverse sequencing primers (adapted from Sandhu *et al.*, 1989).

Plasmids were prepared by the alkaline lysis method of Sambrook *et al.* (1989), or with the Eclipse Mini Kit or QIAprep Spin Miniprep Kit, according to the manufacturer's directions (Gordon Technologies, Qiagen), with the following modifications: After the addition of lysis solution, the cells were incubated at room temperature to lyse for 3-5 minutes prior to the addition of the high-salt buffer. Upon addition of the salt, the preparations were microcentrifuged for 7-10 minutes at 13 krpm to pellet all of the cellular debris.

In order to release their insert, clones in the pCR2.1 vector were analysed by digestion with *EcoRI*, those in pBluescript SK+ were digested with *EcoRI* and *BamHI* and clones in pCRScript SK+ were digested with *EcoRI* and *SacI*. 2 ml of the plasmid prep (125 -400 ng) were completely digested, as directed, in a final volume of 15 µl and then visualized on a 0.7% agarose, 1 X TBE gel run at 10 V/cm.

Colony blots were performed as described in the Hybond-N membrane instruction booklet (Amersham) and probed with a homologous probe prepared using P³²-labelled dCTP using the Pharmacia T7 QuickPrime Kit, as described by the manufacturer. Membranes were sandwiched between two pieces of Whatman 3MM paper and hybridized at 42°C in 50% formamide, 5 x SSC, 0.5% SDS, 50 µg/ml denatured sheared salmon sperm DNA, 5 x Denhardt's solution. The membranes were washed with gentle agitation three times at room temperature in 1 l of 2 x SSC,

0.1% SDS for five minutes, then washed twice for 15 minutes at 65°C in 1 l of 0.1 x SSC, 0.5% SDS.

PCR screening was performed by using a sterile toothpick to inoculate a white colony into a 10 µl PCR mixture and then to inoculate a patch on an LB agar plate. The PCR mixture consisted of 150 µM each dNTP, 1 µM each primer (M13 forward and M13 reverse, from Invitrogen), 0.025 U/ml *Taq* DNA polymerase and 1 x the PCR buffer containing 1.5 mM MgCl₂ supplied with it (Pharmacia or Boehringer-Mannheim). Cycling conditions were: 2 minutes at 94°C followed by 30 cycles of 1 minute at 94°C, 1 minute at 53°C and 1½ minutes at 74°C, and 10 minutes at 74°C. The entire PCR reaction, representing the insert in the cloning vector, was then visualized by agarose gel electrophoresis, as described above.

Single-stranded versions of these RPB1 clones were prepared by a protocol adapted from Garber *et al* (1993). A fresh culture of the clone was inoculated into 2 ml of LB broth with 100 µg/ml ampicillin and 10 mM MgCl₂, incubating it at 37°C for two hours, then infecting these cells with F1 helper phage and then increasing the media by 8 ml, followed by an incubation at 37°C for an additional 8 hours before harvesting the phage. To harvest, the culture was centrifuged and the phage precipitated from the supernatant with ¼ volume of 20% PEG, 2.5 M NaCl at 4°C for 1-12 hours. The suspension was centrifuged and the phage pellet resuspended in 1.25 ml TERPS (10 mM Tris, 1 mM EDTA, 40 µg/ml RNase A, 2.5 µg/ml Proteinase K, 0.5% SDS) and lysed for 2 hours at 65°C. The phage coat was precipitated with 1/5 volume of 5M potassium acetate, the suspension centrifuged, and the phagemid DNA was precipitated from the supernatant with 0.7 volumes of cold isopropanol. A pellet of single-stranded DNA was washed twice with cold 70% ethanol, dried and resuspended in 50 µl of TE.

Sequencing:

Single- and double-stranded versions of the RPBI clones were sequenced according to the Sanger-dideoxy method, by primer walking. Initial sequence was obtained manually using the Pharmacia T7 Sequencing Kit and S³⁵-labelled dATP (Amersham), according to the manufacturer's instructions, with the addition of 3 µl DMSO and annealing at 65 °C to resolve compressions caused by secondary structure in the sequence. These reactions were run on 0.4 - 0.8 mm thick wedge-shaped denaturing polyacrylamide gels (4% acrylamide, 0.5 x TBE, 8M urea) at 1700 V for 5 and 17 hours using electrophoretic apparatus from Owl Scientific. Sequencing gels were vacuum-dried (BioRad) onto Whatman 3MM paper and exposed to Kodak BioMax MR film for 1 - 5 days.

Further sequence was obtained using the ABI Cycle Sequencing Core Kit with AmpliTaq FS or the ABI BigDye Terminator Cycle Sequencing Kit with AmpliTaq FS (ABI, Perkin-Elmer Corporation) according to the manufacturer's instructions. These automated sequences were read by the University of Ottawa Biotechnology Research Institute. Plant-specific degenerate internal sequencing primers were designed from my alignment (prepared using ClustalV under default settings) of the sequenced ends of each clone. CODEHOP (Rose *et al.*, 1998) was used to help design primers. The primer positions correspond to the IIPKQIN motif (5' CCT GAT CAT CCC GAA GCA RAT HAA YHT 3' forward and 5' T GAT CAG GTT GAA GAC CTG CTT NCC NGT CCA 3' reverse) and the GIGDTIAD motif (5' GGI ATI GGI GAY ACI ATI GCI GA 3' forward and 5' TC IGC IAT IGT RTC ICC IAT ICC 3' reverse).

Northern and Southern Hybridizations:

Northern blots of 20 µg of total RNA were prepared using the Ambion NorthernMax kit as directed, with the following modifications. A 0.8% SeaKem agarose formaldehyde gel, prepared as indicated by Sambrook *et al* (1989) was used to fractionate the RNA at 2V/cm for 18 hours. The RNA was then capillary transferred using alkaline 5 x SSC for 5 hours, and the membrane was hybridized and washed as directed (Ambion).

Southern blots were prepared as indicated by Sambrook *et al* (1989), with the following adaptations. 10 µg of total DNA were completely digested singly with *EcoRI* and *BamHI* or *EcoRI*, *PstI* and *XbaI* under standard conditions and electrophoresed on 0.8% SeaKem agarose gels with 1 X TAE buffer at 45 V for 19 hours, then capillary transferred to nylon membranes overnight using 20 x SSC as the transfer medium.

Both Northern and genomic Southern blots were hybridized with their corresponding RPB1 clones. For the Southern blot of PCR products, 10 µl of each PCR was run in a 0.7% agarose gel in 1 x TBE at 10 V/cm for an hour. RNA or DNA was cross-linked to the damp nylon membranes by exposure to ultraviolet light.

Membranes were pre-hybridized for 25 minutes to 2 hours at 42°C in 25 ml of 50% formamide, 5 x SSC, 0.5% SDS, 50 µg/ml denatured sheared salmon sperm DNA, 5 x Denhardt's solution. Hybridization was in 10 ml of this solution at 42°C for 18 - 24 hours. Calf thymus DNA was substituted for salmon sperm DNA in the hybridization solution for the genomic Southern blots. Background hybridization was also avoided by removing unincorporated nucleotides from the probe using Pharmacia ProbeQuant columns (according to the manufacturer's instructions) and by sandwiching the membranes between two pieces of Whatman 3MM paper during the hybridization.

Northern and genomic Southern blots were probed with a homologous probe prepared using P³²-labelled dCTP using the Pharmacia T7 QuickPrime Kit according to the manufacturer's instructions.

Membranes were then washed twice at room temperature in 1 l of 2 x SSC, 0.1% SDS for 10 minutes, then washed twice for 10 minutes at 65°C in 1 l of 2 x SSC, 0.1% SDS. The Southern blot of RPBI PCR products was probed with a heterologous probe (pOSA2, a clone of RPBI from *Arabidopsis*), prepared and hybridized as described above, but only washed twice at room temperature in 1 l of 2 x SSC, 0.1% SDS for fifteen minutes. The membranes were exposed to Kodak Biomax MS film with Kodak Biomax Transcreen HE intensifying screens at -80°C for 2-7 days.

Sequence Analysis:

Plant RPBI D-F nucleotide sequences were edited using GDE and Sequedit v.0.912 (Drouin *et al.*, 1999). Additional eukaryotic RPA1, RPBI and RPC1 nucleotide and amino acid sequences were obtained from Genbank using BLAST (Altschul *et al.*, 1990 and 1997) and Entrez, except for the *Giardia lamblia* RPBI DNA sequence which was obtained from Hans-Peter Klenk. Table 3 provides full identification of the sequences used in alignments in this report. Peter Foster's "pgb_cds" PERL script was used to isolate the coding region of these nucleotide sequences. The amino acid alignment of animal, plant and fungal sequences, from which plant-specific RPBI PCR primers were designed, was prepared using Clustal V (Higgins, 1994). This was done via GDE version 2.0 (Smith *et al.* 1994), using a floating gap penalty of 90 and a fixed gap penalty of 10 besides the default settings, with further manual adjustments made using GDE. My own sequences were identified as plant RPBI D-F using basic BLAST searches (Altschul *et al.*, 1990 and 1997).

Table 3: Codes and sources for sequences included in phylogenetic analyses of type A RNA polymerases.

| Species name | Code | Polymerase | Accession # | Authors |
|----------------------------------|--------------------|------------|---------------------|---|
| <i>Sulfolobus acidocaldarius</i> | <i>Sulf.acido</i> | A' | I133407 | Puhler <i>et al.</i> , 1989b |
| <i>Porphyra yezoensis</i> | <i>Por.yez-II</i> | RPBI | U90208 | Stiller and Hall, 1997 |
| <i>Bonnemaisonia hamifera</i> | <i>Bon.ham-II</i> | RPBI | U90209 | Stiller and Hall, 1997 |
| <i>Drosophila melanogaster</i> | <i>Dro.mel-II</i> | RPBI | M27431 | Jokerst <i>et al.</i> , 1989. |
| <i>Artemia salina</i> | <i>Art.sal-II</i> | RPBI | U10331 | Sidow and Thomas, 1994. |
| <i>Crassostrea gigas</i> | <i>Cra.gig-II</i> | RPBI | U10334 | Sidow and Thomas, 1994. |
| <i>Ilyanassa obsoleta</i> | <i>Ily.obs-II</i> | RPBI | U10338 | Sidow and Thomas, 1994. |
| <i>Helobdella stagnalis</i> | <i>Hel.sta-II</i> | RPBI | U10336 | Sidow and Thomas, 1994. |
| <i>Homo sapiens</i> | <i>Hom.sap-II</i> | RPBI | X63564 | Wintzerith <i>et al.</i> , 1992. |
| <i>Mus musculus</i> | <i>Mus.mus-II</i> | RPBI | M14101 | Ahearn <i>et al.</i> , 1987. |
| <i>Caenorhabditis elegans</i> | <i>Cae.ele-II</i> | RPBI | M29235 | Bird and Riddle, 1989. |
| <i>Spirogyra spp.</i> | <i>Spirogy-II</i> | RPBI | U90210 | Stiller and Hall, 1997. |
| <i>Gingko biloba</i> | <i>Gin.bil-II</i> | RPBI | | this study |
| <i>Cycas revoluta</i> | <i>Cyc.rev-II</i> | RPBI | | this study |
| <i>Magnolia soulangeana</i> | <i>Mag.sou-II</i> | RPBI | | this study |
| <i>Nymphaea odorata</i> | <i>Nym.odo-II</i> | RPBI | | this study |
| <i>Zea mays</i> | <i>Zeamays-II</i> | RPBI | | this study |
| <i>Glycine max</i> | <i>Gly.maxC-II</i> | RPBI | X52495 | Dietrich <i>et al.</i> , 1990. |
| <i>Arabidopsis thaliana</i> | <i>Arab.th-II</i> | RPBI | X52954 | Nawrath <i>et al.</i> , 1990. |
| <i>Ephedra viridis</i> | <i>Eph.vir-II</i> | RPBI | | this study |
| <i>Araucaria heterophylla</i> | <i>Arau.he-II</i> | RPBI | | this study |
| <i>Dictyostelium discoideum</i> | <i>Dic.dis-II</i> | RPBI | S52651, AF058710 | Lam <i>et al.</i> , 1992; Stiller and Hall, 1998. |
| <i>Acanthamoeba castellanii</i> | <i>Aca.cas-II</i> | RPBI | U90211 | Stiller and Hall, 1997. |
| <i>Saccharomyces cerevisiae</i> | <i>Sac.cer-II</i> | RPBI | M11190 | Allison <i>et al.</i> , 1985. |
| <i>Schizosaccharomyces pombe</i> | <i>Sch.pom-II</i> | RPBI | S56564 | Azuma <i>et al.</i> , 1991. |
| <i>Nosema locustae</i> | <i>Nos.loc-II</i> | RPBI | AF061288 | Hirt <i>et al.</i> , 1999. |
| <i>Vairimorpha necatrix</i> | <i>Vai.nec-II</i> | RPBI | AF060234 | Hirt <i>et al.</i> , 1999. |
| <i>Plasmodium falciparum</i> | <i>Pla.fal-II</i> | RPBI | I33328 | Li <i>et al.</i> , 1989. |

| Species name | Code | Polymerase | Accession # | Authors |
|---------------------------------------|-------------------|-------------|-------------|-----------------------------------|
| <i>Mastigamoeba invertens</i> | <i>Mas.inv-II</i> | <i>RPB1</i> | AF083338 | Stiller <i>et al</i> , 1998. |
| <i>Trichomonas vaginalis</i> | <i>Tri.vag-II</i> | <i>RPB1</i> | U20501 | Quon <i>et al</i> , 1996. |
| <i>Trypanosoma brucei</i> | <i>Try.bru-II</i> | <i>RPB1</i> | 133329 | Evers <i>et al</i> , 1989. |
| <i>Giardia lamblia (intestinalis)</i> | <i>Gia.lam-II</i> | <i>RPB1</i> | | H.-P. Klenk, unpublished. |
| <i>Saccharomyces cerevisiae</i> | <i>Sac.ce-III</i> | <i>RPC1</i> | X03129 | Allison <i>et al</i> , 1985. |
| <i>Giardia lamblia (intestinalis)</i> | <i>Gia.la-III</i> | <i>RPC1</i> | P25202 | Lanzendorfer <i>et al</i> , 1992. |
| <i>Trypanosoma brucei</i> | <i>Try.br-III</i> | <i>RPC1</i> | P08968 | Kock <i>et al</i> , 1988. |
| <i>Mus musculus</i> | <i>Mus.musc-I</i> | <i>RPA1</i> | 35134 | Seither <i>et al</i> , 1997. |
| <i>Drosophila melanogaster</i> | <i>Dro.mela-I</i> | <i>RPA1</i> | P91875 | Knackmuss <i>et al</i> , 1996. |
| <i>Trypanosoma brucei</i> | <i>Try.bruc-I</i> | <i>RPA1</i> | P16355 | Jess <i>et al</i> , 1989. |
| <i>Saccharomyces cerevisiae</i> | <i>Sac.cere-I</i> | <i>RPA1</i> | P10964 | Memet <i>et al</i> , 1988. |
| <i>Schizosaccharomyces pombe</i> | <i>Sch.pomb-I</i> | <i>RPA1</i> | P15398 | Yamagishi and Nomura, 1988. |

RPB1 D-F amino acid sequences including my own were later aligned for phylogenetic analysis using ClustalW 1.71 (Thompson *et al.*, 1994) under default settings and then refined manually in GDE, with reference to eukaryotic RPB1 alignments obtained from John Stiller (Stiller and Hall, 1997) and John Logsdon (Hirt *et al.*, 1999). The coding nucleotide sequences were aligned manually with the amino acid alignment using GDE and Sequedit v.0.912. To reduce the effect of compositional bias introduced by synonymous substitutions, only the first and second bases of codons were used in the analysis. "M" and "y" were substituted for the first codon position of leucine and arginine codons, to reduce compositional bias. P. Foster's "pgb_subset" PERL script was used to remove the third base of codons in the nucleotide sequence alignment. The alignments were truncated to contain only the region of the RPB1 fragment within my PCR primers for regions D and F. PHYLIP v.3.573c (Felsenstein, 1993) and PUZZLE v.4.0.2 (Strimmer and von Haeseler, 1996 and 1997) were used to infer phylogenetic trees of both the amino acid and nucleotide sequences.

Phylogenies were first inferred from nucleotide and amino acid sequence alignments using the PROTPARS, DNAPARS, PROTDIST, DNADIST, NEIGHBOR, SEQBOOT, and CONSENSE programs available from PHYLIP v.3.573c (Felsenstein, 1993). The jumble option was used wherever possible to avoid bias from the input order of sequences. The input order of bootstrapped datasets was only randomized once. Global rearrangements were allowed. The PROTDIST algorithm was set to use the Dayhoff PAM model of nucleotide substitution, while both the Kimura and maximum likelihood models and remaining default settings were used for DNADIST. Ten replicates were made of the distance analyses, jumbling once using ten different random seed numbers. Both amino acid and nucleotide multiple sequence alignments were subjected to 1000

bootstrap replicates using SEQBOOT, and consensus distance and parsimony trees were obtained using CONSENSE.

Next, PUZZLE v. 4.0.2 (Strimmer and von Haeseler, 1996 and 1997) was used to infer phylogenetic trees using my plant nucleotide and amino acid alignments, in addition to alignments of RPB1 D-F amino acid sequences amongst animals, plants and fungi, or including RPA1, RPB1 and RPC1 D-F of all eukaryotes. This program uses the quartet puzzling and maximum likelihood approaches to phylogenetic inference. The parameters for the plant amino acid and nucleotide analyses included 10000 puzzling steps, exact parameter estimates, quartet sampling, neighbor-joining tree-fitting, the computation of both clock like and non-clock like branch lengths and use of a mixed model of substitution rate heterogeneity which includes an invariable rate and eight gamma rates. The JTT (Jones *et al.*, 1992) and Blosum 62 (Henikoff and Henikoff, 1992) models of amino acid substitution and the Schoeniger-von Haeseler model of nucleotide substitution were used. Nucleotide analyses were made with this program once by selecting all sites for the analysis; next, selecting only the first and second sites of codons (with "m" or "y" in the first position of leucine and arginine codons); and finally, selecting only the second site of codons. The phylogenetic trees inferred using PUZZLE and PHYLIP were drawn with Treeview and saved as ".wmf" (Windows Metafile) format files, which were then imported into Corel Presentations 8 for further labelling.

Nucleotide frequencies of each sequence at the different codon positions were measured using CODONS 1.4 (Lloyd and Sharp, 1992), which was also used to calculate the amino acid composition. Weighted numbers of synonymous (K_s) and non-synonymous (K_a) substitutions were calculated using Li93 (Li, 1993). The proportion of different synonymous (P_s) and nonsynonymous (P_a) sites was calculated using the equation $P = \frac{3}{4}(1 - e^{-4/3K})$. RRTree v.0.6 (Robinson *et al.*, 1998)

(P_a) sites was calculated using the equation $P = \frac{3}{4}(1 - e^{-(4/3)K})$. RRTree v.0.6 (Robinson *et. al.*, 1998) was used for relative rate test calculations. MS Excel '97 was used to construct all graphs. All sequence analysis programs were used in the context of Windows98, except for pgb_cds, pgb_subset, Li93, GDE and ClustalV, which were used on a Sun Sparc 4 workstation.

RESULTS

Preliminary Analysis of the Region Between Domains F and H of RPB1:

To begin the project, an alignment of published RPB1 amino acid and nucleotide sequences from Animals, Plants and Fungi, with *Dictyostelium discoideum*, a slime mold, as an outgroup, was prepared to use to estimate the usefulness of this gene for inferring the phylogeny of land plants and to use as a guide for designing PCR primers to be used to clone the gene from land plant species. For phylogenetic analysis, the alignment was truncated to include domains F-H of RPB1 (Figure 2), since this was prior to the publication of the sequence for domains A-F from *Dictyostelium* by Stiller and Hall (1998). Figure 3 represents a maximum likelihood phylogenetic tree inferred from amino acid sequences from domain F-H. Neighbor-joining and parsimony trees inferred from the same dataset and from a corresponding alignment of nucleotide sequences yielded very similar trees (data not shown). Analysis of nucleotide composition showed that there was little compositional bias at the first and second codon positions, but that the third positions were biased (data not shown). The consistent and strong statistical support for the relationships in this preliminary analysis indicated that RPB1 would be a good phylogenetic marker with which to study the evolutionary relationships of land plants.

Figure 2: A 606-position amino acid alignment of 13 sequences from RPB1 F-H, with the slime mold *Dicthostelium discoideum* as the outgroup to Animals, Plants and Fungi. Residues identical to that of the outgroup are shown as points (.), while gaps in the alignment are shown as a dash (-). The alignment was prepared with Clustal V with a floating gap penalty of 90 and a fixed gap penalty of 10, with further manual refinement.

| | |
|--------------------|--|
| | 100 |
| <i>Dic.dis-II</i> | AMGGREGLIDTAVKTSETGY IQRCLVKAMEDVSIKYDATV RNSLGDVIOQFAYGEDGIDGC FVENQSIDSLRKDNTELERM YRHQVDKPDYDGGWMDPLVI |
| <i>Cae.ele-II</i> |A.....R.I.....S.MVN.G.....AQMV.LR.....L.M.W.....NMPTMKPN.AVF..D.F.VS.AQNAIKLMDLTDNKF |
| <i>Dro.mel-II</i> |A.....R.I.....S.MVN.G.....V.QL..LR.....LC.E.L..F.NMPTVKLS.KSF.KR.FK-----FDWSNERL |
| <i>Art.sal-II</i> |A.....R.I.....ACMVA.G.....V.Q.....LR.....LA.E.L..F.LPTIKLS.RAF.SK .-----FDPSNERQ |
| <i>Cra.gig-II</i> |A.....R.I.....S.MV.G.....QVEQLV.LR.....L.AT.H..F.TMPT.KPS.RAF..Q.FK-----FDATNERN |
| <i>Hel.sta-II</i> |A.....R.I.....S.MV.G.....QIEQL..LR.....LA.E.W..F.NLP..KPS.KAF.AG.FK-----FDPTNEKH |
| <i>Ily.obs-II</i> |A.....R.I.....S.MV.G.....QVEQLV.LG.....L.A.H..F..LPT.KPSTRAF..Q.F.-----FDPTDERM |
| <i>Mus.mus-II</i> |A.....R.I.S..S.MV.....INQ.V.LR.....LA.E.S..F.NLAT.KPS.KAF.KK.F.-----FDYTNERA |
| <i>Hom.sap-II</i> |A.....R.I.S..S.MV.....INQ.V.LR.....LA.E.S..F.NLAT.KPS.KAF.KK.F.-----FDYTNERA |
| <i>Sac.cer-II</i> |A.....R.....L..IMVH.N.T.....N.....I.....M.AA.HI.K..L.TIGGSDAAF.KR .-----VDLLNTDH |
| <i>Sch.pom-II</i> |A.....R.....MVR.G.....AM..I.....L.AT.L..Y.VF....LSTKQF.KK .-----IDLMEDRS |
| <i>Arab.th-II</i> |R.....IMV.G.....M.AV.WI.S.KL...KMKKS.FD.T.FKYEI.DENWNPTYLSDEHL |
| <i>Gly.maxC-II</i> |R.....IML.G.....M.AI.WI.T.KL.T.KMKK..FD.V.F.YEF.EENWKPNY.LQEPV |
| | domain F |
| | 200 |
| <i>Dic.dis-II</i> | EHRNDSLTRDTLEKEFERI KSDRSLLRNEIIPSGEANWP LPVNLRLKLNNAQKLFNIDI RRVSDLNPAVVVLEIEKLV A RLKIIATADTTEDDENFNRA |
| <i>Cae.ele-II</i> | LRKNYSEDVVREIQESEDG. SLVE.EWSQLEEDRRLLRKD F.RGDAKIVLPCNL.RL.WN AQKIFKVDLRNAVNLSP.HV I.SGVRELSKKLIIVSGNDEI |
| <i>Dro.mel-II</i> | MKKVFTDDVIKEMTDSS.A. QELEAEWDRLVSDRDSLQI F.NGESKVVLPCNLQRM.WN VQKIFHINKRLPTDLSPIRV IKGVKTLLERCVCIVTGND.I |
| <i>Art.sal-II</i> | LRKTFTEDEVTKS.LGDS.V. TEIEKEWEALVKDREALRKV F.SGENKVVLPCNLQRM.WN AQKTFHINKRMPDLSLSP.RV IQGVRDLLKNCVIVNGEDKL |
| <i>Cra.gig-II</i> | MKKCLSEEVIKD.MGDALAV SQLDREWEQLREDRDLRSI F.TGDSKVVLPCNLQRM.WN AQKIFRIDTHKPTDLHPKV VEGVEDLCKRLIVVAGGD.L |
| <i>Hel.sta-II</i> | LKNYL.EDILXS.XXDANV. AEVE.EYKQLEDRTAIRQI F.SGDSKIVLPCNLQRL.WN AQKIFRIHTRKPSNLHPVKI IEDVRELSKKFMIVKGED.L |
| <i>Ily.obs-II</i> | MKRCLKEDVIK.G.RG.HRL. EELE.EWLQL.SDRDS.RQV.FSTGDA.IVLPCLQRM.WN AQKIFRIDKRRPSDLDPINV VSGVRDLQCRLTIV.GEDMI |
| <i>Mus.mus-II</i> | LRRTLQEDLVKDVLSNAHIQ NELEREFERMREDREVLRI F.TGDSKVVLPCNL.RM.WN AQKIFHINPRLPSDLHPKV VEGVKELSKKLVIVNGDDPL |
| <i>Hom.sap-II</i> | LRRTLQEDLVKDVLSNAHIQ NELEREFERMREDREVLRI F.TGDSKVVLPCNL.RM.WN AQKIFHINPRLPSDLHPKV VEGVKELSKKLVIVNGDDPL |
| <i>Sac.cer-II</i> | TLDPSSLSESGEILGLDLKQ VLLDEEYKQLVKDRKFLREV FVDGEANWPLPVNIRRI.QN AQQTFHIDHTKPSDLTIKDI V.GVKDLQENLLVLRGK.EI |
| <i>Sch.pom-II</i> | L---SLYMENSIENDSSVQD LL.EEYTLQVADRELLCKFI F.KGDA.WPLPVNVQRI.QN ALQIFHLE.KKPTDLLPSDI INGLNELIAKLTIFRGS.D.I |
| <i>Arab.th-II</i> | .DLKGIREL..VFDA.YSKL.ET..FQ.GT..ATN.DST.....IK.H.W.....T.K..L.KI..MH.VEI.DAVD..QE .LVPVGD.ALSVEAQK.-- |
| <i>Gly.maxC-II</i> | .DLKTIREF.NVF.A.VQKL.EA..HQ.AI..ASN.DNSL.....K...W.....T.KV.F..P..MH.MEI.EA.D..QE .LVPVGE.ALSVEAQK.-- |

Dic. dis-II WAEVYFNATMLFSILVRSTF ASKRVLTEFRLTEKAPFLWVC GEIESKFLQALAHGEMVGA LAAQSIGEPATQMTLNTFHY AGVSSKNVTLGVPRLNQIIN
Cae. ele-II SKQAQY...L.MN.L...L CT.NMC.KSK.NSE..D.LL...R.Q..I.Q...L...A...KE...
Dro. mel-II SKQANE...L.QC.I...L CT.Y.SE...STE..E.LV...TR.Q..Q.N...L...F...KE...
Art. sal-II SKQANE...L.QC...L CT.LAD.Y.NTE..E.LI...TR.QL.Q.V...I...L...F...KE...
Cra. gig-II SIQANA...L.KC.I...L CA...TE...SSE..E.LI...V.R.Q..Q...L...A...KE...
Hel. stu-II SKTANT...L.MN...L C...IE..H.STE..E.LM...I.RVQ...L...X...KE...
Ily. obs-II SKQANE...L.MKC.I...L CA...AE.H..S.D.N.LL...TR.Q..Q...L...A...KE...
Mus. mus-II SRQAQE...L.N.HL...L C.R.MAE...SGE..D.LL...N.I...L...A...KE...
Hom. sap-II SRQAQE...L.N.HL...L C.R.MAE...SGE..D.LL...N.I...L...A...KE...
Sac. cer-II IQNAQRD.VT.CC.L.RL.TR...Q.Y..KQ..D.L SN.AQ.RSVV...V...F...A.K.S...KE.L.
Sch. pom-II TRD.QN...L.Q.L.K..V...IM.Y..NKV..E.IM..V.AR.Q..VVS...T...A...KE.L.
Arab. th-II -----L.F.N.L...L...E.YK.SRER.E..I...R..S.VA...I.C.V...A...RE...
Gly. maxC-II -----L.N.L...L...E.Y..SRES.E..V...R..S.VV...I.V.V...A...REL...

domain G

Dic. dis-II IAKQVKTPSLTIYLYKPHMAR DMDRAKIVKSQLEYTTLANV TSATEIYYDPDPQNTIISED AEFVNSYFELPDEEIDVH-- S MSPWLLRIELDRG-MVTDK
Cae. ele-II VS.TL...VF.TGAA.K.PEK..D.LCK.H...H...CN.A...K.V.A..E.W.SIFY.M.--DH.LS--RT...KR...
Dro. mel-II .S.KP.A...VF.TGGA...AEK..N.LCR.H..RK..AN.A...R.V...Q...V.Y.M.--DF.PT--RI...KR.T...
Art. sal-II VSRSP.A...VF.TGTA...AEK..Q.LCR.H..KK..AN.A...XIK.V.A..QD..TL.YDM.--DV.TS--R...X...KR.T...
Cra. gig-II .S.KP...V..IGQA...AEK..D.LCR.H..RK..AN.A...V...Q.W.V.Y.M.--DF.S--KI.A...KR.T...
Hel. stu-II VS.KPRA...VI.IGQP...AEK..D.LC..H..RK..EN.A...MH.L.E..Q.W.YI.YDM.--DV.IS--RL...V...KR.T...
Ily. obs-II .S.KP...V..NGPA...AEKC.D.LCR.H..RK..AN.A...M.V.V..Q.W.SI.Y.M.--DF.AS--RI...KR.T...
Mus. mus-II .S.KP...VF.LGQS...AE..DILCR.H..RK..AN.A...N.S.VVA..Q.W.V.Y.M.--DF.A--RI...KH.T.R
Hom. sap-II .S.KP...VF.LGQS...AE..DILCR.H..RK..AN.A...N.S.VVA..Q.W.V.Y.M.--DF.A--RI...KH.T.R
Sac. cer-II V..NM...V..E.GH.A.QEQ..LIR.AI.H..KS..I.S...RS.V.P..E.IIQLH.S.L...AEQSFQ...L...AA.N...
Sch. pom-II V..NI...M.WI.A.N..L..N.QT.I.H..ST...H...D.V.E..KD..EAF.AI...VEENLY.KQ...L...AK.L...
Arab. th-II V..RI...SV..T.EASK.SKEG..T.QCA...RS..Q...VW...MS..E..F...R..Y.M...DVSPD--KI...N.EM...
Gly. maxC-II V..SI...SV...DAGK.TKE...T.QCA...RS..Q...VW...MS..E..VD..M..Y.M...VALE--KI...FV.V..N.EM...

domain G

Dic. dis-II KLTWADITQCVRDFGLSLN CIFSDDNAEKLILIRIMVES QETKGT--DND----DDDQ FLRRIESNMLSEMVLRIKIGK IKKVFMRTEKIPK-----VT
Cae. ele-II . . . EM. ADRIHGG. NDVH T. YT. VF. L. IAGE DKG---EAQEEQVDKME. V . . C. A. . . . DLT. Q. . PA . S. Y. NQPNTDD. KRIII.
Dro. mel-II EQ. AEKINVG. ED. . . N. . . . D. V. . . . IMNN E. N. F-QDE. EAVDKME. M . . C. A. . . . D. T. Q. . EA . G. Y. HLPQTD. KRIVI.
Art. sal-II EQ. SEKITAG. ED. . . N. . . . V. . . . IMN. D. S. F-GEE. EQDKME. M . . C. A. . . . D. T. Q. . EA . S. Y. HLPQTDN. KRIII.
Cra. gig-II EQ. SEKITAG. DD. . . N. . . . V. . . . IMN. D. S. M-QNEEEIVDKME. V . . C. A. L. . D. T. Q. . EA . A. Y. HLPNTED. KRISI.
Hel. sta-II EQ. SEKITAG. DD. . . N. . . . V. V. LMSN . DG. QDQDTEEQIDKMP. T . KH. TD. T. Q. . TS . A. Y. QOPTTDD. KRIIID
Ily. obs-II EQ. SEKITAG. DD. . . N. . . . V. . . . IMN. DDS. M-QDEEEVVDKME. V . . C. . . . L. . D. T. Q. . EA . A. Y. HLPPTDD. KRIHI.
Mus. mus-II EQ. AEKINAG. DD. . . N. . . . V. . . . IMN. D. N. M-QEEEVVDKM. . . V . . C. . . . TD. T. Q. . EQ . S. Y. HLPQTDN. KKIII.
Hom. sap-II EQ. AEKINAG. DD. . . N. . . . V. . . . IMN. D. N. M-QEEEVVDKM. . . V . . C. . . . TD. T. Q. . EQ . S. Y. HLPQTDN. KKIII.
Sac. cer-II D. . . . GQVGERIKQT. KND. F. V. W. E. D. . . . I. C. V. R- --P. SLDAETEA---EE. H. M. KK. NT. ENIT. . VEN . ER. V. MKY-----DRKVP
Sch. pom-II S. S. VAGKIAES. ERD. F. T. W. E. . D. . . . I. C. IIRD DDR. AEDDDNMI---EE. V . . KT. GH. ESIS. . VPV . TR. Y. MEH-----KIVRQI
Arab. th-II S. . . . AEKINLE. DDD. T . . . N. . . . Q. IMND EGP. . ELQDESA---E. V . . KK. . . . T. A. . . . PD . N. . . IK-QVRKSR-----FD
Gly. maxC-II S. . . . A. KINLE. DDD. S . . . N. IMND DAP. . EVQDESA---E. V . . KK. T. . T. . . . PD . N. . . IK-NA. VQ. -----FD

Dic. dis-II ENSGFGVREEWILDTDGVSL LEVMSHPDVDHTRRTSNDIV EIIQVLGIEAVRNALLKELR AVISFDGYSVNYRHLAILAD VMTYRGLHLMATIRHGINRVE
Cae. ele-II PEG. . KSVAD. . E. . TA. . R. L. ERQI. PV. C . FE. K. IER. MD N. L. C. . . AK. YS. Q.
Dro. mel-II . TGE. KAIG. L. E. . T. M. MK. L. ER. . PI. S. . . . C . F. KSVE. MN . LQ. Y. L. L. C. . . AK. QD
Art. sal-II DTGE. RAIA. L. E. . T. . MK. L. ER. . PV. Y. . . . C . FT. KSIE. MN . LQ. Y. L. L. C. . . AK. Y. QD
Cra. gig-II . EGE. KAVA. . E. . TA. MK. L. QR. . PI. . . T . VFSI KSIE. MN H. SL. C. S. . AK. Q.
Hel. sta-II . KGE. KALQD. . E. . TA. RR. L. VEN. . PVK. V . . . VFE. KSIER. MN N. L. C. . . AK. Q.
Ily. obs-II . EGE. KAVA. . E. . S. . MK. L. ER. . PV. Y. T. . . . VFAT K. IER. MV H. L. C. . . AK. Q.
Mus. mus-II . DGE. KALQ. . E. . . . MR. L. EK. . PV. FT. K. ER. Y H. L. C. T. . C. V. . QD
Hom. sap-II . DGE. KALQ. . E. . . . MR. L. EK. . PV. FT. K. ER. Y H. L. C. T. . C. V. . QD
Sac. cer-II PTGEYVKEP. V. E. . . N. S. . . TV. GI. P. . IYT. SEI D. ME. G. A. Y. . VY N. AS. M. L. V. TQ. G. TSV. . . F. . SN
Sch. pom-II . DGT. ERAD. V. E. . IN. T. A. TVEG. . A. . Y. SF. . . L. I. T. S. N. E. L. C. . . S. A.
Arab. th-II . EG. KTS. . M. . E. N. . A. C. E. . PK. HLI . . E. R. D. . . V C. T ND
Gly. maxC-II . . E. . KSN. G. . M. . E. N. . A. C. E. . A. HLI . V. E. RS. . . D. . . I CE T ND

domain H

606

Dic.dis-II
Cae.ele-II
Dro.mel-II
Art.sal-II
Cra.gig-II
Hel.sta-II
Ily.obs-II
Mus.mus-II
Hom.sap-II
Sac.cer-II
Sch.pom-II
Arab.th-II
Gly.maxC-II

TGPLMR
V.A....
..A....
..A....
..A.A.
..V.A.
..A.A.
.....K
.....K
..A....
..A....
.....
...M...
domain H

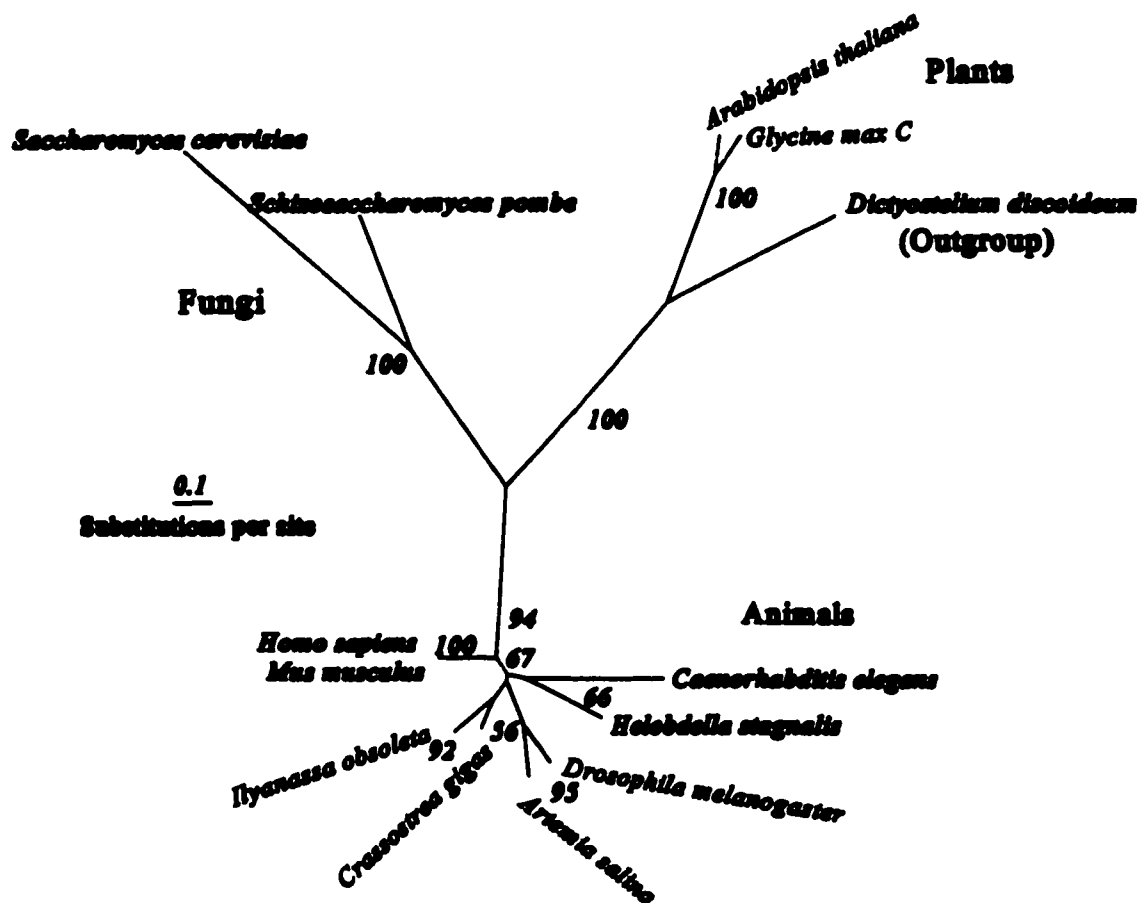


Figure 3: A maximum-likelihood distance phylogenetic tree with non-clocklike branch lengths inferred from a 606-position alignment of RFB1 amino acid sequences of domains F-H from Animals, Plants and Fungi using the JTT model of substitution and an invariable and eight γ -distributed rates of substitution with PUZZLE v4.0.2. The fraction of invariable sites is 0.04 (S.E. 0.03) and the γ parameter $\alpha = 0.67$ (S.E. 0.06). The log L of this tree is -8954.57, while the unlikely tree with clocklike branch lengths had log L = -8991.17. Numbers at the nodes represent the percent of 10000 quartet puzzling steps in support of that relationship, if >50%.

Plant RPBI:

Using the primers for RPBI regions D and F which were described in the “Materials and Methods” section, this region was successfully amplified as a 1.3 kb fragment by PCR from genomic DNA of all of the plants listed in table 1 (Figure 4, gel photo). RPBI D-F PCR products were obtained from both cDNA and genomic DNA templates from *Zea mays*, in which case the nucleotide sequence of the exons was identical. This was confirmed by Southern hybridization (figure not shown) of the PCR products with the *Arabidopsis thaliana* RPBI clone, pOSA2, and by sequencing clones of nine of these PCR products, as summarized in Table 4. Two clones from *Ephedra viridis*, four clones from maize cDNA and from *Cycas* and *Nymphaea*, and three clones from each of the other plants listed in Table 4 were obtained and sequenced on both strands.

Other clones with sequence similarity to RPBI (as determined by BLAST scores) from *Araucaria*, *Marchantia*, *Ephedra*, and *Equisetum* were obtained and partially sequenced (results not shown), but they were not included in further analysis. Even with very clear sequence data, these sequences could not be unambiguously aligned with other plant RPBI D-F coding sequences, due to frame shifts in the inferred amino acid sequence and insertions, deletions, or inversions. Some clones obtained from *Marchantia* and *Ephedra* lacked intron positions conserved amongst other plant and green algal RPBI sequences and had frame shifts and substitutions to other amino acids or stop codons in positions of the inferred amino acid sequences which were perfectly conserved amongst other plant sequences.

PCR amplification of RPBI using the primer for region A vs the region F primer, described in the methods section, was unsuccessful. This was attributed to the abundance of single primer

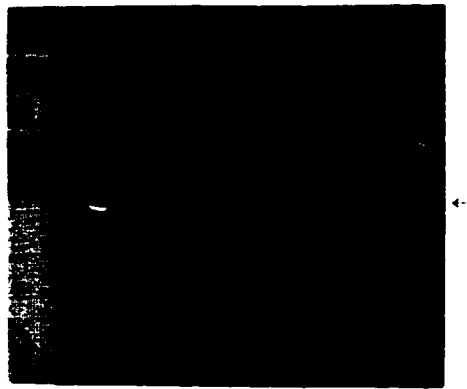


Figure 4: RPBI D-F PCR products from genomic DNA and maize cDNA templates, separated with 0.7% agarose in 1xTBE at 100V for 1 hour. Size markers in the first and last lanes are λ DNA digested with *BstEII*. The arrow corresponds to the major 1.3 kb fragment encoding regions D-F of RPBI. The PCR products (from left to right) are from DNA from the following sources: *Arabidopsis thaliana*, *Araucaria heterophylla*, *Cycas revoluta*, *Ephedra viridis*, *Equisetum hyemale*s, *Ginkgo biloba*, *Gnetum gnemon*, *Magnolia soulangeana*, *Marchantia polymorpha*, *Nymphaea odorata*, *Psilotum nudum*, *Welwitschia mirabilis*, *Zea mays*(total DNA), *Zea mays* (cDNA).

Table 4: Identification of RPB1 PCR products from various plants by hybridization with *Arabidopsis* RPB1 and/or from sequence analysis.

| Plant | Southern hybridization signal | Sequence |
|-------------------------------|-------------------------------|---|
| <i>Araucaria heterophylla</i> | yes | yes |
| <i>Cycas revoluta</i> | no | yes |
| <i>Ephedra viridis</i> | no | yes |
| <i>Ephedra tweediana</i> | no | no |
| <i>Equisetum hyemale</i> | yes | yes |
| <i>Ginkgo biloba</i> | yes | yes |
| <i>Gnetum gnemon</i> | yes | no |
| <i>Magnolia soulangeana</i> | yes | yes |
| <i>Marchantia polymorpha</i> | yes | only from probable pseudogenes |
| <i>Nymphaea odorata</i> | yes | yes |
| <i>Osmunda claytonia</i> | no | no, but putative RPC1 D-F |
| <i>Psilotum nudum</i> | yes | no |
| <i>Welwitschia mirabilis</i> | yes | no |
| <i>Zea mays</i> | no | yes (identical sequences from cDNA and genomic DNA) |

Figure 5: An 834 bp alignment of the RPBI D-F coding sequence of plants, with *Spirogyra* as the outgroup. Sequences for introns and PCR primers are not shown. The positions of introns 7, 8 and 9 are indicated by a ▼ between the boldfaced nucleotides which flank the intron.

| | | | | | |
|--------------------|-------------|--------------------|------------|------------|-------------|
| | | | | | 50 |
| <i>Spirogyra2</i> | CACGTGTGTC | AGACATTCTGA | GACTCGAGCA | GAAACCATGG | AGCTCATGAT |
| <i>Arau.het-II</i> | CACGTACCGC | AGTCCTTCGA | GACTCGGGGC | GAAGTAGCCG | AGCTCATGCT |
| <i>Cyc.rev-II</i> | CATGTACCCC | AGTCATTTGA | GACAAGAGCA | GAAGTCTTGG | AACTAATGAT |
| <i>Eph.vir-II</i> | CATGTTCCGC | AGTCCTTTGA | GACCAGAGCT | GAAATCCTGG | AGCTGATGAT |
| <i>Gin.bil-II</i> | CATGTGCCCC | AGTCATTTGA | GACAAGGGCA | GAAGTGTTGG | AGTTAATGAT |
| <i>Mag.sou-II</i> | CATGTTCCCTC | AGTCATTTGA | AACTAGAGCA | GAAGTCTTGG | AGCTGATGAT |
| <i>Nym.odo-II</i> | CATGTACCTC | AGTGTTTTGA | AACTAGAGCA | GAAGTCTTGG | AGTTGATGAT |
| <i>Zea mays-2</i> | CATGTCCCCC | AGTCATTTGA | GACCAGGGCA | GAAGTTCTGG | AGTTAATGAT |
| <i>Arab.tha-II</i> | CATGTACCAC | AATCATTCTGA | GACCAGAGCC | GAGGTGTTAG | AGCTGATGAT |
| | | | | | 100 |
| | GGTGCCAAAA | TGTGTCGTGA | GTCCCCAATC | CAACAGGCCT | GTGATGGGCA |
| | GGTGCCGAAG | TGCATCGTCT | CGCCGCAGTC | GAATCGGCCG | GTTATCGGTA |
| | GGTTCCAAAG | TGCATCGTCT | CTCCACAGTC | GAATAGACCG | GTTATGGGTA |
| | GGTACCAAAA | TGTATTGTGT | CCCCTCAGTC | CAATAGGCCT | GTTATGGGTA |
| | GGTGCCAAAG | TGCATTGTTT | CTCCAAAGTC | CAACAGGCCT | GTTATGGGTA |
| | GGTACCGAAA | TGCATTGTCT | CGCCTCAATC | AAATCGCCCT | GTTATGGGTA |
| | GGTACCAAAA | TGTATTGTTT | CACCACAGTC | CAATCGGCCG | GTCATGGGTA |
| | GGTGCCAAAA | TGCATTGTCT | CTCCACAATC | AAATAGGCCT | GTAATGGGTA |
| | GGTTCCTAAA | TGTATTGTCT | CCCCCAGGC | GAATCGTCTT | GTGATGGGAA |
| | | | | | 150 |
| | TCGTGCAGGA | TACGCTACTT | GGCTGCAGGA | AAGTCACGAA | AAGGGACACG |
| | TCGTGCAAGA | CACGCTGCTC | GGGTGCCGGA | AGGTGACGAA | GAGGGATAACA |
| | TTGTCCAAGA | CACTCTTTTG | GGTTGCAGAA | AGATCACAAA | GAGAGACACT |
| | TTGTCCAGGA | TACTCTTTTG | GGTTGCAGAA | AGATCACCAA | AAGAGACACA |
| | TTGTTCAGGA | CACTCTTTTG | GGTTGCAGAA | AGGTGACGAA | AAGGGATAACA |
| | TTGTCCAGGA | TACACTCTTA | GGATGCCGGA | AGATAACTAA | AAGAGATAACC |
| | TTGTGCAGGA | TACACTTCTT | GGATGTAGGA | AGATCACAAA | GCGTGACACC |
| | TTGTCCAAGA | CACACTGCTT | GGGTGTCGCA | AAATTACTAA | AAGGGACACT |
| | TTGTGCAGGA | TACCCTCTTG | GGGTGCCGTA | AAATTACAAA | GAGAGATACT |
| | | ▼ | | | 200 |
| | TTTATCGAGA | AAG ACGTTTT | TATGAACATC | CTCATGTGGT | GGGAAGATTT |
| | TTCATCGAAA | AGG ATGTGTT | CATGAACATT | TTGATGTGGT | GGGATGATTT |
| | TTTATAGAGA | AGG ATGTGTT | CATGAACATC | TTAATGTGGT | GGGAAGATTT |
| | TTTATTGAGA | AGG ATGTTTT | CATGAATATC | TTGATGTGGT | GGGAAGACTT |
| | TTTATAGAGA | AGG ATGTCTT | CATGAATATC | TTAATGTGGT | GGGAGGATTT |
| | TTCATTGAGA | AGG ATGTCTT | TATGAATATC | TTGATGTGGT | GGGAGGATTT |
| | TTCATAGAGA | AGG ATGTTTT | CATGAACATT | CTGATGTGGT | GGGAGGATTT |
| | CTAATTGAAA | AGG ATGTATT | TATGAACATC | TTGATGTGGT | GGCAAGATTT |
| | TTCATAGAGA | AGG ATGTATT | CATGAACACA | CTGATGTGGT | GGGAAGACTT |

250

TGAGGGCAAG ATTCCTTCTC CTACGATTTT GAAGCCTCGT CCCCTTTGGA
TGATGGTAAA ATGCCACACC CAGCGATCCT GAAGCCGAGG CCCATTTGGA
TGATGGCAAA ATACCATCTC CGACTATTCT AAAGCCTAGA CCTCTTTGGA
TGATGGGAAA GTACCTGCAC CTGCAATTTT GAAGCCAAGG CCAATTTGGA
CGATGGCAAA ATACCATCCC CAACAATTCT AAAGCCCAGG CCACTTTGGA
TGATGGGAAA ATACCTGCTC CAACTATTCT GAAGCCTAGA CCTCTTTGGA
TGATGGAAAG ATACCTGCTC CTGCCATTAT GAAGCCTAGA CCTTTGTGGA
CGATGGAAAG ATTCCTGCAC CTACCATTTT GAAACCTAGG CCTATTTGGA
CGATGGGAAA GTTCCGGCTC CTGCAATCTT GAAGCCTCGT CCTCTTTGGA

300

CAGGAAAGCA AGTTTTCTCG TTGATCATCC CCAGAGCTGT GAATCTTGAG
CGGGGAAGCA GATTTTTAGC CTGATCATCC CCAAGCAGAT TAACATGACC
CTGGCAAGCA AGTATTCAAT CTTATCATTC CAAGGCAGAT AAACCTTATA
CTGGCAAGCA AGTTTTTAAT CTTATCATTC CAAAACAAAT AAATCTCGTA
CTGGCAAACA AGTATTCAAT CTTATCATTC CAAGGCAAAT AAATCTAATA
CAGGAAAGCA AGTGTTCAAT CTGATCATTC CTAAGCAGAT AAACCTCATA
CAGGCAAACA AGTATTTAAC CTTATAATTC CAAAGCAGAT AAATTTAATT
CTGGGAAACA AGTTTTTAAC TTAATTATCC CCAAGCAAAT AAATTTAATT
CTGGCAAACA AGTTTTTAAT CTTATCATA CAAAACAGAT AAATCTGTTG

350

CGGTACTCTG CATGGCATCC CGATTCTGAG AGGGGAGACT TCTCCCCAGG
CGAACTGCGG CATGGCACAA CGACAGTGAG AGTACGGACG TAACGCCCGG
AGATACTCTG CATGGCATT CAGAGTCTGAA ACAGGATTTA TCACGCCAGG
AGGTACTCTG CCTGGCATAA TGAAAGTGAT AGGGGACACA TGA CTCTGGG
AGGTACTCTG CATGGCATT CAGAGTCTGAA ACAGGATTTA TCACGCCAGG
AGAACCTCGG CATGGCACTC GGAAGCGGAA ACAGGATTTA TCACTCCAGG
AGGTACTCAG CATGGCATT CAGAGTCTGAA ACAGGATTTA TCACTCCAGG
CGGTTTTTCAG CATGGCATT CAGAGTCTGAA ACAGGATTTA TCACTCCAGG
AGGTACTCTG CTTGGCACGC AGATACAGAG ACTGGATTTA TCACTCCGGG

400

AGACACTCAA GTGCGTGTGG AGAAGGGAGA GTTGCTCGCA GGAATTCTCT
TGATAACCAGC GTTCGCATCG AGAAGGGTGA GCTCCTCACT GGTACGCTCT
GGATACTTGT GTTCGGATTG AAAAGGGAGA AGTTCTGTCT GGCACACTCT
AGATACTGTT GTCAGAATTG AAAAGGGAGA AGTTATAACT GGTACTCTCT
GGATACTTGT GTCCGGATCG AAAAAGGAGA GGTCTTTTCT GGCACACTCT
AGATACACAA GTTAGAATAG AGAGAGGCGA GCTGCTTGCT GGCACACTCT
AGATACTTGT GTCCGAATAG AGAGGGGAGA ACTTCTCTCA GGGACCCTCT
TGATACTATG GTCAGGATAG AGAAGGGAGA GCTTCTGTCT GGCACACTTT
GGATACTCAA GTGCGAATTG AAAGAGGGGA ACTTCTTGCC GGAATCTTT

450

GCAAAAAATC ATTGGGAACC TCCGGTGGAA GTCTCGTTCA CATCATATGG
 GTAAGAAGAC GCTGGGGACG TCCGGTGGTA GTCTCATTCA CGTGATATGG
 GTAAAAAAC CCTTGGAACG TCTTCTGGAA GTCTTATTCA TGTGATCTGG
 GCAAGAAAAC ACTTGGTGCA TCTAGTGGTA GTCTTATTCA TGTGATATGG
 GCAAGAAAAC TCTAGGAACG TCTTCTGGGA GTCTGATTCA TGTGATCTGG
 GCAAGAAGAC CCTTGGACA TCTACTGGTA GTCTTATTCA TGTGATCTGG
 GTAAGAAAAC CATGGGCACA TCATCTGGAA GTCTTATTCA CGTTATCTGG
 GCAAAAAGAG TCTTGGACA GGCTCCGGAA GTCTTATCCA TGTGATTTGG
 GCAAAAAGAC CCTTGGTACA TCTAATGGAA GTCTCGTGCA TGTGATTTGG



500

GAAGAAGTAG GTCCTGATGC TGC GCGGAAG TTTCTGGAC ACACACAATG
GAGGAAGTCG GACCGGACGC TGCCCGTAAG TTCTCGGTC ACACGCAGTG
GAGGAGGTTG GTCCAGATGC AGCTCGCAAG TTTTGGGGC ATACACAGTG
GAGGAGGTTG GTCCAGATGC TGC GCGGAAG TTCTGGGTC ATACGCAATG
GAGGAGTCG GGCCAGATGC TGCTCGTAAG TTTTAGGGC ATACACAGTG
GAAGAGGTTG GTCCAGATGC TGCCCGCAA TTCTGGGCC ACACACAGTG
GAAGAGGTTG GTCCAGATGC TGCACGCAAG TTTTGGGGC ATACTCAGTG
GAAGAGGTTG GTCCAGATGC TGCCCGGAAG TTCTTAGGAC ACACACAGTG
GAAGAGGTTG GTCCTGATGC AGCTAGAAA TTCTCGGTC ATACTCAATG

550

GCTCGTGAAT TATTGGCTTC TGCAGCAGGG ATTTAGCATC GGAATCGGTG
 GCTTGTGAAC TACTGGTTGC TCCAACACGG GTTTAGCATC GGTATCGGTG
 GCTTGTTAAC TATTGGCTAT TGCAACAGGG TTTCAGTATT GGTATAGGAG
 GCTTGTTAAT TACTGGTTGC TACAACAGGG TTTCAGTATG GGCATTGGAG
 GCTTGTCAAC TACTGGCTGT TACAGCAGGG TTTCAGTATT GGTATAGGAG
 GCTTGTTAAC TACTGGCTTT TGCAGAATGG ATTTAGTATT GGAATTGGGG
 GCTTGTTAAT TATTGGCTTC TGCAGAATGC TTTTAGCATT GGTATTGGGG
 GCTTGTAAC TACTGGCTTC TTCAAATGG TTTCAGTATT GGAATTGGGG
 GCTTGTCAAT TACTGGCTTC TGCAGAATGG TTTTACCATC GGAATTGGGTG

600

ATACCATTGC CGATGCATCC ACCATGGATA CCATCAATGA AACCATCGCA
 ATACCATCGC CGATTCGGCT ACCATGGAGA AGATCAATGA GACCATCGCT
 ACACCATTGC TGATGCTGCA ACAATGGAAA CAATTAATGA AACAACTCA
 ATACAATTGC TGATGCCACA ACAATGGACA CGATCAATGA AACAAACAA
 ACACAATTGC TGATGCTGCA ACTATGGAAA CGATTAACGA AACAAATTTCA
 ACACAATTGC AGATGCATCA ACTATGGAAA AAATTAATGA AACAAATATCG
 ACACCATTGC AGATGCTTCT ACCATGGAAA AGATTAATGA GACGATCTCT
 ATACTATTGC AGATGCATCC ACCATGGAAA CAATTAATGA TACAATATCT
 ACACAATTGC CGATTCATCA ACAATGGAGA AAATTAATGA AACTATTTCC

42

650

| | | | | |
|------------|------------|------------|------------|------------|
| AAGGCTAAGA | CTGAAGTGAA | AGATCTCATC | GAAGCAGCTT | GTGAGAAACA |
| AAGGCAAAGC | AGGACGTGAA | GGAAGTATC | AAGGCGTCCC | AAGAAAAGAG |
| AAAGCAAAGG | CTGAAGTCAA | TCAACTTATT | CAGCTTGCTC | ACCAGAAAGC |
| GATGCTAAAA | TTAAAGTGCA | AGAAGTATA | GAAAAGTACA | TGGCACACAA |
| AAAGCAAAGA | ATGAGGTCAA | ACAAGTATT | AAGGCTGCTC | AAGAGAAGGC |
| AAAGCGAAGA | ATGAAGTGAA | GGAAGTATT | AAGGCTGCCC | AAGAAAAGCA |
| AAAGCAAAAA | ATGAGGTGAA | AGAGCTTATC | AAGGCTGCCC | AGGAGAAACA |
| AAAGCTAAGA | ATGCTGTGAA | GGAGCTTATT | AAAAAAGCTC | ATGAGAAGCA |
| AATGCAAAAA | CTGCTGTGAA | AGATCTTATC | CGGCAGTTCC | AGGGAAAGGA |

700

| | | | | |
|------------|------------|------------|------------|------------|
| GTTGGAAGCT | CAACCCGGTC | GTACCCTCAT | GGAGTCCTTT | GAGAATCGTG |
| CCTCGAACCT | CAACCCGGTC | GCACGCTGAT | GGAATCATT | GAGAACAAAG |
| ATTAGAGGCA | GAACCTGGGC | GCACGATGAT | GGAATCATT | GAAAACAGAG |
| GCTAGAACAA | GAGCCAGGTC | GAACGCTATT | AGAATCATT | GAAAATCAAG |
| ATTGGAGGCA | GAACCTGGCC | GTACAATGAT | GGAGTCTTTT | GAAAACAGAG |
| GCTAGAGGCA | GAACCTGGGC | GAAGTATGAT | GGAATCATT | GAGAATAGGG |
| GTTGGAAGCT | GAACCTGGGC | GAACCATGAT | GGAGTCATT | GAGAATAGAG |
| GTTGGAAGCT | GAGCCAGGAC | GCACTATGAT | GGAATCATT | GAAAACAGAG |
| ATTGGACCCT | GAGCCTGGCC | GAAGTATGAG | AGATACATT | GAGAACAGGG |



750

| | | | | |
|---------------------|------------|------------|------------|-------------|
| TCAATCAG GGT | GTTGAACAAG | GCCCGTGACG | ATGCAGGTCG | AGCTGCTCAA |
| TGAACCAG GGT | TCTGAACAGG | GCTCGTGATG | AAGCCGGAAG | TAGTGCTCAG |
| TCAATCAG GGT | GTTGAATAAG | GCCCGTGATG | ATGCAGGAAG | TAGTGCTCAA |
| TCAACCAG GGT | TTTAAATAAG | GCTCGAGATG | ATGCAGGTAA | CAGTGCACAA |
| TGAATCAG GGT | GTTGAACAAG | GCTCGTGATG | ATGCGGGAAG | TAGTGCTCAA |
| TGAACCAG GGT | GCTGAATAAG | GCTCGTGATG | ATGCTGGGAG | TAGTGACACAG |
| TCAATCAG GGT | GTTGAACAAA | GCTCGTGATG | ACGCTGGTAG | TAGTGCTCAG |
| TGAACCAG GGT | TCTTAACAAA | GCCCGTGATG | ATGCTGGTAG | CAGTGCTCAG |
| TTAACCAG GGT | TTTGAATAAA | GCTCGTGATG | ATGCTGGAAG | TAGTGCTCAA |

800

| | | | | |
|------------|------------|------------|------------|------------|
| TCCAGTCTTT | CAGAGAGCAA | CAACGTCAA | GCCATGGTGA | CTGCGGGATC |
| AAGAGTCTAT | CTGAGAGCAA | CAACGTCAAG | GCGATGGTGA | CGGCCGGTTC |
| AAAAGCTTAT | CAGAGAGTAA | TAATTTGAAG | GCAATGGTGA | CTGCTGGCTC |
| AGAAGTTTAT | CTGAGAGCAA | TAATTTGAAA | GCCATGGTGA | CTGCTGGTTC |
| CGAAGCTTAT | CAGAGAGTAA | TAATTTGAAG | GCAATGGTGA | CAGCTGGGTC |
| AAGAGCTTAT | CAGAAAGTAA | CAATCTGAAG | GCTATGGTGA | CAGCAGGATC |
| AGGAGCTTAT | CTGAAAGTAA | TAACCTGAAG | GCTATGGTCA | CTGCAGGATC |
| AATAGCTTGT | CTGAAAGCAA | CAATTTGAAG | GCTATGGTCA | CTGCAGGTTT |
| AAGAGTTTAT | CAGAAACCAA | TAACCTTAA | GCCATGGTGA | CAGCAGGATC |

| | | | |
|------------|------------|------------|------|
| AAAAGGTTCC | TTCATCAACA | TCTCTCAAAT | GATT |
| CAAGGGATCG | TTCATCAACA | TCTCGCAGAT | GATT |
| AAAAGGAAGT | TTTATTAATA | TATCACAGAT | GACT |
| AAAAGGAAGT | TTTATCAATA | TATCACAGAT | GACT |
| AAAAGGAAGT | TTTATTAACA | TATCACAGAT | GACA |
| AAAGGGAAGT | TTTATCAACA | TTTCACAGAT | GACT |
| TAAAGGAAGT | TTTATCAATA | TTTCACAGAT | GACT |
| AAAAGGCAGT | TTCATTAACA | TTTCACAAAT | GACT |
| CAAAGGAAGT | TTCATCAATA | TTTCTCAAAT | GACA |

artefacts (results not shown) obtained by PCR with the region A primer, which corresponds to a zinc finger domain and so it is likely to be similar to sequences found in genes encoding other proteins which also have DNA-binding domains.

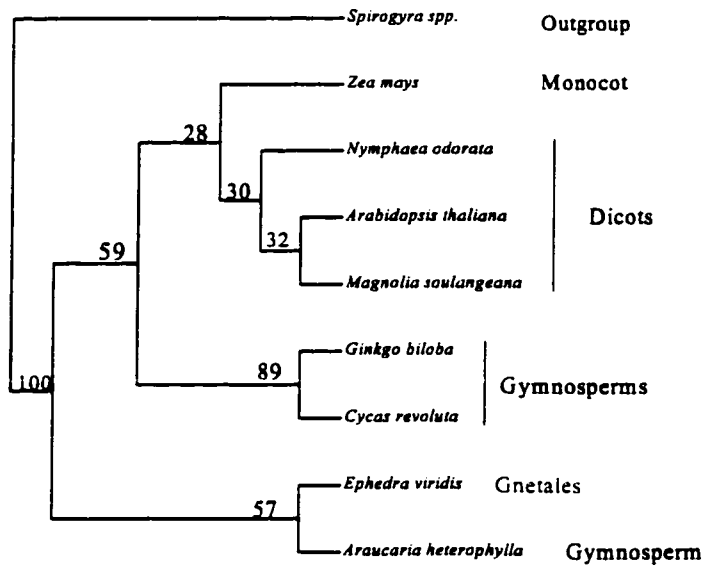
The positions of introns 7, 8 and 9, conserved in *Spirogyra* and *Arabidopsis*, are also conserved in the seven RPBI-coding sequences for regions D-F. This is highlighted in the nucleotide sequence alignment of plant RPBI D-F coding sequences in Figure 5. In the partial sequence of *Equisetum hyemale* RPBI D-F, the positions of introns 7 and 9 are conserved as well (results not shown). Dinucleotide and trinucleotide repeats in these introns in *Equisetum hyemale* RPBI D-F have hampered further sequencing efforts. While intron position is conserved, the intron sequences and lengths vary from species to species amongst plants. Exon sizes, however, are perfectly conserved amongst these plants and *Spirogyra*.

Phylogenetic Analysis of Plant RPBI D-F DNA Sequences:

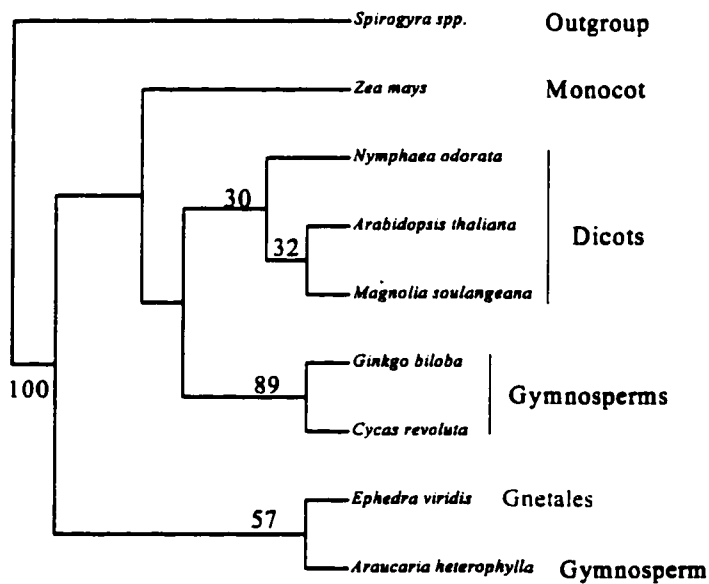
Figures 6, 7 and 8 represent phylogenetic trees inferred from the nucleotide sequences by parsimony, neighbor-joining and maximum-likelihood methods respectively. The alignment used for these analyses was derived from figure 5 and is shown in Appendix 1. The trees in figures 6-8 consistently show *Araucaria*, a Gymnosperm, as one of the most basal taxa. However, two equally parsimonious trees place *Zea mays*, a monocot, as either a sister group to a clade of dicots or to a clade of dicots and gymnosperms (*Ginkgo* and *Cycas*). In both neighbor-joining and parsimony trees, Gnetales (*Ephedra*) are shown as a sister group to a mixed up clade of angiosperms and the *Cycas/Ginkgo* clade. In figure 7, the neighbor-joining tree shows maize as a sister group to a clade of water lilies (*Nymphaea*) and Gymnosperms (*Cycas* and *Ginkgo*). In figure 8, the PUZZLE Schoeniger-von Haeseler model maximum-likelihood tree resolves fewer relationships than the

previous trees did, with a major difference also being that it grouped Gnetales (*Ephedra*) with an Angiosperm (*Arabidopsis*), while the parsimony method had placed Gnetales with a Gymnosperm (*Araucaria*) and neighbor-joining placed Gnetales as a sister clade to angiosperms, *Cycas* and *Ginkgo*, albeit with less statistical support. *Cycas* and *Ginkgo* are consistently grouped together by each of the three methods as well. The phylogenetic trees based upon nucleotide sequences, however, are characterized by inconsistencies and low statistical support, if any, for most of the relationships except for the strong support of the *Cycas* and *Ginkgo* clade.

In all of the parsimony and distance phylogenetic trees shown here, no bootstrap values are indicated on the true trees where relationships in the true trees were not seen in the consensus trees obtained after bootstrapping.



A



B

Figure 6: Two equally parsimonious phylogenetic trees of nucleotide sequences encoding plant RPBI D-F, inferred using DNAPARS in PHYLIP 3.573c. Bootstrap support (% of 1000 replicates) is indicated at the nodes. The third base of codons were removed and the first codon position of the sixfold degenerate arginine and leucine codons were substituted with "m" or "y".

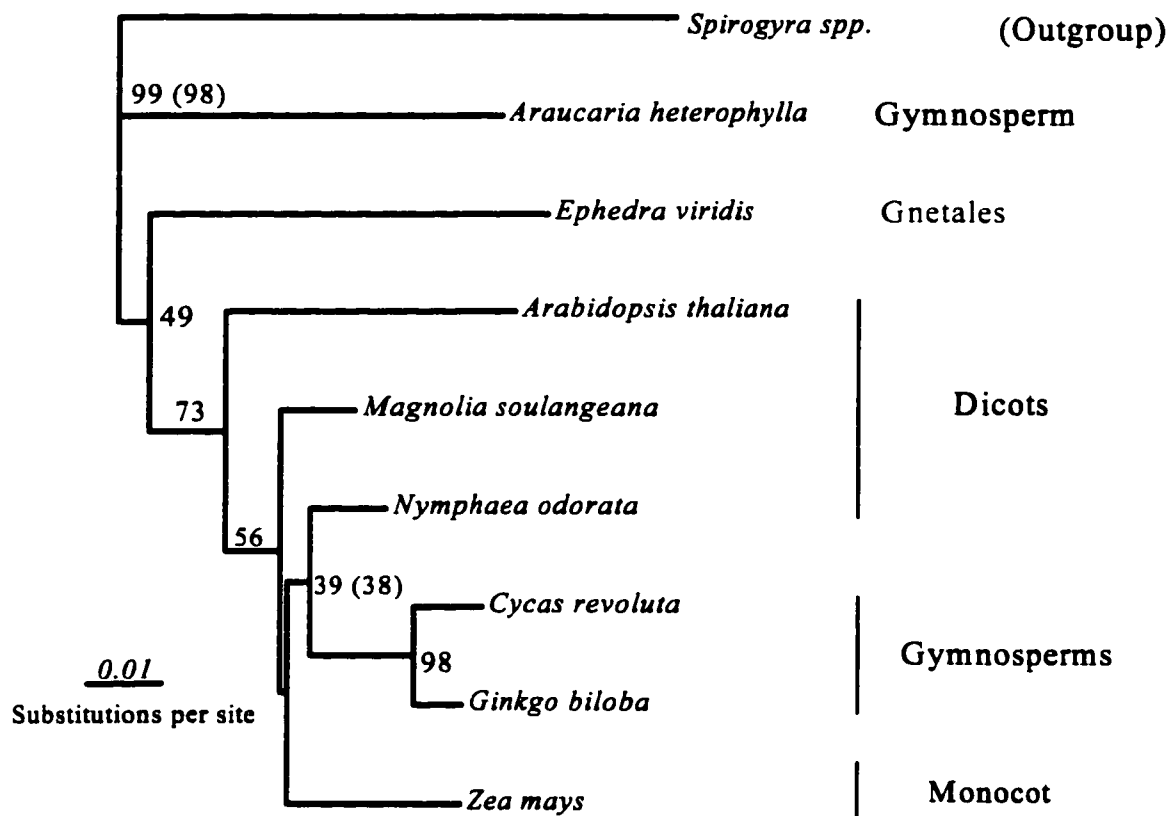


Figure 7: Neighbor-joining phylogenetic tree of nucleotide sequences encoding plant RPB1 D-F. The distances were estimated with DNADIST of PHYLIP 3.573c using the Kimura 2-parameter model of substitution. The % bootstrap values obtained from 1000 replicates are shown at the nodes, with the bootstrap values obtained using the maximum likelihood model of substitution given in parentheses, when different from those obtained by the Kimura model. The third base of codons were removed, and the first base of the sixfold degenerate codons for arginine and leucine were substituted with “m” or “y”

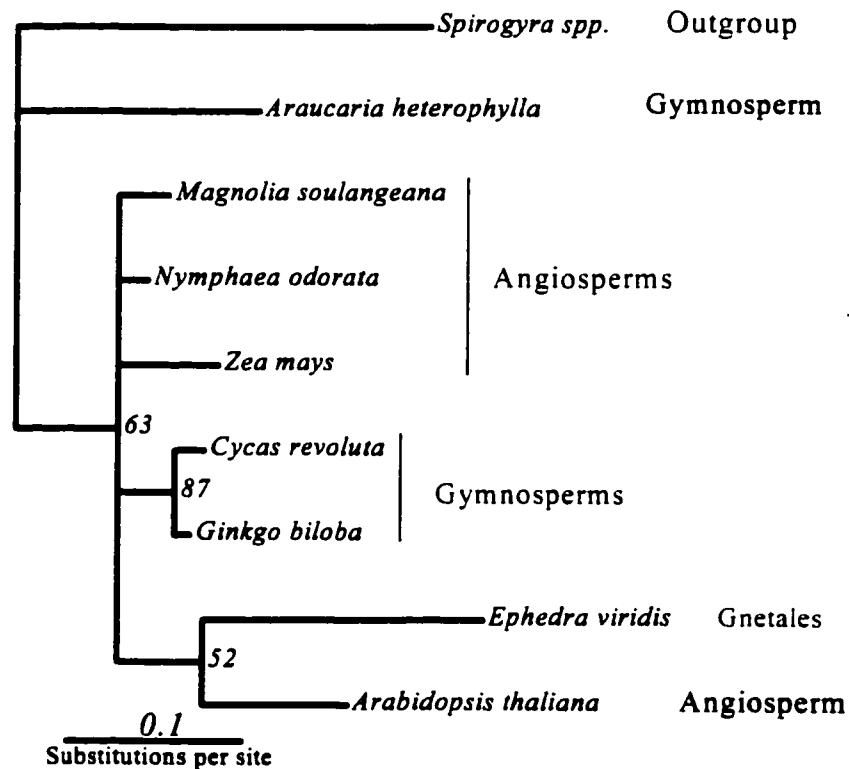


Figure 8: Phylogenetic tree of plant nucleotide sequences encoding RPB1 D-F inferred using maximum likelihood distances with the Schoeniger-von Haeseler model and an invariable and eight γ -distributed rates of substitution, using PUZZLE 4.0.2. The numbers at the nodes indicate the percent of 10000 quartet puzzling steps supporting that relationship, if $>50\%$. The third base of codons were removed and "m" or "y" were substituted for the first base of the sixfold degenerate codons for arginine and leucine. The fraction of invariable sites is 0 (S.E. = 0) and the γ distribution parameter $\alpha = 0.17$ (S.E. 0.03). The log L = -1392.46 for the tree with nonclocklike branchlengths shown here.

Phylogenetic Analysis of Plant RPBI D-F Amino Acid Sequences:

An alignment of the amino acid sequences inferred from the nucleotide sequence alignment in Figure 5 was also used for phylogenetic analysis by the parsimony, neighbor-joining and maximum likelihood methods. The boldfaced section of Figure 9, when considered without any gaps or deletions, represents the amino acid alignment used to infer the phylogenetic trees shown in Figures 10, 11, 12 and 13, with *Spirogyra*, a green alga, as the outgroup (see Appendix 2). Figures 10 and 11 were inferred by the parsimony and neighbor-joining methods, respectively, while Figures 12 and 13 are both maximum likelihood trees, inferred using the Blosum62 and JTT models, respectively.

Four equally parsimonious trees are found in Figure 10. These trees consistently group *Cycas* and *Ginkgo* together as sister groups, with strong bootstrap support. Similarly, these trees also consistently place *Araucaria* as the most basal plant and Gnetales (*Ephedra*) as a sister group to angiosperms and the *Cycas/Ginkgo* clade, albeit with low bootstrap support. In three of the four equally parsimonious trees, *Ephedra* is nestled amongst the gymnosperms, and the *Cycas/Ginkgo* clade is a sister group to angiosperms. The parsimony trees are inconsistent about the relationships of angiosperms amongst one another, with monocot, herbaceous or eudicot origins of angiosperms shown as being equally likely.

The neighbor-joining tree in Figure 11 shows some relationships similar to those in the parsimony analysis. *Cycas* and *Ginkgo* are grouped together as sister groups, with strong statistical support. *Araucaria* is again placed as the most basal plant, and Gnetales (*Ephedra*) as a sister group to angiosperms and the *Cycas/Ginkgo* clade, with strong statistical support. Relationships amongst angiosperms remain unresolved, except that *Arabidopsis* is placed as basal to both angiosperms and

Figure 9: 557 position amino acid alignment of 38 D-F sequences from type A RNA polymerases, with a single Archaeal sequence from *Sulfolobus acidocaldarius* as the outgroup to sequences for eukaryotic RNA polymerases I, II and III. Residues identical to that of the outgroup are shown as points (.), while gaps are shown as a dash (-). Plant sequences are boldfaced. Sequences corresponding to my PCR primers for regions D and F are not shown. Positions of the alignment which were included for phylogenetic analysis of animals, plants and fungi are identified by a box of dotted lines, while positions included for phylogenetic analyses of all of the taxa in the alignment are identified by solid boxes. No gaps or deletions were introduced into the alignment for analyses of plant taxa alone. Residues flanking the three intron positions in plants are underlined. Appendices 2-4 show the alignments used for phylogenetic analysis.

101

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---------------------------|----------|-------------|------------|-----------|--------------|------------|-----------|-------------|-------------|-------------|-------------|-------------|--------------|--------------|-------------|-------------|-------------|------------|-------------|------------|------------|--------|------|------------|-----|-------|-----|-------|-----|-----|-----|-----|
| <i>Sulf. acido</i> | NLHV | PQSE | EAI | AEARE | LMLVH | KNIIT | PRYGG | PIIGG | QDQYI | SGAY | LLSVK | TKLLT | VEEVAT | ILG | VTD | --- | | | | | | | | | | | | | | | | |
| <i>Por. yez-II</i> | | THATR | ..VM | ..M.P | RC.VS | QGNK | VM | IV | ..TL | L.CM | FTYRD | F | RRDVTMS | L | LHVEG | --- | | | | | | | | | | | | | | | | |
| <i>Bon. ham-II</i> | | THQTR | ..VQ |P | HC.VS | QGNK | VM | IV | ..TL | L.CM | FTQRD | F | ERDLMN | L | MHVG | --- | | | | | | | | | | | | | | | | |
| <i>Dro. mel-II</i> | | M | TR | VENI | HITP | RQ |Q | ANR | VM | IV | ..TL | TAVR | KMTKRDVFI | R | Q.MN | L.MFLPT | | | | | | | | | | | | | | | | |
| <i>Art. sal-II</i> | | A | L | TR | LENI | HITP | RQ |Q | ANR | VM | IV | ..TL | CAVR | KMTKRDVF | EK | QMM | L.MYLPT | | | | | | | | | | | | | | | |
| <i>Cra. gig-II</i> | | L | L | TK | ISN | A | ..P | RM |Q | ANR | VM | IV | ..TL | TAVR | KMTKRDVF | DRGX | MN | L.MFLPR | | | | | | | | | | | | | | |
| <i>Ily. obs-II</i> | | L | L | TR | IMN | CA | ..P | RM |Q | ANR | VM | IV | ..TL | TAVR | KMTKRDVF | RAQ | MMH | L.MFLPT | | | | | | | | | | | | | | |
| <i>Hel. sta-II</i> | | L | A | L | TR | ISQ | AS | K | RM |Q | ANR | VM | IV | ..SL | TAVN | KMTRRD | FI | KD | IMN | ..MYLPC | | | | | | | | | | | | |
| <i>Cae. ele-II</i> | | L | L | L | TR | IE | IAM | ..P | RQL |Q | ANR | VM | IV | ..TL | CAVR | MMTKRDVF | IDW | PFMM | D | L.MYLPT | | | | | | | | | | | | |
| <i>Mus. mus-II</i> | | L | L | L | TR | IQ | AM | ..P | RM | ..V | QSNR | VM | IV | ..TL | TAVR | KFTKRDVF | ERG | ..MN | L | MFLST | --- | | | | | | | | | | | |
| <i>Spirogy. II</i> | M | C | TF | TR | TM | ..M | P | CWVS | QSNR | VM | IV | ..TL | L | CRKVT | KRD | FIEK | QVFN | ..M | WED | --- | --- | | | | | | | | | | | |
| <i>Arau. he-II</i> | M | ..F | TRG | VA |P | C | VS | QSNR | V | ..IV | ..TL | L | CRKVT | KRD | FIEK | QVFN | ..M | W | D | --- | --- | | | | | | | | | | | |
| <i>Cyc. rev-II</i> | M | ..F | TR | VL |M | P | C | VS | QSNR | VM | IV | ..TL | L | CRKVT | KRD | FIEK | QVFN | ..M | WED | --- | --- | | | | | | | | | | | |
| <i>Eph. vir-II</i> | M | ..F | TR | IL |M | P | C | VS | QSNR | VM | IV | ..TL | L | CRKVT | KRD | FIEK | QVFN | ..M | WED | --- | --- | | | | | | | | | | | |
| <i>Gin. bil-II</i> | M | ..F | TR | VL |M | P | C | VS | KSNR | VM | IV | ..TL | L | CRKVT | KRD | FIEK | QVFN | ..M | WED | --- | --- | | | | | | | | | | | |
| <i>Mag. sou-II</i> | M | ..F | TR | VL |M | P | C | VS | QSNR | VM | IV | ..TL | L | CRKVT | KRD | FIEK | QVFN | ..M | WED | --- | --- | | | | | | | | | | | |
| <i>Nym. odo-II</i> | M | ..CF | TR | VL |M | P | C | VS | QSNR | VM | IV | ..TL | L | CRKVT | KRD | FIEK | QVFN | ..M | WED | --- | --- | | | | | | | | | | | |
| <i>Zeemays-II</i> | M | ..F | TR | VL |M | P | C | VS | QSNR | VM | IV | ..TL | L | CRKVT | KRD | IEK | QVFN | ..M | WQD | --- | --- | | | | | | | | | | | |
| <i>Arab. th-II</i> | M | ..F | TR | VL |M | P | C | VS | QANR | VM | IV | ..TL | L | CRKVT | KRD | FIEK | QVFN | T | MWED | --- | --- | | | | | | | | | | | |
| <i>Sac. cer-II</i> | | TR | LSQ | CA | ..P | LQ | VS | QSNK | CM | IV | ..TL | C | IRK | TLRD | FIELDQ | LN | M | YWVPI | --- | --- | --- | --- | | | | | | | | | | |
| <i>Sch. pom-II</i> | | TR | IQ | ITM | ..P | Q | VS | QSNK | VM | IV | ..TL | A | VRKF | LRDNF | ..RNA | MN | MLWVPI | --- | --- | --- | --- | | | | | | | | | | | |
| <i>Dic. dis-II</i> | | TL | TR | VI | ..I | M | ..P | RQ | VS | QSNR | VM | IV | ..TL | L | SR | FTKRDC | FMEKDL | ..M | W | --- | --- | --- | | | | | | | | | | |
| <i>Aca. cas-II</i> | M | ..TPG | ..R | ..VI | ..M | ..P | Q | V | AQSNK | V | ..IV | ..TL | L | GC | ..TQRD | FIEKDVMMN | ..MWLES | --- | --- | --- | --- | | | | | | | | | | | |
| <i>Nos. loc-II</i> | ..M | ..YNSK | ..LS |S | H | ..S | QSNK | VM | II | ..TL | L | VCR | TS | GVFIKR | ..FCN | L | VYASNI | --- | --- | --- | --- | | | | | | | | | | | |
| <i>Vai. nec-II</i> | ..M | ..YNSK | ..LE | ..C | ..S | QVLS | QSNK | VM | IV | ..SL | TALR | FTLRD | SFFDRR | ..TMQ | L | YSVNI | INNYEFTDSSK | L | IM | THDDS | F | GNL | LHTE | --- | | | | | | | | |
| <i>Pla. fal-II</i> | ..LA | ..H | ..TRS | ..IKH | ..I | ..Q | ..RQ | ..VS | QGNK | VM | IV | ..SL | LAIRK | FTRRDNF | ..K | ..MS | L | ..IWIPX | --- | --- | --- | --- | --- | --- | | | | | | | | |
| <i>Mas. inv-II</i> | M | ..N | ..L | ..S | VKN | ..A | ..P | ..FQ | ..V | ..QKNS | CM | ..VV | ..SL | L | ..CS | ..I | ..RRD | ..F | ..EDVMMN | LAL | ..ISYD | --- | --- | --- | | | | | | | | |
| <i>Tri. vag-II</i> | | QT | RT | V | HI | ..A | ..P | ..Q | ..S | QANK | V | ..LV | ..AL | L | ..CR | ..F | ..RNQ | ..MN | LMWLPT | --- | --- | --- | --- | --- | | | | | | | | |
| <i>Try. bru-II</i> | | LLTK | ..LI | ..M | ..M | ..P | ..FVS | NKSA | CM | IV | ..SL | L | ..S | ..R | ..TD | ..D | ..F | DKYF | QS | VALWL | ..I | --- | --- | --- | | | | | | | | |
| <i>Gia. lam-II</i> | ..AL | ..L | ..V | ..IS | ..CM | ..S | NQVSVK | DNR | ..VYIV | ..VL | L | ..C | ..FTG | DVTIP | PFAR | ..CE | YIMM | GF | SRTSDPT | YAHK | LQ | RYDKGS | VNQ | RHSSASMSYD | --- | | | | | | | |
| <i>Sac. ce-III</i> | | T | ..R | ..IN | ..G | ..K | ..N | ..LL | ..KS | ..E | ..AAT | ..F | ..T | ..S | ..I | ..H | DSFY | D | RATL | TQ | L | SMMSIG | --- | --- | | | | | | | | |
| <i>Gia. la-III</i> | IFYM | ..GQ | ..R | ..GI | ..GS | E | ..S | ..H | ECM | ..LT | ..FL | ..T | ..I | ..G | ..GIEM | RQ | YQC | HVS | YGC | GF | GF | ATY | GV | SLYN | Y | REFIK | SI | HDRRG | TEE | --- | | |
| <i>Try. br-III</i> | V | ..FV | ..T | ..K | ..R | ..LQ | ..STA | R | ..SAKN | E | ..ACT | ..FL | ..AA | ..V | TSRD | VFFDRG | ..FSQ | MV | SHWLQ | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | | |
| <i>Mus. musc-I</i> | A | ..F | ..LGR | ..YV | ACTD | QYLV | KD | ..Q | ..LA | ..LI | ..HM | VSGAN | MTIRG | ..CFF | ..R | ..QYME | L | VYRGL | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | | |
| <i>Dro. mela-I</i> | A | ..Y | ..V | ..R | ..YN | ..VN | ..A | ..S | YLV | KD | ..T | ..LG | ..LI | ..HV | ISGVK | ..IRGR | FFNR | ..DYQQ | L | VFOGL | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | |
| <i>Try. bruc-I</i> | V | ..I | ..TR | ..VET | ..DAN | ..I | ..YLV | ..TS | ..R | ..LI | ..HV | AAGV | ..VTLRDK | FFD | ..STFVQ | LVYNG | V | GPYIQEN | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | |
| <i>Sac. cere-I</i> | M | ..F | ..N | ..N | ..R | ..LN | ..ANTD | ..SQYL | ..TS | ..S | ..VR | ..LI | ..H | ..AGVW | ..TS | ..DSFF | ..R | ..QYQQ | YI | YG | CI | RI | PE | --- | --- | --- | --- | --- | --- | --- | --- | |
| <i>Sch. pomb-I</i> | M | ..F | ..TN | ..RS | ..QFI | ..ANTD | ..SQYL | ..TS | ..D | ..LR | ..LI | ..HV | VMGVW | ..TC | ..D | ..FYIRD | ..YQQ | L | ..FQALP | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

| | | | | | | | | | |
|-------------|-----|--------------|---------------------------|-----|--------|-----|----------------------|-----|--------------------|
| Sulf. acido | --- | FHEHPANISKG | PRACKD | --- | EICP | --- | DSFVIKNG | --- | LLLEGVFDKKAIGNC |
| Por. yez-II | --- | VNLVRYC | THPD DEST | --- | D.S | --- | TRVL.VG. | --- | IT.IV.RTV.SA |
| Bon.ham-II | --- | VNI IRENI | THPD DEKT | --- | D.S | --- | TKVY.SR. | --- | IC.IV.RTV.SA |
| Dro.mel-II | --- | VNMIRTH | THPD EEDEGP | --- | YKW.S | --- | TKVMVEH. | --- | IM.ILC.SL.TS |
| Art.sal-II | --- | VNMKTH | THPD DEDDGP | --- | YKW.S | --- | TKVMVEH. | --- | IM.ILC.TL.AX |
| Cra.gig-II | --- | TNCIRTH | THPD EEDKGP | --- | YKW.S | --- | TKVL.ED. | --- | IS.ILC.TL.TS |
| Ily.obs-II | --- | VNVIRTH | THPD GEDSGP | --- | YKW.S | --- | TKVL.ED. | --- | IS.ILC.TL.TS |
| Hel.sta-II | --- | INVIRTH | THPD DEDRGP | --- | HKW.S | --- | TKVLVED. | --- | S.ILC.SL.AS |
| Cae.ele-II | --- | VNVLRTH | THPD SEDSGP | --- | YKW.S | --- | TKVI.EH. | --- | S.IVCS.TV.KS |
| Mus.mus-II | --- | INCIRTH | THPD DEDSGP | --- | YKH.S | --- | TKV.VE.. | --- | IM.ILC.SL.TS |
| Spirogy. II | --- | VNLERSAW | HPD SE | --- | RGDFS | --- | TQVRVEK. | --- | A.ILC.SL.TS |
| Arau.he-II | --- | INMTRT | AHND SE | --- | STDVT | --- | TSVR.EK. | --- | T.TLC.TL.TS |
| Cyc.rev-II | --- | INLIRYSAW | HE SE | --- | TGF.T | --- | TCVR.EK. | --- | V.S.TLC.TL.TS |
| Eph.vir-II | --- | INLIRYSAW | HE SD | --- | RGHMT | --- | TVVR.EK. | --- | VIT.TLC.TL.AS |
| Gin.bil-II | --- | INLIRYSAW | HE SE | --- | TGF.T | --- | TCVR.EK. | --- | V.S.TLC.TL.TS |
| Mag.sou-II | --- | INLIRYSAW | HE AE | --- | TGF.T | --- | TCVR.EK. | --- | A.TLC.TL.TS |
| Nym.odo-II | --- | INLIRYSAW | HE SE | --- | TGF.T | --- | TCVR.EK. | --- | S.TLC.TM.TS |
| Zeamays-II | --- | INLIRYSAW | HE EE | --- | KGF.T | --- | TCVR.EK. | --- | S.TLC.SL.TS |
| Arab.th-II | --- | INLIRYSAW | HAD TE | --- | TGF.T | --- | TCVR.EK. | --- | A.TLC.TL.TS |
| Sac.cer-II | --- | I.IQREFD | ---E. TT | --- | LLS | --- | NGML.ID. | --- | IIF.VE.TV.SS |
| Sch.pom-II | --- | INLIRDD | ---DK QS | --- | LSN | --- | GML.E.. | --- | IIV.V.TV.AS |
| Dic.dis-II | --- | INLIRPT | STHND KKP | --- | CSA | --- | PRVI.ERS. | --- | A.ILC.RSL.AA |
| Aca.cas-II | --- | TN.VNAD | ---D EEP | --- | DMSF | --- | TKVL.EE. | --- | VS.ILN.TL.TS |
| Nos.loc-II | --- | IYNGKSN | ---EH NEE | --- | KLNFF | --- | KV..MD. | --- | N..IC.SV.TA |
| Vai.nec-II | --- | FKRFNMVK | INLMRDSST.SK DDNP | --- | DLENV | --- | YVI.R.. | --- | I.S.II...V.ST |
| Pla.fal-II | --- | NNNNNNNN | NNNIGGGINS | --- | YCSI | --- | N.GKVI..N | --- | S.IIC.RTV.SS |
| Mas.inv-II | --- | KCG.Y.NAEQSV | RE VA.AQE | --- | EYMNSL | --- | .K.C | --- | HT..ITN..V.KS |
| Tri.vag-II | --- | ISHNSYSQ | DANK MDQN | --- | S.PL | --- | A.RHVI.RD. | --- | A.ILG.TVARS |
| Try.bru-II | --- | EVNHPATP | ---QD RPP | --- | FPH | --- | N.VVM.RR. | --- | C.PIT.SIV.AA |
| Gia.lam-II | --- | FYRSGDPK | DISA-HAD K | --- | Y.SFLJ | --- | VMCGRIA | --- | VVG.-----TESA |
| Sac.ce-III | --- | NSPVINLD | AKNKVFVP .KSKSL | --- | PN.MSQ | --- | N.G.VI.RGS | --- | I.S.M.SVL.DGMKH |
| Gia.la-III | --- | FDSVINLE | HGDKTY.K DSDR | --- | RALSV | --- | N.DY.I.Q.S | --- | H.V.RLT.TFLASSK-NC |
| Try.br-III | --- | EVDVLLS | SFEAPTK FYTR.G | --- | KHD.A | --- | REGYVAFDLS | --- | FIS.RL..LL.GGAKDG |
| Mus.musc-I | --- | EDYAPLNL | S.K.K.GSK AWVKEKPRP-IPDFU | --- | PDS | --- | M.C.QVI.RE. | --- | C.L..AHY.SS |
| Dro.mela-I | --- | EGYERINL | DSF.K.AGK NNNVSRPRPPICGTN | --- | PEGNDL | --- | S.E.QVQ.R..S | --- | V.L..QQY.AT |
| Try.bruc-I | --- | REIEGGITL | K.TSQ.QPS AFDRIPAG | --- | SC | --- | DAVRAKSGAVU..TVMFA.S | --- | IT.FMC.QL.AS |
| Sac.cere-I | --- | PDMPGINL | ISKNK.KNE YWG | --- | S | --- | EENEVLF.D. | --- | C.II..SQY.AS |
| Sch.pomb-I | --- | SDRPGNL | LKSK.KVPGK YWS | --- | S | --- | REGSVLFD. | --- | C.II..SSF.AS |

| | | | | | |
|-------------|-----------------------|---------------------|--------------------|-----------------------|------------------------|
| Sulf. acido | NLHWSIREYGTGYGKWLMDN | VFKMEIRLEHR-GFTMTLE | DITIPDEAQNEITTKIKE | GYSQVDEYIRKFNEG | OLEPIPG |
| Por.yez-II | LI.VTWK.K.P.RTCV.ISA | IQVLVNHVYVIM | .QSIGIG | .TIADAHTDANVRAT.TG | E.TLL |
| Bon.ham-II | LI.ITWK.F.PKITDT.ISQ | IQVLVNHVYILO | .RQSIGIG | .TIAD.ATMRNVIDT.QG | E.VLL |
| Dro.mel-II | L.ICFL.L.HDIAGREYF. | IQTVINNW.LFE | .HSIGIG | .TIADPQY...QQA.K | E...T |
| Art.sal-II | I.IIFL.L.HDICGKFXG. | IQTVVNNW.LYB | .HSIGIG | .TIADPQY.S.Q.T.K | E...T |
| Cra.gig-II | LA.VVFM...WVIAGEMYGH | IQTLVNNW.LLBT | .HSIGIG | .TIADPQYID.QDT.K | D...T |
| Ily.obs-II | LV.IVFL.M.F.VAGE.YG. | IQTVVNNW.LLBT | .HSIGIG | .TIADQQTYH.QET.RK | E...T |
| Hel.sta-II | LQ.IIHH.L.SDATADEFAY | IQMVTNHW.LVTF | .H.IGIA | .TIADAKTYSD.Q.A.K | E...M |
| Cae.ele-II | L.VVTL.L.Y.IAANFYSH | IQTVINAW.IRE | .H.IGIG | .TIADQATYLD.QNT.RK | D...T |
| Mus.mus-II | LV.I.YL.M.HDITRLEFYS. | IQTVINNW.LIE | .H.IGIG | .SIADSKTYQD.QNT.K | E...T |
| Spirogy. II | LV.IIWE.V.PDAARKFELGH | QTMLVNYW.LQD | .SIGIG | .TIADASTMDT.NET.AK | ...AQ |
| Arau.he-II | LI.VIME.V.PDAARKFELGH | QTMLVNYW.LQH | .SIGIG | .TIADSAITMEK.NET.AK | S...Q |
| Cyc.rev-II | LI.VIME.V.PDAARKFELGH | QTMLVNYW.LQD | .SIGIG | .TIADAAITMET.NET.SK | A.AE |
| Eph.vir-II | LI.VIME.V.PDAARKFELGH | QTMLVNYW.LQD | .S.GIG | .TIADATMDT.NET.QD | K.QE |
| Gin.bil-II | LI.VIME.V.PDAARKFELGH | QTMLVNYW.LQD | .SIGIG | .TIADAAITMET.NET.SK | A.AE |
| Mag.sou-II | LI.VIME.V.PDAARKFELGH | QTMLVNYW.LQD | .SIGIG | .TIADASTMEK.NET.SK | ...AE |
| Nym.odo-II | LI.VIME.V.PDAARKFELGH | QTMLVNYW.LQD | .A.SIGIG | .TIADASTMEK.NET.SK | ...AE |
| Zeemays-II | LI.VIME.V.PDAARKFELGH | QTMLVNYW.LQD | .SIGIG | .TIADASTMET.NDT.SK | ...AE |
| Arab.th-II | LV.VIME.V.PDAARKFELGH | QTMLVNYW.LQD | .IGIG | .TIADSSITMEK.NET.SN | E.D.E |
| Sac.cer-II | LI.VVT...K.PQVCAK.FG. | IQ.VVNF.W.L.N | .STGIG | .TIADGPTMR...ET.A. | L.TAKH |
| Sch.pom-II | LV.TIWK.K.P.IC.GFENG | IQRVVNYW.L.N | .SIGIG | .TIADADTMK.V.RTV... | R.K.E. |
| Dic.dis-II | LT.VVMN.H.HDTCRLEFI.Q | QTVVNHM.IN | .GIG | .TIADSAITMAKV.LT.SS | F.CQ |
| Aca.cas-II | LV.VIWN.H.S.VC.HFLNQ | .QHVVNYW.L.H | .SVGVG | .TIAD.TLAK..QT.RK | ...RQ |
| Nos.loc-II | LI.IIAND..H.EITRFI.S | LQ.LISTY.TLIST | .SVGIG | .TISSP.TMAH.SRA.GD | G |
| Vai.nec-II | LI.IIANDF.PDRVTCEF.D | AQ..MNLYFATINA | .SIGIG | .AIADK.TMSQVORS.ET | S.KL |
| Pla.fal-II | LI.VLWH.M.PDKT.DFLSA | LQ.VTNNW..YV | ...VSCS | ..IASNKVLGKVREILDK | K..RL |
| Mas.inv-II | II.ILWKDQ.PMAARDFLSR | .QLLTNAYILT | ...SVGT | .TLADPDTLQAVKAE.EG | E..CQ |
| Tri.vag-II | LI.VV.NS.N.NIA.DFLNQ | QTLIVNNW | ...SIG.I | .CVV..FVLQ.VKHE.DD | R.KVQA |
| Try.bru-II | LI.VIFN.H.SDEVARFNG | .QRVTFE.LNF | ...SVGVQ | .TVADSDTLRQMNDVIVK | ...RM |
| Gia.lam-II | LI.ILF.D..I.PCRAFI.. | QORVVC..MLDH | ...SVGMG | .MVSSEHTERKVAEIQTK | T.NRKA |
| Sac.ce-III | VEYFIL.D..POEAANA.NR | MA.LCA..GN | ...SIGIN | .V.PA.DLKQKKEELVEI | KIKLA.. |
| Gia.la-III | IFYFLVQN..PVSAAARI.LR | FA.VAA...MNY | ...IGID | .VMPSQRVLGKKEVIVQ. | E..TQ |
| Try.br-III | LFARLHTIA.GG.TARV.SR | IAQFTS.Y.TNY | ...SLG.G | .VAPT.P.LNKQKAAVLAR | K.A.A |
| Mus.musc-I | LV.CCYEI..G.TSGRVLTC | LARL.TAY.QLYR | ...LGV. | ..LVKPN.DVVRQRI.E | RMI.L |
| Dro.mela-I | LI.CMYEL..GDVSTL.VTA | FT.V.TF..QL.E | ...LGVK | ..LVT.V.DRKRRI.R | --EIQGKWQDAHLSKQDRDF |
| Try.bruc-I | AP.HVYEL..PHRTGQ.FAA | FGRVLLLA.RKE | ..LSLAMD | .MFLV...ERRCDLLRK | --ELVEKMEAAAYV--KDSKFR |
| Sac.cere-I | IV.SLHEV..P.VAAKVLSV | LGR.L.TNYITA.TA | .CGMD | .LRLTA.GNKWR.DIL.TSV | DPELLKRLQEI LR--DNKKS |
| Sch.pomb-I | LV.SVHEL..PDIAGR.LSV | LSRL.TAYAQM.R | ...CRMD | .LRLDEQGD.WRRQLLENSEK | --LLNANLEEVYR--DDEKL |

| | | |
|--------------------|---|-------------------------------|
| <i>Sulf. acido</i> | RTIEESLESYILDTLDKLRK VAGEIATKYLDPPFNN | VYIM AITGARGSELNITOMT |
| <i>Por. yez-II</i> | KSMM. F. VEVNKV. NGA. D TS. SS. QLS. LKS. . . | IKR. VSA. SK. FI. . S. IC |
| <i>Bon. ham-II</i> | KGMM. F. TEVNKV. NGA. D KS. AS. QRS. LKS. . . | IKR. VSA. SK. FI. . S. IC |
| <i>Dro. mel-II</i> | N. LRQTF. NKNRI. NDAHD KT. GS. K. S. TEY. . . | LKA. VVS. SK. NI. . S. VI |
| <i>Art. sal-II</i> | N. LRQTF. NQVRI. NDA. D KT. GS. KNS. TEY. . . | LKA. VVS. SK. NI. . S. VI |
| <i>Cra. gig-II</i> | N. LRQTF. NMVRI. NDA. D KT. SK. Q. S. SDY. . . | FKA. VVA. SK. KI. . S. VI |
| <i>Ily. obs-II</i> | N. LRQTF. NQVRI. NDA. D KT. SK. Q. S. SE. . . | FKA. VVA. SK. KI. . S. VI |
| <i>Hel. sta-II</i> | N. LRQTF. NQVRI. NDA. D KT. SL. Q. S. SE. . . | FKS. VVA. SK. NKI. . S. VI |
| <i>Cae. ele-II</i> | N. LRQTF. NKNQI. NDA. D RT. SS. Q. S. SE. . . | FKS. VVS. SK. KI. . S. VI |
| <i>Mus. mus-II</i> | N. LRQTF. NQVRI. NDA. D KT. SS. Q. S. SEY. . . | FKS. VVS. K. KI. . S. VI |
| Spirogy. II | . LM. F. NRVNOV. N. A. D D. RA. QSS. SES. . . | . KA. VTA. SK. FI. S. I |
| Arau. he-II | . LM. F. NKNQV. NRA. D E. SS. Q. S. SES. . . | . KA. VTA. SK. FI. S. I |
| Cyc. rev-II | . MM. F. NRVNOV. N. A. D D. SS. Q. S. SES. . . | LKA. VTA. SK. FI. S. . . |
| Eph. vir-II | . LL. F. NQVNOV. N. A. D D. NS. QRS. SES. . . | LKA. VTA. SK. FI. S. . . |
| Gin. bil-II | . MM. F. NRVNOV. N. A. D D. SS. QRS. SES. . . | LKA. VTA. SK. FI. S. . . |
| Mag. sou-II | . MM. F. NRVNOV. N. A. D D. SS. Q. S. SES. . . | LKA. VTA. SK. FI. S. . . |
| Nym. odo-II | . MM. F. NRVNOV. N. A. D D. SS. QRS. SES. . . | LKA. VTA. SK. FI. S. . . |
| Zeanays-II | . MM. F. NRVNOV. N. A. D D. SS. QNS. SES. . . | LKA. VTA. SK. FI. S. . . |
| Arab. th-II | . MRDTF. NRVNOV. N. A. D D. SS. Q. S. AET. . . | LKA. VTA. SK. FI. S. . . |
| <i>Sac. cer-II</i> | M. LR. F. DNVRV. NEA. D K. RL. EVN. KDL. . . | . KQ. VMA. SK. FI. A. . S |
| <i>Sch. pom-II</i> | M. LR. F. AKVSRI. NQA. D N. RS. EHS. KDS. . . | . KQ. VAA. SK. FI. S. . S |
| <i>Dic. dis-II</i> | KSVI. FF. QKYNQV. N. A. D D. T. SS. QDS. SED. . . | LKA. VTA. SK. FI. S. . M |
| <i>Aca. cas-II</i> | . MM. F. FV. NQI. N. A. D D. NS. Q. S. RRS. . . | FKA. V. A. SK. AI. . S. VL |
| <i>Nos. loc-II</i> | MNLQ. T. . K. NLA. N. A. D IS. TR. VES. NHL. G. . . | LKQ. LKA. SK. YI. . S. I. |
| <i>Vai. nec-II</i> | MSMR. F. QVNYI. N. P. D IS. AS. S. S. SFC. . . | MRT. VLA. SK. FI. . S. V. |
| <i>Pla. fal-II</i> | KSLY. F. TRVNNE. NCA. E. M. KV. SES. . ER. . . | IFS. VAS. SK. II. . S. II |
| <i>Mas. inv-II</i> | . SLV. F. AKTNKS. QDALS N. KKSLS. KYD. . . | FKL. IES. SK. . M. C. I. |
| <i>Tri. vag-II</i> | MSYMQ. F. TEVNLTNEILS KTYKVINAKIRGD. S. . . | LSE. LSA. SK. ADT. MS. II |
| <i>Try. bru-II</i> | M. LLQ. F. ADVNSA. N. C. E. E. AKK. LSNVRR. T. S. . . | FKV. IEA. SK. TD. . C. IA |
| <i>Gia. lam-II</i> | QSRNDAF. QEVISKVSGTSL ALEKVI. DAAPHR. A. . . | LLV. INA. SK. KKF. MM. IS |
| <i>Sac. ce-III</i> | CNE. QT. . AK. GGL. S. V. E. EV. DVCINE. . NW. A. . . | PL. . . TC. SK. T. . VS. . V |
| <i>Gia. la-III</i> | S. VQ. T. . ATLNQI. SNV. E. SCAQ. L. E. HFT. K. . . | PL. . . SLC. SK. PI. . A. . I |
| <i>Try. br-III</i> | L. VKQ. . . ARLNTE. S. V. D. EC. TA. VQT. SIH. . . | PL. . . VQS. SK. A. . A. . M |
| <i>Mus. musc-I</i> | NM. DMKFKEEVNHYSNEIN. ACMPLGLHRQF. E. . . | LQM. VQS. K. TV. TM. IS |
| <i>Dro. mela-I</i> | VILLDRKYK. LLDGYTNDINS. TCLPRGLITKF. S. . . | LQL. VLS. K. MV. TM. IS |
| <i>Try. bruc-I</i> | -E--ATAAPM. A. YAT. IQQ. EFVPQRMVLPF. K. H. . . | LLL. T. S. K. N. A. . . S |
| <i>Sac. cere-I</i> | GILDAVTS. KVNAITSQVVS. KCVPDG. MKKE. C. S. . . | MOA. LS. K. NV. VS. IM |
| <i>Sch. pomb-I</i> | QGLDAAMKGMNGLTSSIIN. KCI PDGLLTKE. Y. H. . . | MQT. TVS. K. NV. VS. IS |

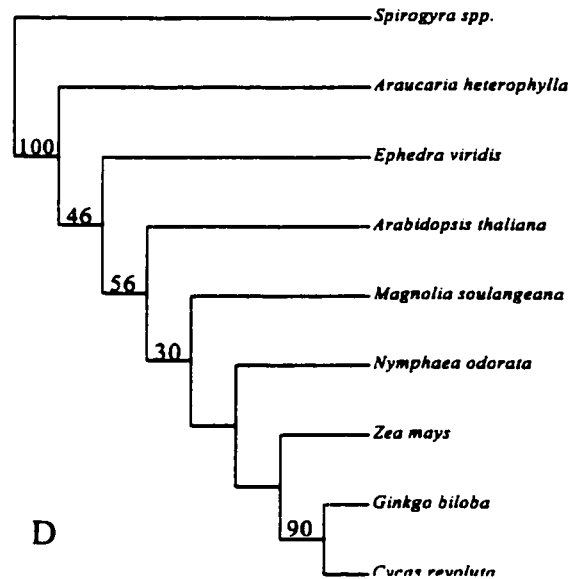
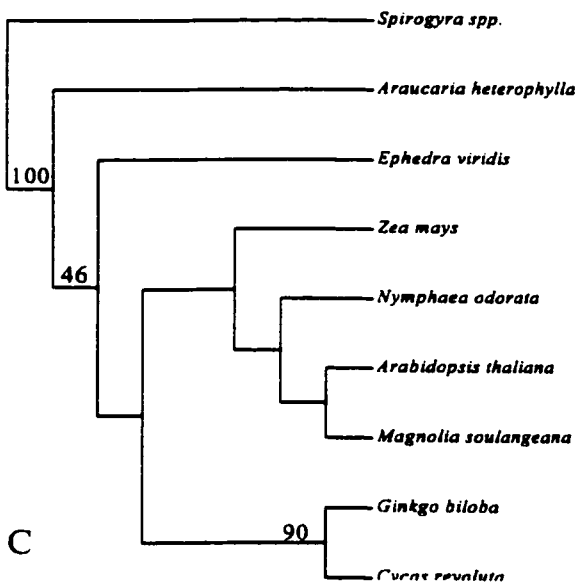
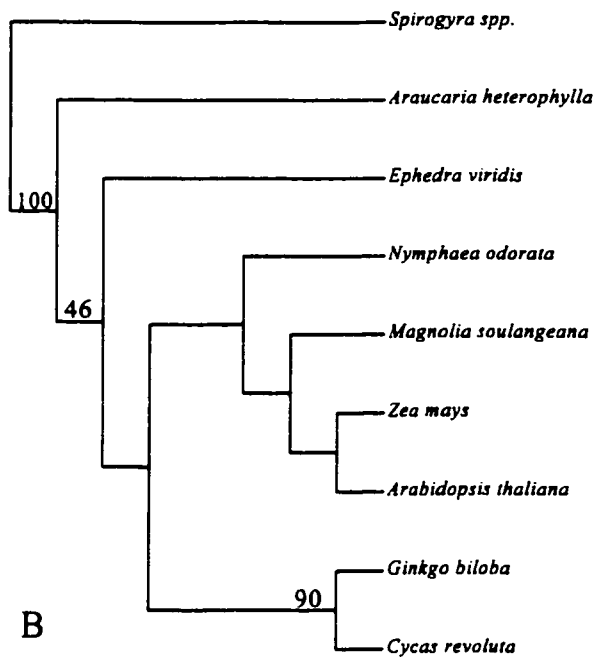
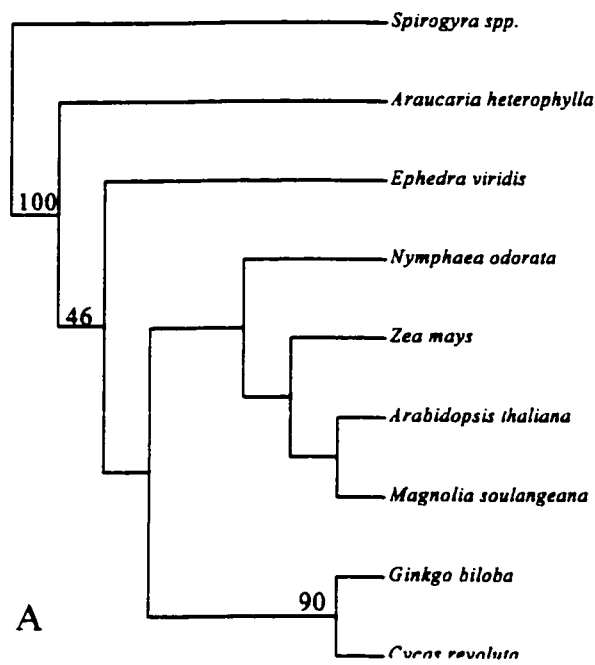


Figure 10: Four equally parsimonious phylogenetic trees inferred from a 280-position alignment of plant RPBI D-F amino acid sequences using PROTPARS in PHYLIP 3.573c. Bootstrap support (% of 1000 replicates) is indicated at the nodes.

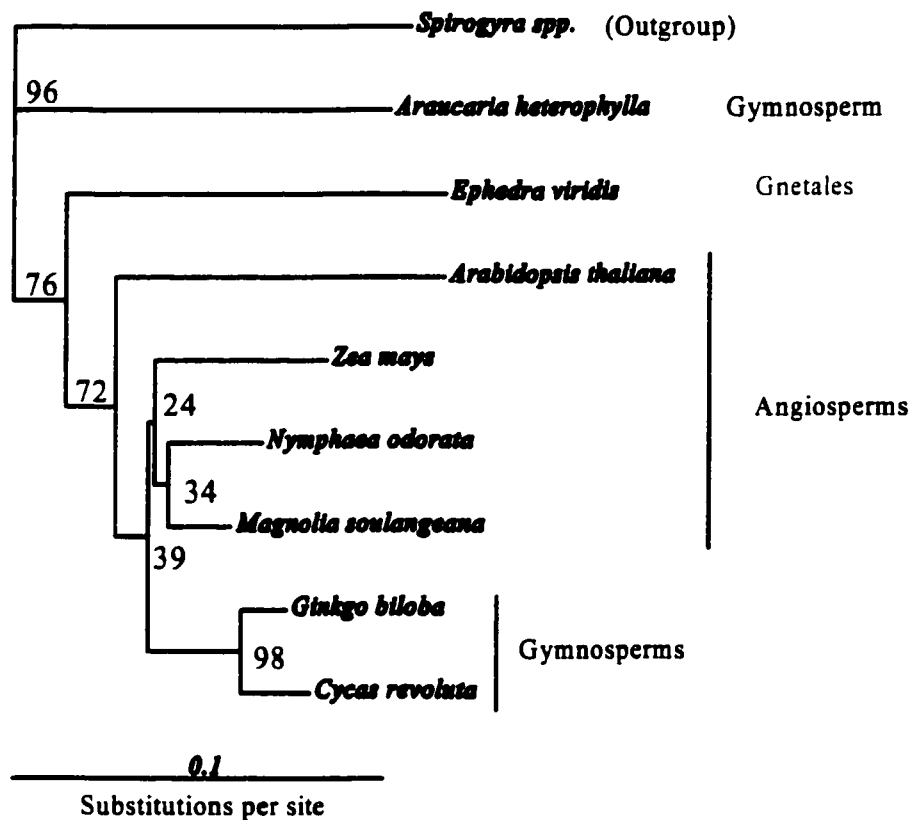


Figure 11: Neighbor-joining distance phylogenetic tree inferred from a 280-position alignment of plant RPB1 D-F amino acid sequences using the Dayhoff PAM model with PROTDIST and NEIGHBOR in PHYLIP 3.573c. Bootstrap support (% of 1000 replicates) is indicated at the nodes.

the *Cycas/Ginkgo* clade, with strong bootstrap support, which is consistent with the neighbor-joining analysis of nucleotide sequences and with only one of the equally parsimonious amino acid trees. In contrast, the two equally parsimonious trees from nucleotide data both placed *Zea mays*, a monocot, as the most basal angiosperm clade. Relationships amongst *Zea mays*, *Nymphaea*, and *Magnolia*, or between these three and their gymnosperm “sister group”, are not resolved with bootstrap values greater than 50%. While this neighbor-joining analysis provided greater statistical support for the relationships which it provided, the relationships of most angiosperms are still unresolved.

Figure 12 (A and B) represent trees with non-clocklike and clocklike branch lengths obtained with PUZZLE using the Blosum62 model of maximum-likelihood analysis. In these trees, dicots are all grouped together, with *Arabidopsis* as a sister group to a clade that consists of the woody and herbaceous dicots. *Cycas* and *Ginkgo* are also placed together with strong statistical support, although the exact relationship of this clade to the dicot clade and the monocot is unresolved. Again, Gnetales are placed at the base of the clade consisting of angiosperms and the *Cycas/Ginkgo* clade with strong statistical support and *Araucaria* is placed as the most basal plant lineage. A most basal angiosperm clade is not established from this analysis. The clocklike branch lengths shown in Figure 12B were rejected by PUZZLE, thus considered to be less likely than the non-clocklike branch lengths. In other words, PUZZLE found it unlikely that all of the sequences had identical substitution rates relative to one another. The analysis represented in Figure 12 (A and B) is characterized by its greater statistical support for relationships amongst the ingroups than was found in the previous phylogenetic analyses (Figures 6-8, 10, 11).

Figure 13 (A and B) represent trees with non-clocklike and clocklike branch lengths obtained with PUZZLE using the JTT model of maximum-likelihood analysis. In this set of trees, angiosperms are clearly separated from the remaining taxa with moderate to high statistical support, and the Gnetales (*Ephedra*) are found near the base of the tree, amongst the gymnosperm taxa, with high statistical support. *Zea mays*, a monocot, is found as the most basal angiosperm (i.e., a sister group to dicots), and the *Cycas/Ginkgo* clade is found as a sister group to angiosperms. Also, *Nymphaea*, a herbaceous dicot, is shown as a sister group to *Magnolia* and *Arabidopsis*. The clocklike branch lengths shown in Figure 13B were rejected by PUZZLE in this analysis as well. The relationships shown from this analysis are markedly different from the maximum-likelihood analysis of nucleotide data in Figure 8, which showed Gnetales as a sister group to *Arabidopsis*, an angiosperm, and did not resolve the relationships amongst angiosperms. However, this analysis did consistently place *Araucaria* as the most basal plant, and place *Cycas* and *Ginkgo* as sister taxa to one another. While some parsimony methods did group angiosperms together, Figure 13 represents the only analysis of these sequences in which the angiosperm taxa are grouped exclusively together and also related to one another with moderate to high statistical support.

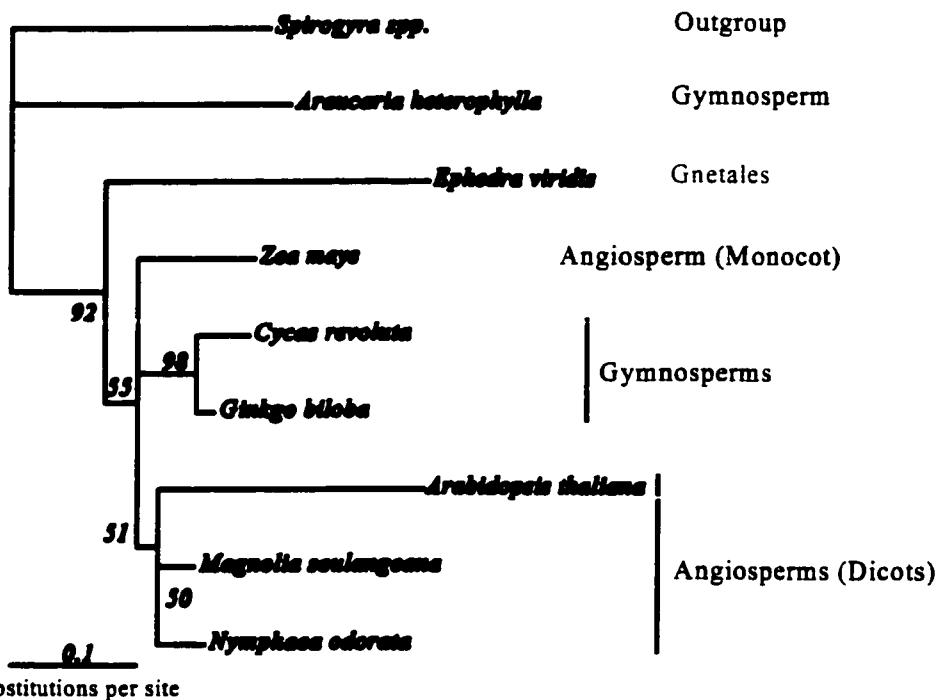


Figure 12A: Maximum-likelihood phylogeny inferred from a 280-position alignment of plant RPB1 D-F amino acid sequences using the Blosum62 model of substitution with an invariable rate and eight γ -distributed rates of substitution and maximum likelihood distances, using PUZZLE 4.0.2. The numbers at the nodes indicate the percent of 10000 quartet puzzling trees supporting that relationship, if >50%. The fraction of invariable sites is 0.41, the γ parameter $\alpha = 1.01$ and $\log L = -1776.73$ for the tree with nonclocklike branch lengths shown here.

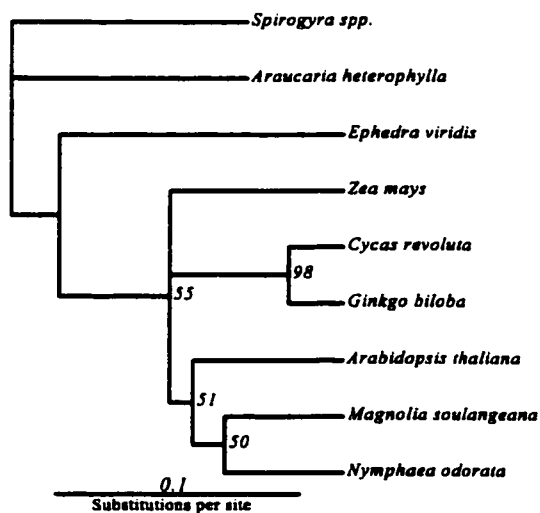


Fig.12B: Same as above, except with clocklike branch lengths. With $\log L = -1805.46$, PUZZLE rejected this tree

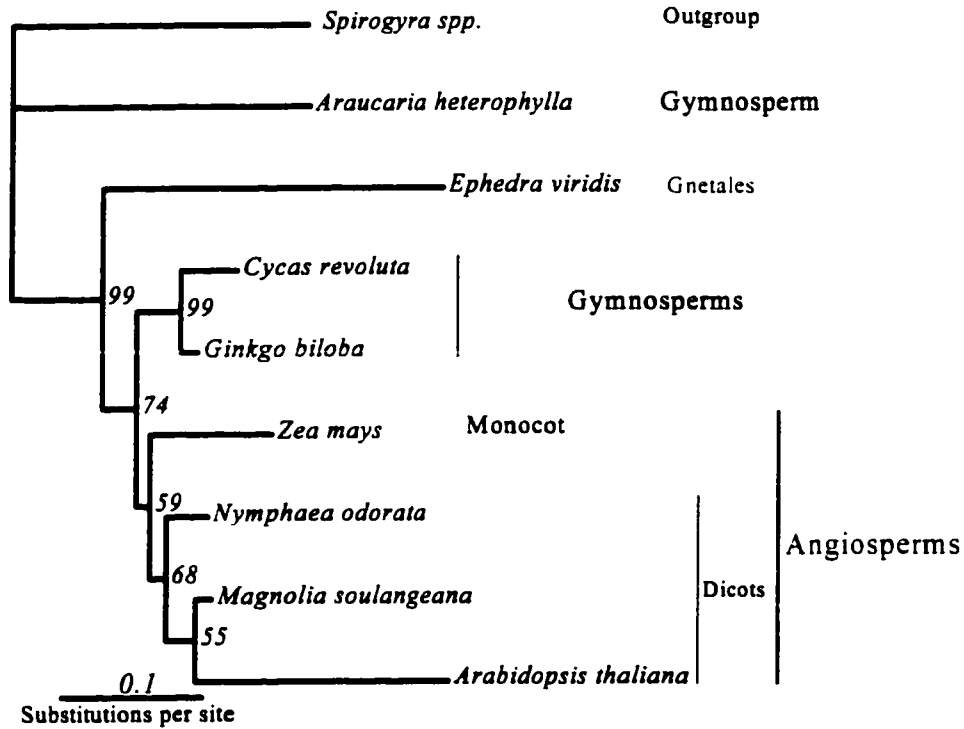


Figure 13A: A maximum likelihood distance phylogenetic tree inferred from a 280-position alignment of plant RPB1 D-F amino acid sequences using an invariable rate and eight γ -distributed rates of substitution according to the JTT model. Numbers at the nodes represent the percent of 10000 quartet puzzling steps supporting that relationship, if >50%. The frequency of invariable sites is 0.41, the γ parameter $\alpha = 1.01$ and $\log L = -1754.94$ for the nonclocklike branchlengths shown here.

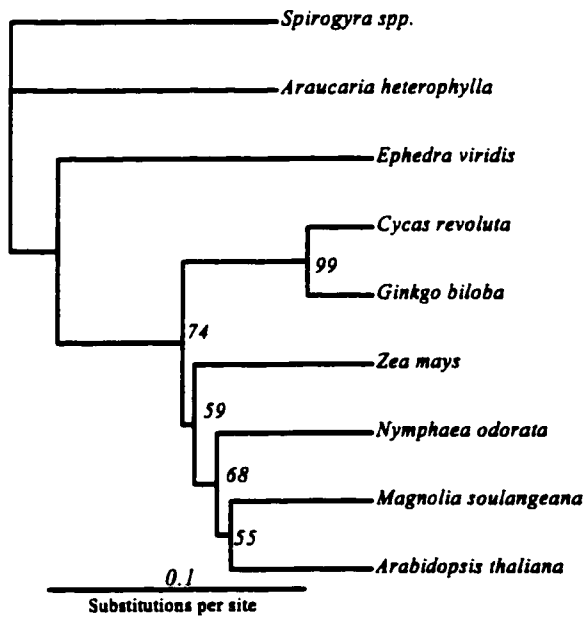


Fig. 13B: A tree with clocklike branch lengths, inferred using the JTT model and otherwise as described above. PUZZLE rejected this tree with $\log L = -1789.17$.

RPB1 in other eukaryotes:

Figures 14, 15 and 16 represent phylogenetic trees of 19 RPB1 D-F amino acid sequences in animals, plants and fungi inferred by parsimony, neighbor-joining and maximum-likelihood methods for contrast with Figure 3, which used regions F-H of RPB1 from only 13 sequences, with another outgroup. The 282-position alignment represented by boxes in dotted lines in Figure 9 was used for this analysis, with the red alga *Porphyra* as the outgroup (see Appendix 3).

In Figure 14, only one of three equally parsimonious trees is shown since they only differed in their groupings of angiosperms amongst one another, none of which was supported by >50% bootstrap support. However, angiosperms were consistently grouped together exclusively, with the *Ginkgo* and *Cycad* clade as a sister group, and Gnetales as a sister group to that clade. However, *Araucaria*, which is a gymnosperm, was put in a more basal position than the green alga *Spirogyra* amongst plants. In the animal clade, both arthropods (*Drosophila* and *Artemia*) were grouped together with <50% bootstrap support, and both molluscs (*Crassostrea* and *Ilyanassa*) were grouped together, with bootstrap support. The nematode (*C. elegans*) was grouped as a sister to the arthropods with strong bootstrap support, and this clade was then grouped as a sister to the molluscs with <50% bootstrap support. The leech (*Helobdella*) was found basal to this group. The mouse represented all vertebrates and was placed as a sister group to all of the invertebrates mentioned above, with 100% bootstrap support. While RPB1 sequences from other vertebrates are available in Genbank, they are also both from placental mammals (human and Chinese hamster) and are identical to mouse at the amino acid level in domains D-F of RPB1, thus they were not included in this analysis. Animals and plants were related as sister groups, in contrast to figure 3.

The neighbor-joining analysis in Figure 15 also grouped animals and plants together. The topology of the plant branches of the tree are consistent with the neighbor-joining analysis in Figure 11, except that *Spirogyra* is placed as an ingroup in the second-most basal position amongst land plants, consistent with the parsimony results shown in Figure 14. The topology of the animal branches of the tree are consistent with Figure 3 only in the relationships amongst arthropods and molluscs. In contrast to the parsimony tree(s), *C. elegans* is placed as a sister group to the arthropod-mollusc clade, albeit with <50% bootstrap support. In Figure 3, *C.elegans* forms a clade with *Helobdella* and together these are a sister group to the arthropod-mollusc clade. However, mouse (*Mus*) is placed between *C.elegans* and *Helobdella* in Figure 15, making *Helobdella* the most basal animal clade, while mouse took that position in Figures 14 and 3. Fungi are grouped together, near the outgroup.

Figure 16 represents the tree with non-clocklike branch lengths that was obtained using the JTT model of maximum-likelihood analysis in PUZZLE. Animals and plants are related as sister groups in this tree. The relationships amongst the major groups of animals and plants, however, are largely unresolved. Amongst animals, arthropods are grouped together and molluscs are grouped together, but no relationships are made with these or the leech, nematode or vertebrate. Amongst plants, *Araucaria* and *Spirogyra* are grouped together and placed as the most basal clade of plants, sister to a mixed clade of angiosperms and gymnosperms. In this group, *Nymphaea* and *Magnolia* (two angiosperms) are grouped together, and *Cycas* and *Ginkgo* (two gymnosperms) are grouped together, but the relationship of these two groups with one another or with the remaining angiosperms and Gnetales are not resolved, in contrast with figures 12 and 13.

Figures 17 and 18 represent phylogenetic trees of 38 sequences of regions D-F of the largest subunit of RNA polymerases I, II and III (RPA1, RPB1 and RPC1) from a wider assortment of eukaryotes, inferred using the JTT and Blosum62 maximum likelihood methods in PUZZLE. The 252-position alignment represented by boxes in solid lines in Figure 9 was used for this analysis, with the Crenarchaeote archaeon *Sulfolobus acidocaldarius* RPA' as the outgroup (see Appendix 4). While this is a very distant outgroup, it was chosen in order to more clearly determine whether RPA1 or RPC1 are the sister group to RPB1. Both trees place RPA1 as the sister group to RPB1. In Figure 17, the relationships of the major eukaryotic groups to one another are not resolved by either RPB1 or by RPC1 sequences by the JTT model of substitution. The relationships seen in RPB1 sequences within animal taxa and within plants are similar to those in Figure 16, except that *Spirogyra* and *Araucaria* are not separated from other plants here. Other RPB1 clades seen in figure 17 are: red algae, fungi, Microsporidia, slime molds (*Dictyostelium* and *Acanthamoeba*) and a group with an Apicomplexan and a Parabasalid (*Plasmodium* and *Trichomonas*). In the RPA1 branch of the tree, animals and fungi are related as sister groups, in the absence of the newly published *Arabidopsis* RPA1 sequence.

The relationships seen in Figure 18 are similar to those in Figure 17, with a few notable exceptions. The relationship of *Ginkgo* and *Cycas* is lost in Figure 18, and the slime molds form a sister group to all plants. A Kinetoplastid - Archamoeba clade is formed with *Mastigamoeba* and *Trypanosoma* RPB1. Within the RPC1 sequences, a fungus and a Diplomonad (*Giardia*) form a clade, with a Kinetoplastid (*Trypanosoma*) as a sister group. Other relationships are essentially the same as those seen in Figure 17.

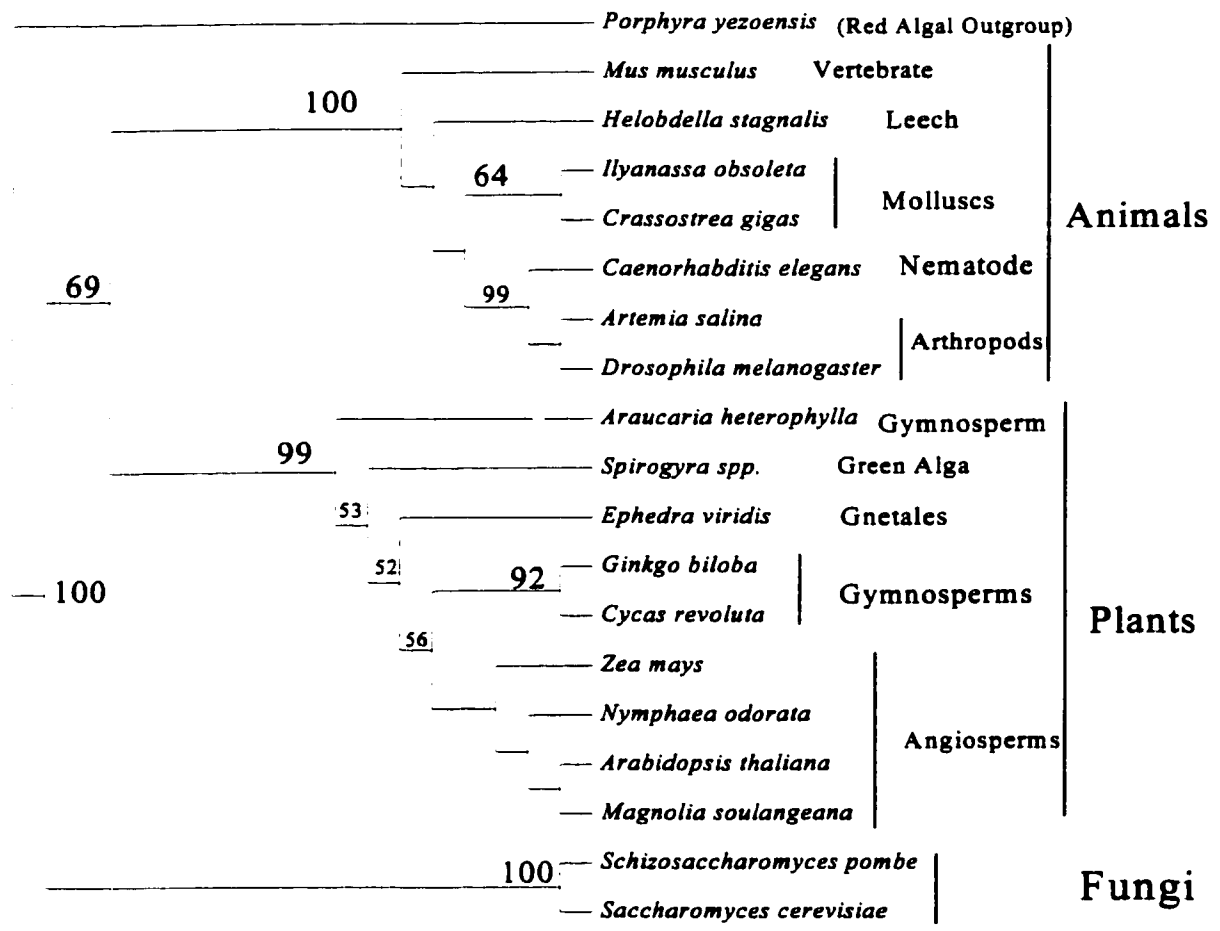


Figure 14: One of three equally most parsimonious phylogenetic trees inferred from a 282-position alignment of amino acid sequences of domains D-F of RPB1 from Animals, Plants and Fungi, using PROTPARS of PHYLIP 3.573c. The topology matching that of the consensus bootstrapped tree is shown, with the percent of 1000 bootstrap replicates supporting that topology shown at the nodes, if above 50%. *Porphyra*, a red alga, is the outgroup. The other two topologies differed only in that, amongst angiosperms, *Arabidopsis* remained in the same position, paired with either *Nymphaea* or *Zea*.

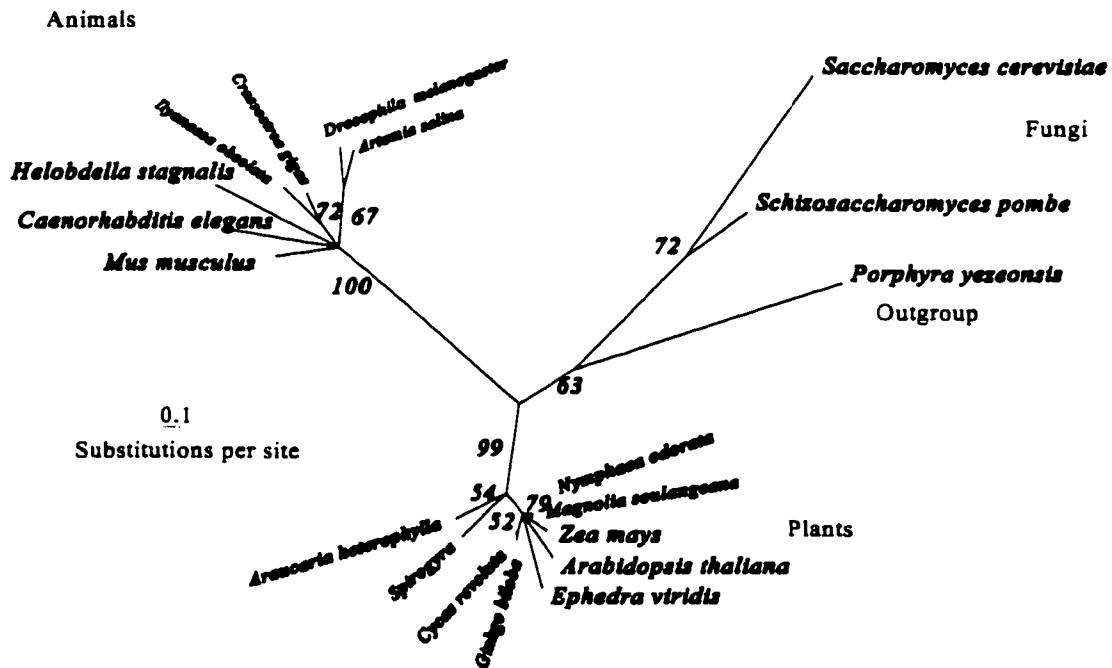


Figure 16: Maximum-likelihood distance phylogenetic tree with non-clocklike branch lengths of a 282-position alignment of RPB1 D-F amino acid sequences from Animals, Plants and Fungi, estimated using the JTT model of substitution and an invariable and eight γ -distributed rates of substitution with PUZZLE v.4.0.2. Numbers at the nodes represent the percent of 10000 quartet puzzling steps supporting that relationship, if >50%. The outgroup is a red alga. The γ parameter $\alpha = 0.87$ and the fraction of invariable sites is 0.06. The log L of the tree shown here with non-clocklike branch lengths is -4849.98. The tree with clocklike branch lengths had a log L of -4897.91, and was rejected by PUZZLE and considered to be unlikely.

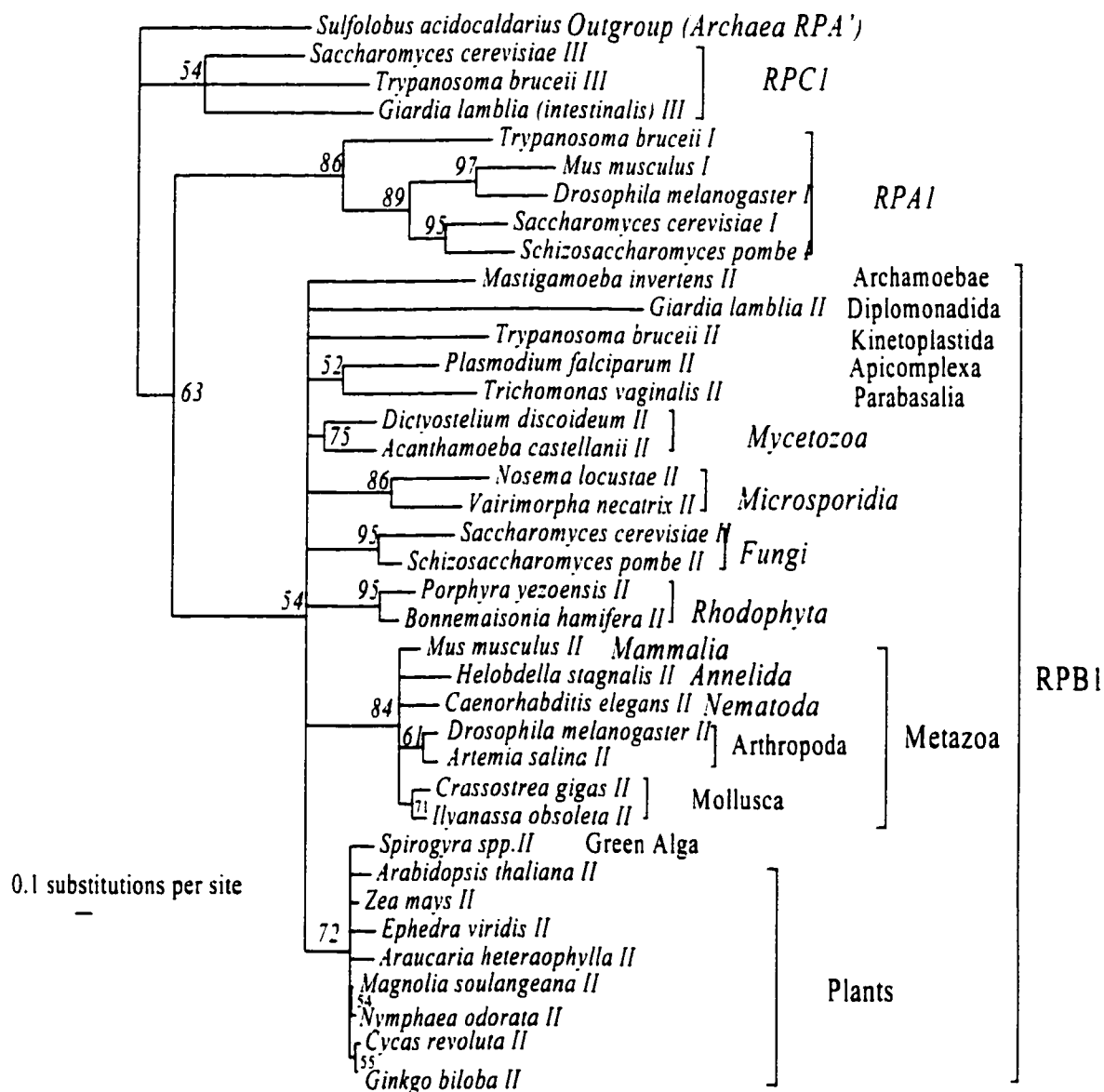


Figure 17: A maximum-likelihood distance phylogenetic tree inferred from a 252-position alignment of the amino acid sequences for domains D-F of the largest subunit of all eukaryotic RNA polymerases, inferred using the JTT model of substitution and considering an invariable and 8 γ -distributed rates of substitution, with the PUZZLE 4.0.2 program. The numbers at the nodes represent the percent of 10000 quartet puzzling steps supporting that topology, if above 50%. The fraction of invariable sites is 0.07 (S.E. 0.02) and γ parameter $\alpha = 1.72$ (S.E. 0.18). The log L of th tree with non-clocklike branchlengths, shown here, is -12657.44, while the log L of the tree with clocklike branchlengths was -12788.63 and rejected by PUZZLE.

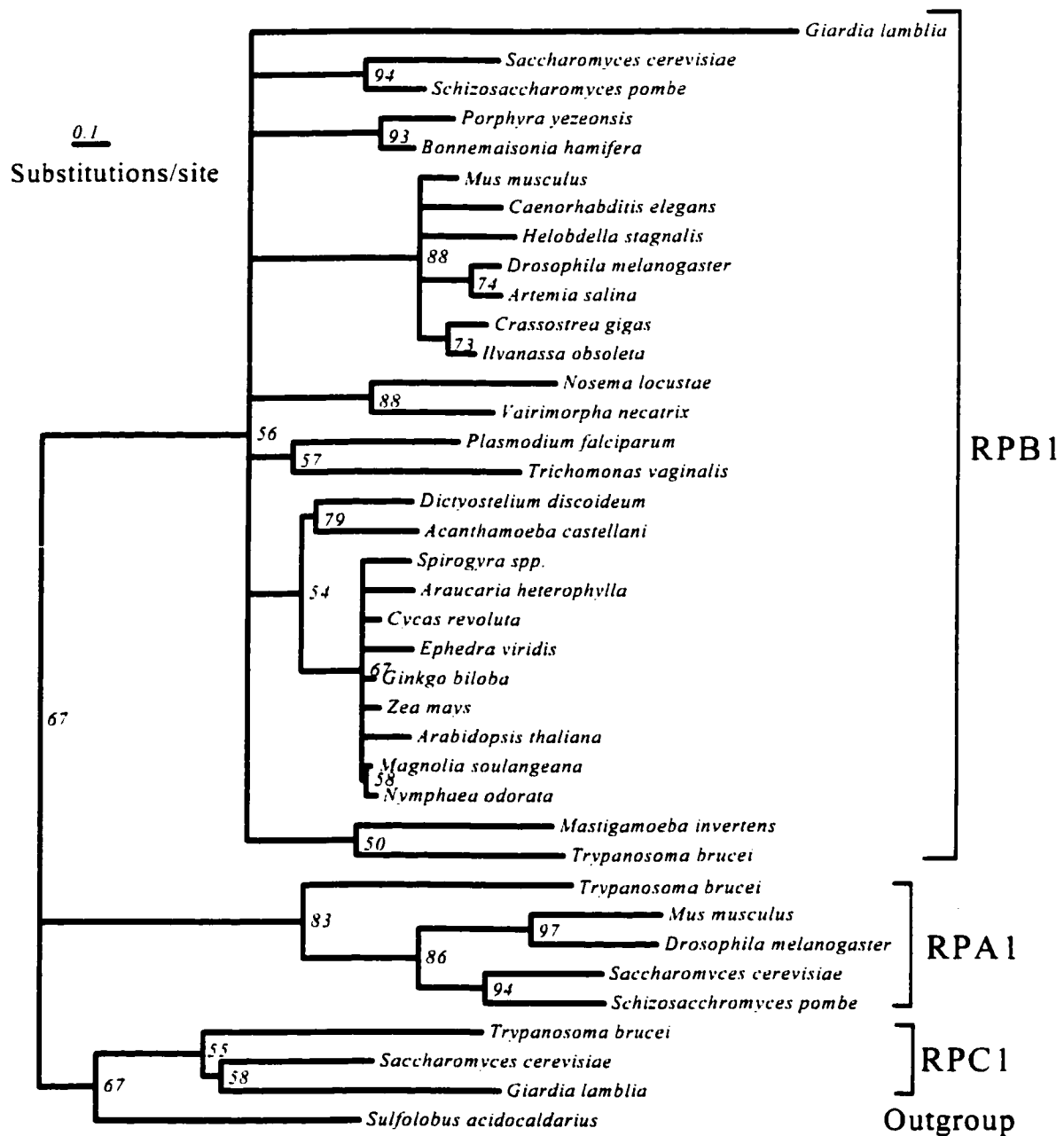


Figure 18: Maximum-likelihood distance phylogenetic tree inferred from a 252 residue alignment of amino acid sequences of domains D-F of the largest subunit of RNA polymerases I, II and III from all eukaryotes, with the archaeal homologue (RPA') from *Sulfolobus acidocaldarius* as the outgroup. The Blosum62 model of substitution was used in Puzzle 4.0.2, also considering an invariable rate as well as 8 gamma-distributed rates of substitution. Gaps and positions which could not be unambiguously aligned in more than three taxa were not included in the analysis. Numbers at the nodes represent the percent of 10000 quartet puzzling steps supporting that relationship, if over 50%. The log L for this tree with nonclocklike branch lengths is -12580.61 and the clocklike tree was rejected. The frequency of invariable sites is 0.07 and the γ distribution parameter $\alpha = 2.23$.

Compositional Bias:

The G+C content of the region of the RPB1 gene encoding domains D-F varies from 41% to 54% in plants and from 36% to 60% in a survey including most other eukaryotes (Figure 19). In plants and in other eukaryotes, this compositional bias is seen the most at the third base of codons, and at the first base of codons to a lesser degree. The nucleotide composition at the second base of codons is not affected by compositional bias. The compositional bias seen at the first base of codons is mostly due to bias in the first base of the sixfold degenerate arginine and leucine codons.

The analysis of G+C content at the third codon position of RPB1 regions D-F shown in Figure 20 illustrates the codon usage patterns seen amongst plants and other eukaryotes. For example, Figure 20B shows that most animals, most plants, and both fungi each share similar codon usage patterns. Figure 21 also shows the diversity of G+C vs A+T content at the third base of codons amongst plants for this region of RPB1. Within plants, the four angiosperms also share more similar nucleotide composition together than they do with the other plants. Fungi are 30-35% G+C at the third codon position, while most plants are 35-44% G+C and most animals fell into a broader range of 44-70% G+C, at the same codon position.

Figure 22 illustrates the effects of G+C compositional bias on amino acid usage. The G+C content in the coding region is compared to the content of amino acids with G+C rich codons (G, A, R, P) and to the content of amino acids with A+T rich codons (F, Y, M, I, N, K). Whether closely examining plants alone, or looking at a more broad range of eukaryotes, there appears to be little effect on amino acid usage from nucleotide compositional bias.

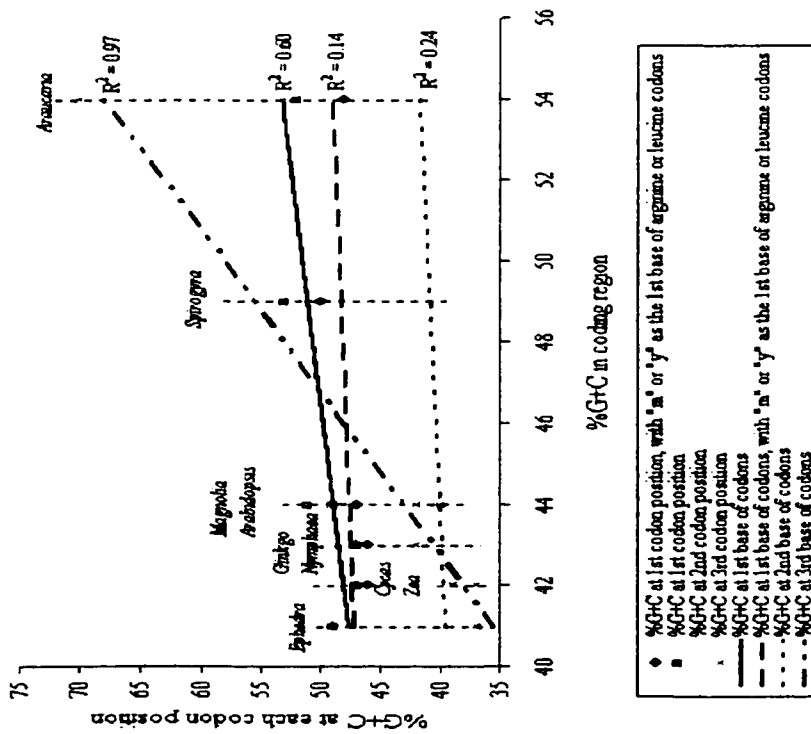


Figure 19. % G+C in the coding region vs % G+C at each codon position of RPB1 D-F in plants, demonstrating the effects of compositional bias at the first and third bases of codons

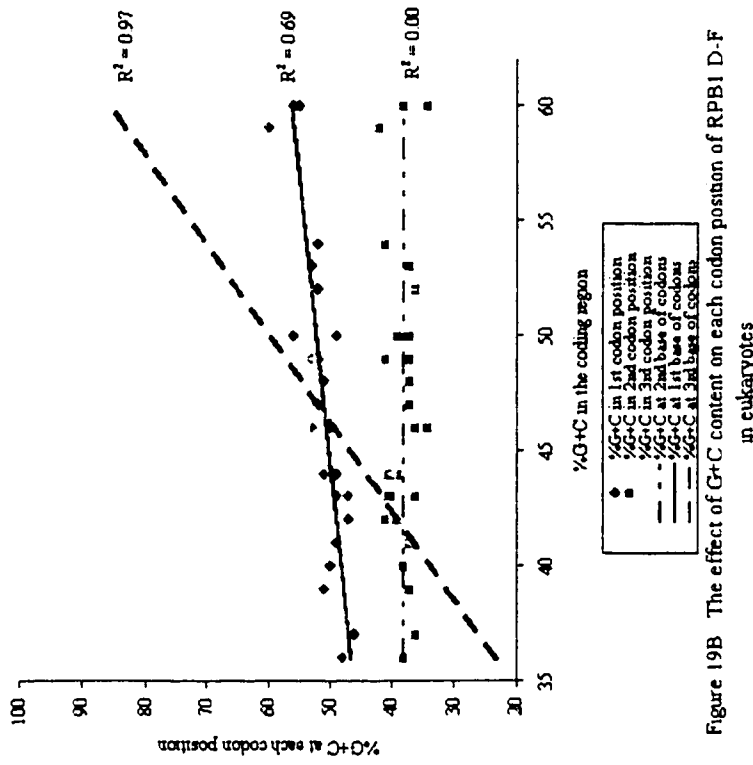


Figure 19B The effect of G+C content on each codon position of RPB1 D-F in eukaryotes

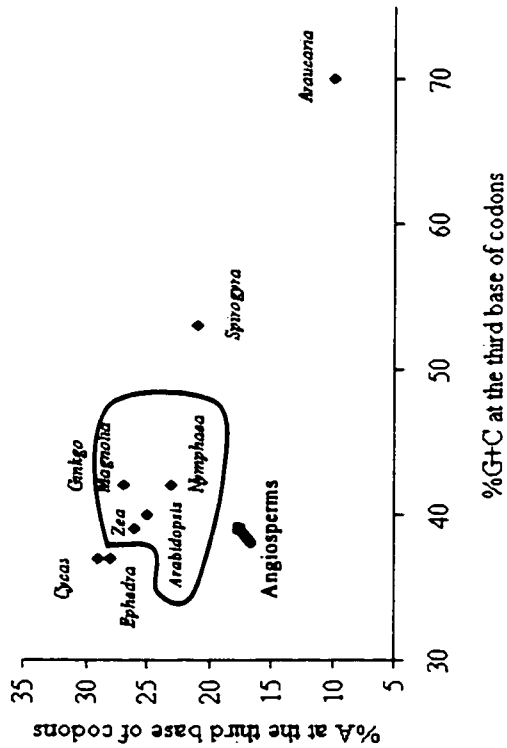
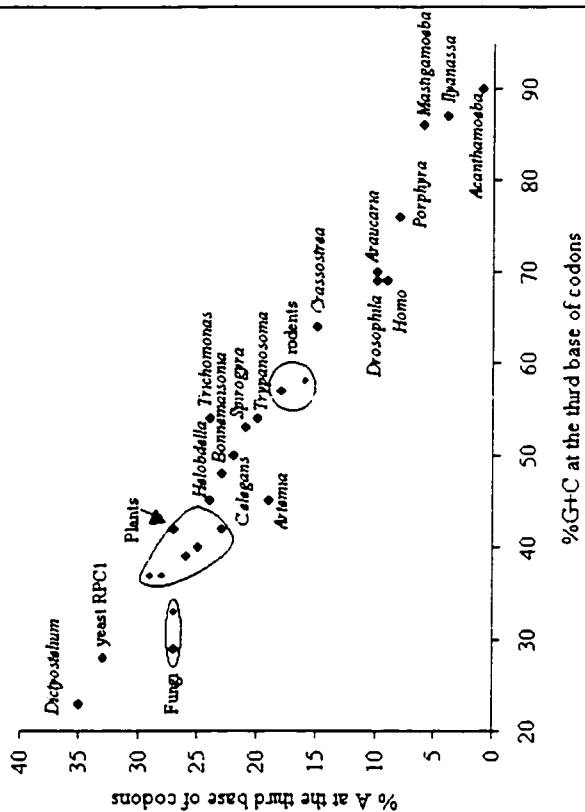


Figure 20: %G+C vs %A at the third base of codons of RPB1 D-F in plants

Figure 20B: G+C content at the third base of codons reflects codon usage patterns in RPB1 D-F amongst eukaryotes



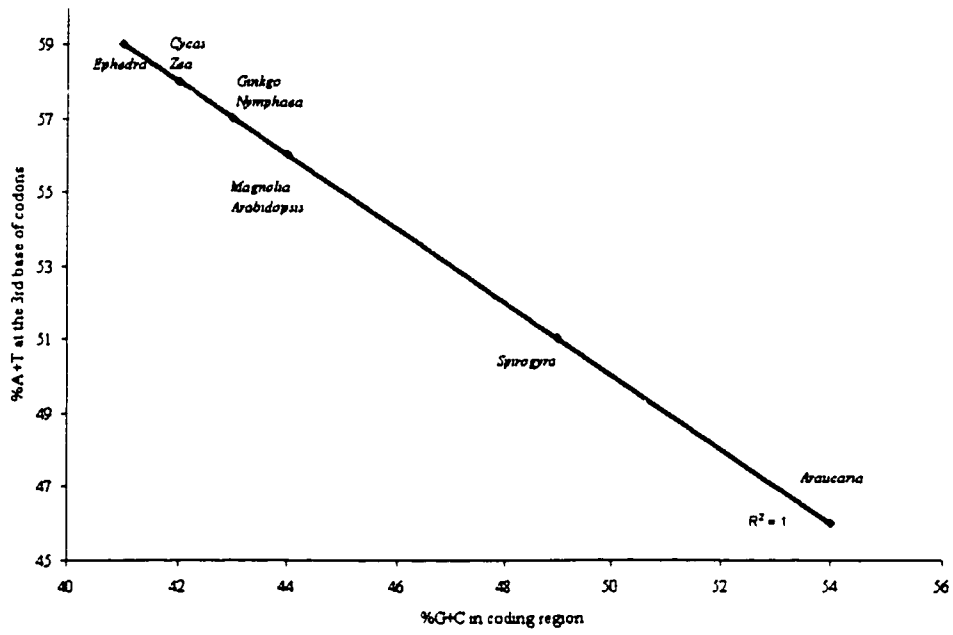
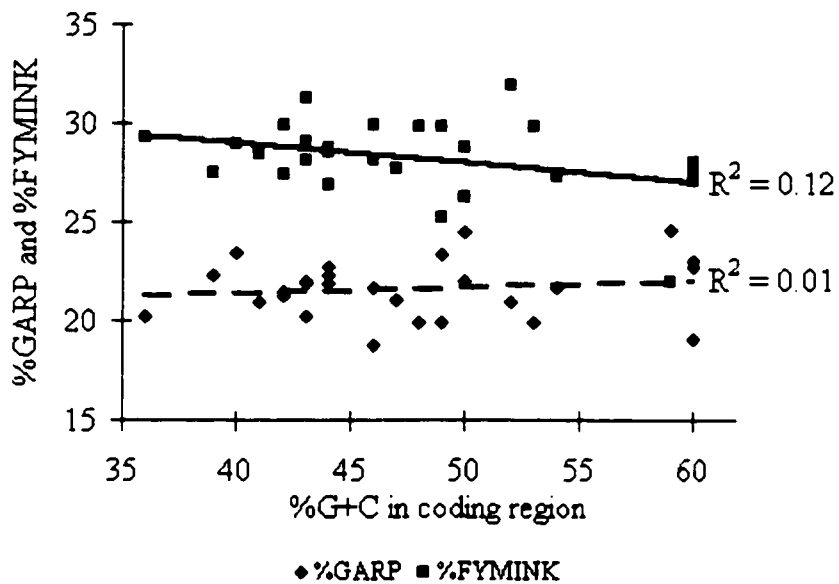
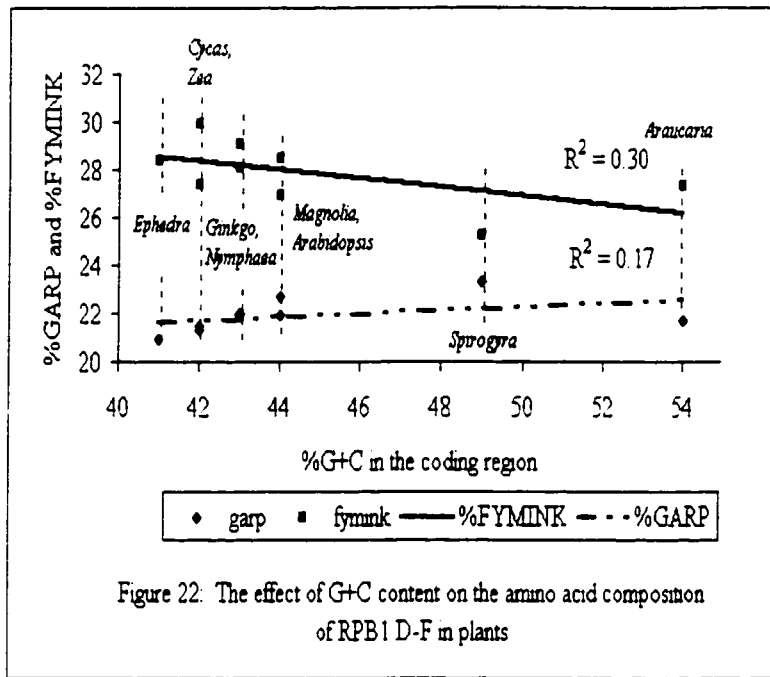


Figure 21: %G+C of coding region vs % A+T at the third codon position in RPB1 D-F of plants, demonstrating the variation in G+C content among plant RPB1 D-F sequences.



Substitution Analysis of Region D-F of RPBI in Land Plants:

Tables 5 and 6 show the average nucleotide substitutions per site amongst the 834 bp alignment of plant sequences shown in figure 5, as calculated using the Li93 program. Table 5 shows the average nonsynonymous substitutions per site (Ka), when they were calculated to be less than 3.24 by the Li93 method. According to the Jukes-Cantor formula $K = -\frac{3}{4} \ln(1 - \frac{4}{3}P)$, if the number of substitutions observed (P) is ≥ 0.75 , then K is undefined. In this situation, Li93 cannot calculate K. Those numbers greater than 3.24 result from the presence of too many multiple substitutions. The Jukes-Cantor formula was used to estimate how many of the 556 nonsynonymous sites in the alignment had multiple substitutions. This is summarized in Table 7. Table 6 shows the average synonymous substitutions per site when they were calculated to be less than 3.24 since other values represented the presence of too many multiple substitutions. The estimate of the number of 278 synonymous sites in the alignment that had multiple substitutions is summarized in Table 8. These results show that the synonymous sites are saturated with multiple substitutions, while the incidence of multiple substitutions at nonsynonymous sites is approximately 100 times less.

The Ka of *Nymphaea* and *Magnolia* vs maize and vs *Araucaria*, *Ginkgo* and *Cycas*, as well as the Ka of *Spirogyra* vs *Araucaria* and *Nymphaea* compared with the approximate divergence times of these groups is shown in Figure 23. Other comparisons are not shown in this graph if either the divergence time is uncertain or the Ka (see Table 5) is undefined. Major conclusions should not be made from this graph without the inclusion of more data from vascular plants at least.

Table 5: The weighted average of nonsynonymous substitutions per site (K_a), with standard error, calculated for the nine plant RPB1 D-F nucleotide sequences by Li93. The $P_a = \frac{3}{4} (1 - e^{-4/3 K_a})$ is shown underneath in parentheses. Dashes represent undefined terms.

| | | Standard Error | | | | | | | | |
|--------------------------------------|-------------------------------|--------------------------------|---|---------------------------------|----------------------------------|--------------------------------|---------------------------------------|-----------------------------------|-----------------|---------------------------------------|
| | | <i>Spirogyra</i> <i>sp.</i> | <i>Araucaria</i> <i>heterophylla</i> | <i>Cycas</i> <i>revoluta</i> | <i>Ephedra</i> <i>viridis</i> | <i>Ginkgo</i> <i>biloba</i> | <i>Magnolia</i> <i>soulangeana</i> | <i>Nymphaea</i> <i>odorata</i> | <i>Zea mays</i> | <i>Arabidopsis</i> <i>thaliana</i> |
| Nonsynonymous substitutions per site | <i>Spirogyra sp.</i> | | 0.016 | - | 0.016 | - | - | 0.014 | - | - |
| | <i>Araucaria heterophylla</i> | 0.136 (0.124) | | - | - | - | 0.013 | - | - | - |
| | <i>Cycas revoluta</i> | - | - | | 0.013 | 0.005 | 0.008 | 0.008 | 0.01 | 0.011 |
| | <i>Ephedra viridis</i> | 0.137 (0.125) | - | 0.097 (0.091) | | 0.013 | 0.013 | 0.013 | 0.013 | 0.014 |
| | <i>Ginkgo biloba</i> | - | - | 0.018 (0.018) | 0.093 (0.087) | | 0.007 | 0.007 | 0.01 | 0.012 |
| | <i>Magnolia soulangeana</i> | - | 0.093 (0.087) | 0.041 (0.040) | 0.094 (0.088) | 0.035 (0.034) | | 0.006 | 0.01 | 0.009 |
| | <i>Nymphaea odorata</i> | 0.114 (0.106) | - | 0.039 (0.038) | 0.095 (0.089) | 0.029 (0.028) | 0.022 (0.022) | | 0.01 | 0.01 |
| | <i>Zea mays</i> | - | - | 0.049 (0.047) | 0.101 (0.094) | 0.047 (0.046) | 0.034 (0.033) | 0.043 (0.042) | | - |
| | <i>Arabidopsis thaliana</i> | - | - | 0.077 (0.073) | 0.105 (0.098) | 0.080 (0.076) | 0.053 (0.051) | 0.065 (0.062) | - | |

Table 6: The weighted average of synonymous substitutions per site (K_s), with standard error, calculated for the nine plant RPB1 D-F nucleotide sequences by Li93. P_s is shown underneath in parentheses. Dashes represent undefined terms.

| | | Standard Error | | | | | | | | |
|---|-------------------------------|--------------------------------|---|---------------------------------|----------------------------------|--------------------------------|---------------------------------------|-----------------------------------|-----------------|---------------------------------------|
| | | <i>Spirogyra</i> <i>sp.</i> | <i>Araucaria</i> <i>heterophylla</i> | <i>Cycas</i> <i>revoluta</i> | <i>Ephedra</i> <i>viridis</i> | <i>Ginkgo</i> <i>biloba</i> | <i>Magnolia</i> <i>soulangeana</i> | <i>Nymphaea</i> <i>odorata</i> | <i>Zea mays</i> | <i>Arabidopsis</i> <i>thaliana</i> |
| Synonymous substitutions per site (K_s) | <i>Spirogyra sp.</i> | | - | - | - | - | - | - | - | - |
| | <i>Araucaria heterophylla</i> | - | | - | - | - | - | - | - | - |
| | <i>Cycas revoluta</i> | - | - | | 0.258 | 0.077 | 0.226 | 0.276 | 3.882 | 2.52 |
| | <i>Ephedra viridis</i> | - | - | 1.323 (0.621) | | 0.261 | 0.824 | - | - | - |
| | <i>Ginkgo biloba</i> | - | - | 0.562 (0.395) | 1.451 (0.642) | | 0.267 | 0.68 | - | - |
| | <i>Magnolia soulangeana</i> | - | - | 1.31 (0.62) | 1.96 (0.70) | 1.566 (0.657) | | 0.209 | 0.242 | 0.507 |
| | <i>Nymphaea odorata</i> | - | - | 1.456 (0.642) | - | 2.28 (0.71) | 1.318 (0.621) | | 0.346 | 0.769 |
| | <i>Zea mays</i> | - | - | 2.835 (0.733) | - | - | 1.476 (0.645) | 1.669 (0.669) | | - |
| | <i>Arabidopsis thaliana</i> | - | - | 2.988 (0.736) | - | - | 1.889 (0.690) | 2.162 (0.708) | - | |

Table 7: Estimate of the number of multiple substitutions at nonsynonymous (n.s.) sites in plant RPBI D-F sequences shown in Figure 5, calculated using the difference between the K_a and P_a values shown in Table 5.

| species 1 | species 2 | Inferred frequency of n.s. substitutions (K_a) | Observed frequency of n.s. substitutions (P_a) | % multiple substitutions at n.s. sites | # of multiple substitutions in 556 n.s. sites |
|------------------|--------------------|--|--|--|---|
| <i>Araucaria</i> | <i>Spirogyra</i> | 0.136 | 0.124 | 1.2 % | 7 |
| <i>Ephedra</i> | <i>Spirogyra</i> | 0.137 | 0.125 | 1.2 % | 7 |
| <i>Nymphaea</i> | <i>Spirogyra</i> | 0.114 | 0.106 | 0.8 % | 4 |
| <i>Araucaria</i> | <i>Magnolia</i> | 0.093 | 0.087 | 0.6 % | 3 |
| <i>Cycas</i> | <i>Ephedra</i> | 0.097 | 0.091 | 0.6 % | 3 |
| <i>Cycas</i> | <i>Ginkgo</i> | 0.018 | 0.018 | 0 % | 0 |
| <i>Cycas</i> | <i>Magnolia</i> | 0.041 | 0.04 | 0.1 % | 1 |
| <i>Cycas</i> | <i>Nymphaea</i> | 0.039 | 0.038 | 0.1 % | 1 |
| <i>Cycas</i> | <i>Zea</i> | 0.049 | 0.047 | 0.2 % | 1 |
| <i>Cycas</i> | <i>Arabidopsis</i> | 0.077 | 0.073 | 0.4 % | 2 |
| <i>Ephedra</i> | <i>Ginkgo</i> | 0.093 | 0.087 | 0.6 % | 3 |
| <i>Ephedra</i> | <i>Magnolia</i> | 0.094 | 0.088 | 0.6 % | 3 |
| <i>Ephedra</i> | <i>Nymphaea</i> | 0.095 | 0.089 | 0.6 % | 3 |
| <i>Ephedra</i> | <i>Zea</i> | 0.101 | 0.094 | 0.7 % | 4 |
| <i>Ephedra</i> | <i>Arabidopsis</i> | 0.105 | 0.098 | 0.7 % | 4 |
| <i>Ginkgo</i> | <i>Magnolia</i> | 0.035 | 0.034 | 0.1 % | 1 |
| <i>Ginkgo</i> | <i>Nymphaea</i> | 0.029 | 0.028 | 0.1 % | 1 |
| <i>Ginkgo</i> | <i>Zea</i> | 0.047 | 0.046 | 0.1 % | 1 |
| <i>Ginkgo</i> | <i>Arabidopsis</i> | 0.08 | 0.076 | 0.4 % | 2 |
| <i>Magnolia</i> | <i>Nymphaea</i> | 0.022 | 0.022 | 0 % | 0 |
| <i>Magnolia</i> | <i>Zea</i> | 0.034 | 0.033 | 0.1 % | 1 |
| <i>Magnolia</i> | <i>Arabidopsis</i> | 0.053 | 0.051 | 0.2 % | 1 |
| <i>Nymphaea</i> | <i>Zea</i> | 0.043 | 0.042 | 0.1 % | 1 |
| <i>Nymphaea</i> | <i>Arabidopsis</i> | 0.065 | 0.062 | 0.3 % | 2 |

Table 8: Estimate of the number of multiple substitutions at synonymous (syn.) sites in plant RPBI D-F sequences shown in Figure 5, calculated using the difference between the K_s and P_s values shown in Table 6.

| species 1 | species 2 | Inferred frequency of syn. substitutions (K_s) | Observed frequency of syn. substitutions (P_s) | % multiple substitutions at syn. sites | # of multiple substitutions in 278 syn. sites |
|-----------------|--------------------|--|--|--|---|
| <i>Cycas</i> | <i>Ephedra</i> | 1.323 | 0.621 | 70.2 % | 195 |
| <i>Cycas</i> | <i>Ginkgo</i> | 0.562 | 0.395 | 16.7 % | 46 |
| <i>Cycas</i> | <i>Magnolia</i> | 1.31 | 0.62 | 69 % | 192 |
| <i>Cycas</i> | <i>Nymphaea</i> | 1.456 | 0.642 | 81.4 % | 226 |
| <i>Cycas</i> | <i>Zea</i> | 2.835 | 0.733 | 210.2 % | 584 |
| <i>Cycas</i> | <i>Arabidopsis</i> | 2.988 | 0.736 | 225.2 % | 626 |
| <i>Ephedra</i> | <i>Ginkgo</i> | 1.451 | 0.642 | 80.9 % | 225 |
| <i>Ephedra</i> | <i>Magnolia</i> | 1.96 | 0.7 | 126 % | 350 |
| <i>Ginkgo</i> | <i>Magnolia</i> | 1.566 | 0.657 | 90.9 % | 253 |
| <i>Ginkgo</i> | <i>Nymphaea</i> | 2.28 | 0.71 | 157 % | 436 |
| <i>Magnolia</i> | <i>Nymphaea</i> | 1.318 | 0.621 | 69.7 % | 194 |
| <i>Magnolia</i> | <i>Zea</i> | 1.476 | 0.645 | 83.1 % | 231 |
| <i>Magnolia</i> | <i>Arabidopsis</i> | 1.889 | 0.69 | 119.9 % | 333 |
| <i>Nymphaea</i> | <i>Zea</i> | 1.669 | 0.669 | 100 % | 278 |
| <i>Nymphaea</i> | <i>Arabidopsis</i> | 2.162 | 0.708 | 145.4 % | 404 |

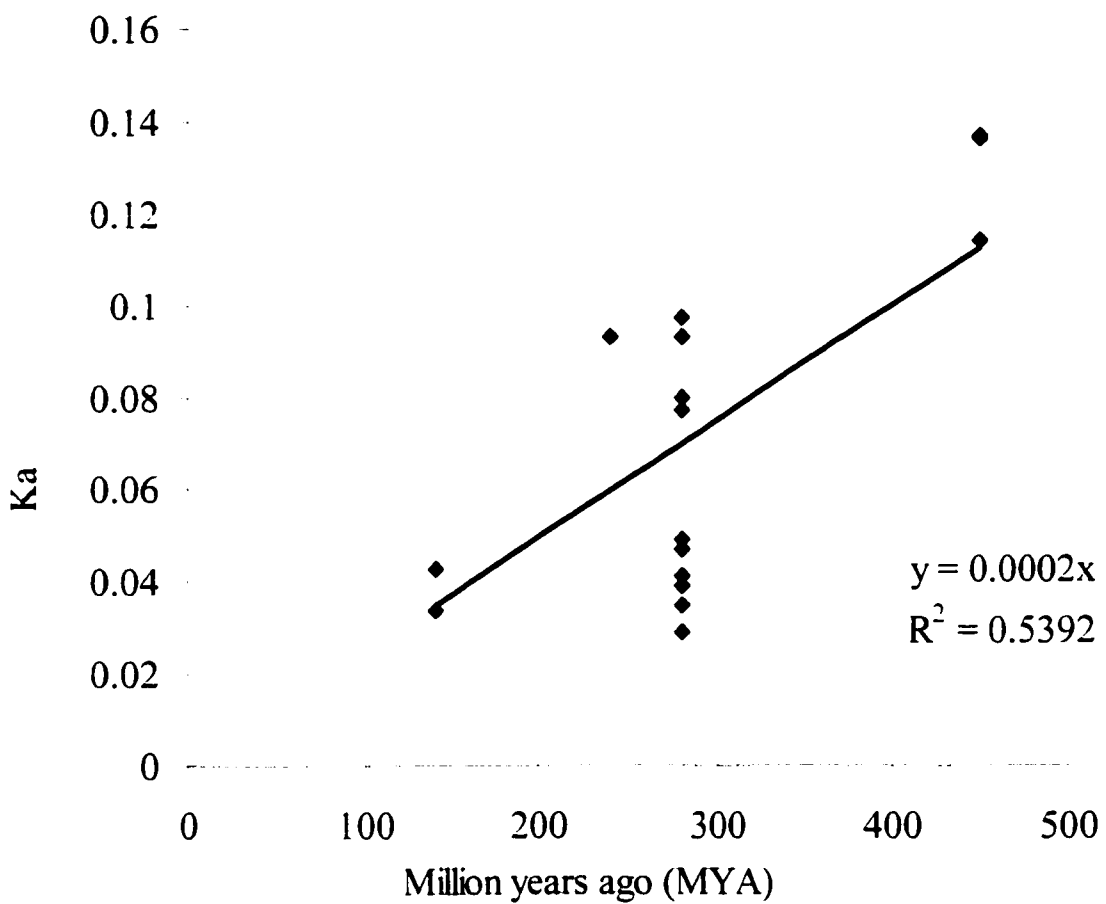


Figure 23: Divergence times vs nonsynonymous substitutions per site (Ka) in plant RPB1D-F sequences. Ka values (see Table 7) of *Magnolia* and *Nymphaea* vs maize, *Cycas* and *Ginkgo*, of *Araucaria* vs *Spirogyra* and *Magnolia*, of *Spirogyra* vs *Nymphaea* and *Ephedra*, and of maize, *Ephedra* and *Arabidopsis* vs *Cycas* and *Ginkgo* were compared to divergence times adapted from Crane *et al* (1995), Cronquist (1981) and Stewart and Rothwell (1993). Some comparisons (e.g., *Magnolia* vs *Nymphaea*) were not plotted because of unknown divergence times, while others were not plotted because the Ka was undefined.

Relative Rate Tests:

Relative rate tests were done using the RRTree program with amino acid sequence alignments of plants, of animals, plants and fungi, and of RPA1, RPB1 and RPC1. These results are summarized in Tables 9-12. Each plant sequence was compared individually (Table 9), while in the other analyses, groups of sequences were compared to each other using the group's mean substitution rate (K), with each taxon in the group equally weighted (Tables 10-12). The p-value for the exact probability that the two groups have the same substitution rates was tested against $\alpha = 0.05/(\text{number of comparisons})$. No significantly different rates of substitution were found.

Table 9: Results of relative rate tests of substitution rates (K) of amino acid sequences of RPB1 regions D-F amongst plants, based upon the alignments shown in Figure 9 and also in Appendix 2. Since 28 pairwise comparisons were made in this analysis, $\alpha = 0.05/28 = 0.002$. *Spirogyra* is the outgroup. P-values are given in the upper upper corner of the table, while the difference in substitution rate is given in the bottom of the table.

| | p-value | | | | | | | |
|--|-------------------------------|-----------------------|------------------------|----------------------|-----------------------------|-------------------------|-----------------|-----------------------------|
| | <i>Araucaria heterophylla</i> | <i>Cycas revoluta</i> | <i>Ephedra viridis</i> | <i>Ginkgo biloba</i> | <i>Magnolia soulangeana</i> | <i>Nymphaea odorata</i> | <i>Zea mays</i> | <i>Arabidopsis thaliana</i> |
| difference in substitution rate (ΔK) | <i>Araucaria heterophylla</i> | 0.765 | 0.895 | 0.678 | 0.715 | 0.833 | 0.891 | 0.894 |
| <i>Cycas revoluta</i> | 0.04 | | 0.611 | 0.809 | 0.958 | 0.874 | 0.804 | 0.611 |
| <i>Ephedra viridis</i> | -0.018 | -0.058 | | 0.542 | 0.607 | 0.704 | 0.763 | 0 |
| <i>Ginkgo biloba</i> | 0.053 | 0.013 | 0.071 | | 0.909 | 0.717 | 0.689 | 0.542 |
| <i>Magnolia soulangeana</i> | 0.044 | 0.004 | 0.062 | -0.009 | | 0.767 | 0.729 | 0.516 |
| <i>Nymphaea odorata</i> | 0.027 | -0.013 | 0.045 | -0.026 | -0.017 | | 0.914 | 0.668 |
| <i>Zea mays</i> | 0.018 | -0.022 | 0.036 | -0.035 | -0.026 | -0.009 | | 0.756 |
| <i>Arabidopsis thaliana</i> | -0.018 | -0.058 | 0 | -0.071 | -0.062 | -0.045 | -0.036 | |

Table 10: Results of relative rate tests of substitution rates (K) of amino acid sequences of RPB1 regions D-F between animals, plants and fungi, based upon the alignments shown in Figure 9 and also in Appendix 3. Since three pairwise comparisons were made in this analysis, $\alpha = 0.05/3 = 0.017$. *Porphyra* is the outgroup.

| Group 1 | Group 2 | ΔK | p-value |
|----------------|----------------|------------------------------|----------------|
| animals | plants | -0.136 | 0.397 |
| animals | fungi | -0.027 | 0.878 |
| plants | fungi | 0.108 | 0.526 |

Table 11: Results of relative rate tests of substitution rates (K) of amino acid sequences of RPB1 regions F-H between animals, plants and fungi, based upon the alignment shown in Figure 2. Since three pairwise comparisons were made in this analysis, $\alpha = 0.05/3 = 0.017$. *Dictyostelium* is the outgroup.

| Group 1 | Group 2 | ΔK | p-value |
|----------------|----------------|------------------------------|----------------|
| animals | plants | 0.177 | 0.264 |
| animals | fungi | -0.065 | 0.642 |
| plants | fungi | 0.242 | 0.142 |

Table 12: Results of relative rate tests of substitution rates (K) of amino acid sequences of regions D-F of RPA1, RPB1 and RPC1 in eukaryotes, based upon the alignments shown in figures 9 and in Appendix 4. Since three pairwise comparisons were made in this analysis, $\alpha = 0.05/3 = 0.017$. *Sulfolobus* is the outgroup.

| Group 1 | Group 2 | ΔK | p-value |
|----------------|----------------|------------------------------|----------------|
| RPA1 | RPB1 | -0.221 | 0.34 |
| RPA1 | RPC1 | -0.268 | 0.253 |
| RPB1 | RPC1 | 0.048 | 0.825 |

Copy Number of Plant RPB1s:

Southern analysis of *Magnolia* and *Nymphaea* genomic DNA probed with homologous RPB1 D-F clones by G. Drouin showed the presence of multiple signals. The presence of hybridization signals 7-11 kb in size verifies the source of the cloned RPB1 PCR products. The largest multiple bands indicate that multiple copies of RPB1 genes exist in these genomes. The presence of several smaller bands in lanes 1 and 3 of Figure 24A and in lane 2 of Figure 24B may indicate that some of the copies of RPB1 were cut internally. Table 13 shows the ploidy of the plants used in this study. Both *Magnolia soulangeana* and *Nymphaea odorata* are polyploid, with chromosome numbers of $6n = 114$ and $12n = 84$, respectively. It is possible that the multiple copies of RPB1 are a result of this polyploidy and/or gene duplication(s). Since the blots were washed under moderate stringency, the degree of identity of the hybridization signals with the probe is not known, so RPB1 pseudogenes may have hybridized if they were present. These two Southern blots could not specifically rule out the possibility that paralogous RPB1 sequences from *Magnolia soulangeana* and *Nymphaea odorata* were being used for phylogenetic analysis.

Initial hybridizations with the Northern blots (by B.M.) did not work (results not shown) and were not repeated, but the blots have been stored for future use. Southern blots of *Cycas* and *Ginkgo* (by B.M., results not shown) total DNA digested with *EcoRI*, *BamHI*, or a double digest of both enzymes, and probed with homologous probes of RPB1 D-F cloned PCR products revealed hybridization of 8-kb bands, which could contain the entire coding region and intron sequences of RPB1. However, these membranes had a lot of background hybridization and it is difficult to confirm or rule out the presence of smaller bands such as those seen in Figure 24. Southern analysis of *Araucaria* total DNA did not work as a result of prevailing difficulty in obtaining sufficiently

clean DNA for electrophoresis of a small, concentrated sample volume. Southern blots were prepared for *Osmunda*, *Psilotum*, *Marchantia* and *Equisetum* but have not been hybridized and are stored for later use. Southern analysis of *Ephedra* and maize did not work and needs to be repeated.

It is important to note three points in favor of orthology of other RPBI sequences used in this analysis. First of all, RPBI is known to exist as a single copy in all eukaryotes besides *Trypanosoma brucei* and *Glycine max*, although this has not been investigated further in diploid RPBI sequences cloned in this study. The expressed genes from *Arabidopsis thaliana*, *Zea mays*, and *Spirogyra* were used for the phylogenetic analysis of RPBI in plants. The three intron positions and the length of the coding sequence in regions D-F found previously in *Arabidopsis* and in *Spirogyra* were perfectly conserved in all of the plant RPBI sequences.

Table 13 identifies several (*) species from which RPBI PCR products were not cloned. Prevailing difficulty was experienced in isolating sufficiently clean and abundant samples of DNA from *Osmunda*, *Equisetum*, *Marchantia* and *Psilotum*, but most of all with *Osmunda*. Cloning RPBI D-F was also complicated by the PCR amplification of multiple PCR products, which probably indicate gene duplications compounded by the ancient polyploidy (Stebbins, 1950) of these species. RPBI-like clones were never obtained from *Gnetum* and *Welwitschia*.

Figure 24 A and B: Southern blots of *Magnolia soulangeana* (A) and *Nymphaea odorata* (B) total genomic DNA digests probed with homologous clones of regions D-F of RPBI. Lanes 1, 2 and 3 represent single digests with *EcoRI*, *PstI* and *XbaI*. The λ BstEII size marker is in the centre.

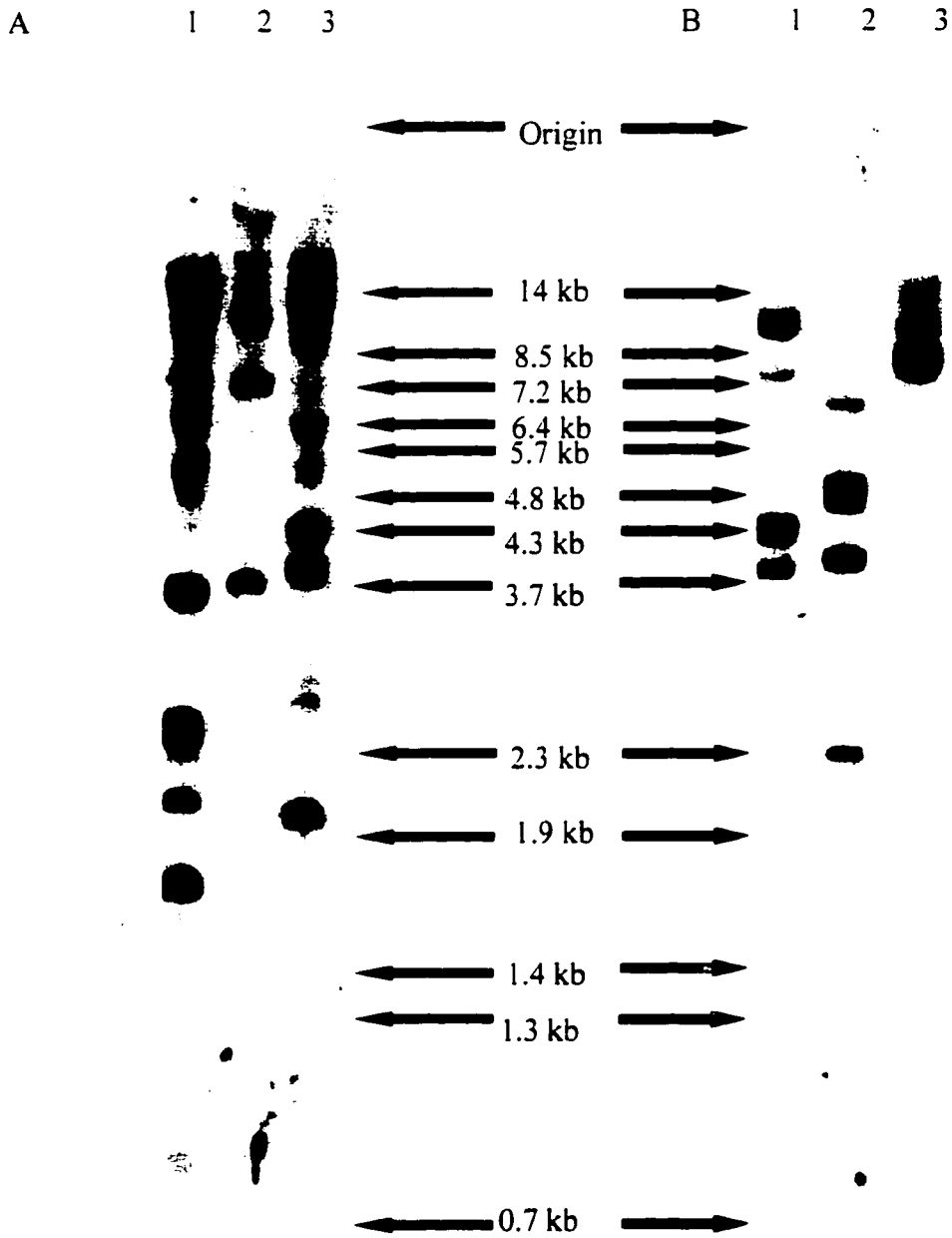


Table 13: Ploidy of representatives of some major plant phyla (Bennett *et.al.*, 1982; Bennet and Smith, 1976; Darlington and Ammal, 1945; Sparrow *et.al.*, 1972; Stebbins, 1950).

| <u>SPECIMEN</u> | <u>CHROMOSOME NUMBER</u> |
|--------------------------------|---|
| <i>Arabidopsis thaliana</i> | n=5 (2n) |
| <i>Araucaria heterophylla</i> | n=13 (those species listed were all 2n) |
| <i>Cycas revoluta</i> | n=11,12 (2n) |
| <i>Ephedra viridis</i> | n=7 (2n or 4n) |
| <i>Equisetum arvense*</i> | n=108, polyploid |
| <i>Ginkgo biloba</i> | n=12 (2n) |
| <i>Gnetum gnemon*</i> | n=? (those <i>Gnetum</i> species listed were 2n=24) |
| <i>Magnolia soulangeana</i> | n=19 (6n) |
| <i>Marchantia polymorpha*</i> | n=9, ploidy uncertain |
| <i>Nymphaea odorata</i> | n=7, 14 (12n=84) |
| <i>Osmunda cinnamomea*</i> | n=22, likely polyploid |
| <i>Psilotum nudum*</i> | n=104, polyploid |
| <i>Welwitschia mirabilis *</i> | n=7 (6n or 12n found) |
| <i>Zea mays</i> | n=10 (2n) |

*species from which RPBI was not cloned.

DISCUSSION

Phylogenetic Analysis of DNA Sequences of Regions D-F of RPB1 in Plants:

Only one relationship amongst plants was consistently resolved using three different methods of phylogenetic analysis of RPB1 D-F nucleotide sequences. *Cycas* and *Ginkgo* were always grouped together, with strong support. On the basis of this analysis, due to the conflicting phylogenies obtained using different methods to analyse the same dataset, neither the woody nor herbaceous origin of angiosperms can be supported. Similarly, the position of the Gnetales relative to gymnosperms and angiosperms is also unclear. The grouping of *Cycas* and *Ginkgo* is consistent with results found by Chaw *et al* (1997) with an 18S rRNA dataset and by Goremykin *et al* (1996) with an *rbcL* and chloroplast intergenic transcribed spacer (cpITS) sequence dataset.

The data about nucleotide compositional bias in plant RPB1 D-F sequences, described in figures 19 through 21, reveals that compositional bias likely had little effect upon this set of phylogenetic analyses. The trend of nucleotide composition of the coding region had negligible impact upon the second base of codons in this sequence, nor does it impact the first codon position after the neutral symbols "m" and "y" were substituted for the first base of arginine and leucine codons. It is also unlikely that the inconsistencies in the trees in figures 6-8 result from the similar codon usage patterns shared by seven of the eight ingroup plant taxa, which are simply a unifying characteristic of plant RPB1 D-F. The short length of this sequence (556 nonsynonymous sites) and the conserved nature of this part of the coding region of RPB1 probably did not provide enough informative variation to clearly infer the phylogeny of these plants (Nei, 1996).

codons. It is also unlikely that the inconsistencies in the trees in figures 6-8 result from the similar codon usage patterns shared by seven of the eight ingroup plant taxa, which are simply a unifying characteristic of plant RPB1 D-F. The short length of this sequence (556 nonsynonymous sites) and the conserved nature of this part of the coding region of RPB1 probably did not provide enough informative variation to clearly infer the phylogeny of these plants (Nei, 1996).

It also appears that the distance and parsimony trees constructed with the RPB1 nucleotide dataset were more prone to “long branch attraction” than was the PUZZLE maximum-likelihood (ML) method. With the former two methods, the most different sequences fell to the bottom of the trees. The latter method separated *Ephedra* and *Arabidopsis*, which both had long branches, away from the bottom of the tree (where both of the other long branches lie) and positioned them closer within the ingroups. However, even there, these two taxa were still grouped together, albeit with low support. The results shown in figure 8 which use the ML method and account for substitution rate heterogeneity in different parts of the alignment were the most robust results obtained from the nucleotide data. However, this method is not enough to resolve relationships in the absence of sufficient informative variation.

Phylogenetic Analysis of Amino Acid Sequences of Regions D-F of RPB1 in Plants:

Figure 13(A) represents the most robust tree obtained from analysis of RPB1 D-F in plants, since it inferred specific relationships of seven of the eight ingroups with one another, each with 55% or greater statistical support, unlike any of the previous trees in this study. This shows the difference made by using a method which accounts for substitution rate heterogeneity amongst sequences.

sequences by Goremykin *et al* (1996) and Hasebe (1992) and by different analyses by Hansen *et al* (1999) of a 9149 bp alignment which included several chloroplast genes and nuclear-encoded 18S rRNA. This relationship is also supported by the analysis of the MADS-box multigene family by Winter *et al* (1999) and by the analysis of a concatenated dataset of genes encoded in the nucleus, mitochondria and chloroplasts by Qiu *et al* (1999).

Figure 13 also weakly supports the herbaceous origin of dicots, but the position of monocots as the most basal angiosperms conflicts with most other molecular and morphological datasets. Parsimony analysis (figure 10 A and B) also weakly supports a herbaceous origin of angiosperms in two of four equally parsimonious trees, while neighbor-joining (figure 11) cannot resolve this relationship. RPB2 analysis (Denton *et al*, 1998) also supports a herbaceous origin of angiosperms but further comparisons of plant RPB2 and RPB1 analyses are limited by the different taxa used in each dataset. With *Hordeum* representing monocots in the RPB2 tree, rather than *Zea* in the RPB1 tree, RPB2 placed monocots amongst the woody dicots rather than at the base of the angiosperm tree. Recent analyses of sequences encoded in the nucleus, mitochondria and chloroplasts agree that angiosperms had a herbaceous origin, placing *Nymphaea* as one of the earliest angiosperms (Mathews and Donoghue, 1999; Parkinson *et al*, 1999; Qiu *et al*, 1999; Soltis *et al*, 1999; see Appendix 5 and discussion below).

As with the parsimony and neighbor-joining phylogenetic trees inferred from plant RPB1 nucleotide data, the parsimony and neighbor-joining amino acid trees also appear to show “long branch attraction”. This flaw is shaken apart somewhat using the Blosum62 model of substitution in the maximum-likelihood analysis, which draws *Arabidopsis* amongst other angiosperms rather than stuck amongst gymnosperms. The JTT (Jones, Taylor, Thornton) model of substitution not

only does this, however, but also resolves more relationships from this dataset, with high statistical support.

The discrepancy between distance and parsimony methods is discussed in detail by Kuhner and Felsenstein (1994) as well as by Tateno *et al* (1994) concluding that the difference in substitution rates amongst different taxa can cause inaccuracy in parsimony analysis, and generally that in the case of high transition/transversion ratios, maximum likelihood was more accurate than neighbour-joining, and these methods were more efficient than parsimony methods.

While the Blosum62 and JTT models of substitution used with the ML method in PUZZLE are both suitable for amino acid data encoded by nuclear genes, the Blosum62 matrix is intended for use with more distantly related amino acid sequences (Strimmer and von Haeseler, 1999). The differences in the trees shown in figures 12 and 13 are probably due to the difference in the suitability of the model chosen to analyse that dataset. The relationships of plant RPB1 sequences are not as decisively resolved by the Blosum62 model as they are by the JTT model, so these sequences may be too closely related (or too conserved) for the Blosum62 model to work properly. In contrast, when the JTT and Blosum62 models are applied to a dataset of much more distantly related RPB1 D-F sequences, as shown in figures 17 and 18, only the Blosum62 model is able to relate slime molds as a sister group to green algae and plants. In this case, the JTT model could not resolve the relationships of any of the major groups of eukaryotes with one another.

Several other points of criticism arise with respect to the phylogenetic tree topologies reconstructed in this study. It was found previously (Wainright *et al*, 1993) that the absence of early branching lineages within each of the major groups to be classified, the use of partial sequences, and not using a closely and distinctively related outgroup (Rodrigo *et al*, 1994) all

detracted from the support for the grouping of animals and fungi as a monophyletic group and animals as monophyletic using 18S rRNA sequences as a phylogenetic marker. The phylogenetic trees presented here using regions D-F of RPB1 as a phylogenetic marker suffer from these problems.

Three criteria need to be improved to obtain a better estimate of the evolutionary relationships of land plants using RPB1. First of all, *Spirogyra* is not a good outgroup to the existing plant RPB1 D-F dataset because it is too distantly related to the ingroups. This is demonstrated in tables 5 and 6, which show that too much homoplasy exists between the outgroup and several ingroups. Second, this analysis lacks vascular plants such as *Marchantia*, *Equisetum* and *Psilotum*, which are earlier-diverging land plants and intermediate to the current outgroup and ingroups. Inclusion of representatives of other major lineages of gymnosperms and angiosperms, such as the Pinaceae, Rosidae and Asteridae would also help to resolve the relationships of the most basal angiosperm clades. Third, the use of short, partial sequences detracts from this analysis. Nei (1996) reiterates that “methods such as neighbor joining, likelihood, and parsimony methods produce reasonably good phylogenetic trees when a sufficiently large number of nucleotides or amino acids are used”. Examination of the alignments in Appendices 1 and 2 shows few synapomorphies. This is particularly detrimental to parsimony analysis and resulted in multiple equally most parsimonious trees. A longer sequence from RPB1 would provide additional informative sites for analysis.

Since there is evidence that multiple copies of the RPB1 gene exist in *Magnolia* and *Nymphaea*, it is possible that this phylogenetic analysis suffers from the use of non-orthologous genes. These two plants are polyploids (6n and 12n), and it is unknown which copies of RPB1 identified by Southern analysis are expressed or which might be pseudogenes. The *Arabidopsis* and

Zea mays RPBI sequences used here are from the expressed RPBI gene, and *Arabidopsis* only has a single copy of RPBI (Dietrich *et al*, 1990). All of the plant species used in this analysis are diploid except for *Magnolia* and *Nymphaea*. While the use of cDNA sequences of RPBI for phylogenetic analysis would not exclude the possibility of comparing non-orthologous genes, it would be a good start and could be followed by a collection of more RPBI sequence data from the same plants so that orthologous sequences might be found.

Several recent plant molecular evolution studies included a broad sampling of angiosperm taxa and several closely-related outgroups. They corroborated one another in finding that herbaceous dicots, including *Amborella* and *Nymphaea*, are the earliest extant angiosperm species. The parsimony analyses of nuclear-encoded phytochrome A and C genes from 26 angiosperms (Mathews and Dongohue, 1999) supported an herbaceous origin of angiosperms whether phytochrome A, C or A and C sequences were analysed. Parkinson *et al* (1999) performed a maximum-likelihood analysis, accounting for site-to-site rate heterogeneity, of a 6564 bp alignment of *rbcL* genes encoded in chloroplasts, small subunit rRNA genes encoded in the nucleus and small subunit rRNA, *coxI* and *rps2* genes encoded in mitochondria. This analysis included 51 taxa and had several gymnosperms (*Cycas*, *Ginkgo* and Coniferales) as the outgroups to angiosperms. Qiu *et al* (1999) also performed a parsimony analysis of a 8733 bp alignment of 18S rRNA genes encoded in the nucleus, *atpI* and *matR* genes encoded in mitochondria and *atpB* and *rbcL* genes encoded in chloroplasts of 105 angiosperm and gymnosperm species. Gnetales and several gymnosperms, including *Cycas*, *Ginkgo* and Coniferales were specified as the outgroup to angiosperms. Soltis *et al* (1999) performed an analysis similar to that of Qiu *et al* (1999), but less

the mitochondrial genes, which enabled them to compare sequences from 567 angiosperm and gymnosperm species.

While these most recent and large-scale studies provide very important information, it is clear that additional data from a protein-coding nuclear gene is required to corroborate these findings. Phylogenetic analysis of regions A through H of RPB1 would provide much longer sequences than the combined phytochrome dataset (Mathews and Donoghue, 1999) and would be informative when analysed independently as well as in combination with other nuclear, mitochondrial and chloroplast-encoded genes.

Substitution Analysis:

The substitution analysis of synonymous and non-synonymous sites of plant nucleotide sequences for regions D-F of RPB1 showed that the 278 synonymous sites had large levels of multiple substitutions, but the 556 non-synonymous sites had few if any (0-7) multiple substitutions. So, as long as synonymous sites are excluded, as they are here, little if any homoplasy is introduced into the analysis. The fact that there is very little multiple substitution at the nonsynonymous sites suggests that these sequences would be useful to resolve the evolutionary relationships of land plants if longer sequences were available.

Relative Rate Tests:

While the phylogenetic methods which account for substitution rate heterogeneity offer more resolution to the inference of the evolutionary relationships of plants with RPB1 D-F sequences, relative rate tests of this dataset do not indicate any significantly different rates of amino acid

substitution. Relative rate tests also show that regions D-F do not have significantly different substitution rates when comparing RPB1 of animals to plants to fungi, nor when comparing RPA1 to RPB1 to RPC1 in all eukaryotes.

Phylogenetic Analyses of Regions D-F of RPB1 in Other Eukaryotes:

Both the analyses of domains F-H and of domains D-F of RPB1 in animals, plants and fungi shown in figures 3 and in figures 14-16 consistently grouped animals and plants as sister groups, by parsimony, neighbor-joining, and PUZZLE maximum-likelihood methods. These results are contrary to most of the recent molecular phylogenetic studies of these groups, which support the evolutionary hypothesis that animals are more closely related to fungi than to plants.

Plants and animals were found to form a monophyletic group in molecular phylogenetic analyses of the 90 kDa heat shock protein Hsp90 (Gupta, 1995), multiple rRNA and tRNA species (Gouy and Li, 1989) and Cu-Zn superoxide dismutases (SOD) by Fitch and Ayala (1994). The analysis of Hsp90 was made by using neighbor-joining and parsimony methods to analyse a 620 amino acid alignment of 31 sequences, with approximately 40% identical sites. Gouy and Li's (1989) analysis was also made using neighbor-joining and parsimony methods. The analysis of SOD sequences by Fitch and Ayala (1994) using Fitch's ANCESTOR program and accounting for rate heterogeneity was made using an alignment of 67 amino acid sequences 118 residues in length, with the invariable sites removed. Fitch and Ayala point out that without examining the possibility that subsets of amino acids used in a phylogenetic analysis may evolve at different rates, errors can be made in estimating divergences by assuming a molecular clock. Gupta (1995) and Gouy and Li (1989) did not take this possibility into account while inferring those phylogenetic trees. A possible

problem with the SOD analysis, however, could be the shortage of informative sites in this relatively short and conserved sequence. The SOD dataset is even shorter in length than the RPBI dataset presented in this study.

Sidow and Thomas (1994) also showed that animals and plants grouped together, excluding fungi, with distance, parsimony and maximum-likelihood analyses of RPBI nucleotide and amino acid sequences. Alignments corresponding to the first and second bases of codons and the inferred amino acid sequence of 2038 codons of RPBI were used for their analyses. These analyses were limited by a sparse sampling (6-10 species) of diverse eukaryotic taxa, however, and the authors did not report that they accounted for rate heterogeneity in their phylogenetic analysis.

Morris (1998) and references therein support some of the relationships seen amongst animals in figures 3 and 14, which are the maximum-likelihood tree of regions F-H of RPBI and a parsimony tree of regions D-F of RPBI in animals, plants and fungi. Morris' paper shows the consensus of fossil and molecular systematic evidence of the phylogenetic relationships of animals. The RPBI F-H tree as well as the parsimony tree of the D-F data separate the deuterostomes (*Mus*) from the protostomes (the other animals). Figure 14 shows arthropods and nematodes grouping together, as a sister group to molluscs, which is also supported by Morris' paper, which also states that these lineages diverged more than 500 million years ago. Since this predates the origin of land plants, this is additional support to the argument that a longer sequence from RPBI is needed to provide additional informative sites from this conserved protein sequence to sort out the evolutionary relationships of plants, since the D-F region does not even provide enough information to resolve the (earlier) relationships of animals with consensus using different methods.

Recent phylogenetic analyses of small subunit ribosomal RNA (SSU rRNA) sequences by maximum-likelihood methods (Wainright *et al.*, 1993 and Van de Peer and De Wachter, 1997) and distance methods, accounting for rate heterogeneity amongst sites (Kumar and Rzhetsky, 1996), in addition to studies of numerous protein-coding genes all support a close relationship of animals with fungi, excluding plants from this group. Baldauf and Palmer (1993) found animals and fungi to be most closely related by both parsimony and distance analyses of sequences for EF-1 α (an elongation factor), α and β -tubulin and actin. This relationship was also found in later analyses of these sequences, i.e. with neighbor-joining analysis of actin by Drouin *et al.* (1995), with neighbor-joining analysis of α and β -tubulin by Keeling and Doolittle (1996) and with EF-1 α using distance, parsimony and maximum-likelihood methods by Baldauf and Doolittle (1997) and Roger *et al.* (1999). The latter two studies utilised the JTT amino acid transition probability matrix adjusted for amino acid frequencies and the PROTML program for the maximum-likelihood analysis but did not account for rate heterogeneity. While these proteins also have highly conserved sequences, greater length (approximately 100-150 amino acids longer than this RPBI D-F alignment) provided additional informative sites. Analysis of other elongation factors (EF-2 by Hirt *et al.* (1999) and EF-Tu and EF-G by Baldauf *et al.* (1996)) also supported the grouping of animals with fungi. Both analyses used distance, parsimony and PROTML maximum-likelihood methods, and account for rate heterogeneity within sequences, while Hirt *et al.* (1999) also accounted for base composition heterogeneity between sequences in their analyses.

Sidow and Thomas (1994) point out that genes encoding elongation factors may undergo concerted evolution, so the evolutionary relationships of sequences of unrelated proteins should also be studied. Animals and fungi were also found to be sister groups in analyses of genes encoding the

60 kDa chaperonin (Cpn60) protein (Clark and Roger (1995), Roger *et al* (1998)), the 70 kDa heat-shock protein, Hsp70 (Borchiellini *et al* (1998)), aminoacyl-tRNA synthetases (Brown and Doolittle, 1995) and translation initiation factor eIF-2 γ (Keeling *et al* (1998)).

These are only a few of some major examinations of the relationships amongst plants, fungi, and animals using molecular data. It is difficult to group two of these groups together with consensus because they have emerged during very proximal periods of time. This is also reflected by the weak support for many internal nodes in the analyses of RPBI shown in figures 14-18.

The analyses represented by figures 14-18 were dominated by taxa which were very closely related amongst their own groups, and each of these groups were very distantly related to one another. This makes it likely that there is a shortage of synapomorphies amongst closely related taxa, while homoplasies have accumulated between distantly related groups. The topology of the RPBI tree might be different if a more broad representation of sequences were included, including more early-branching lineages in each of the major eukaryotic kingdoms that would effectively be intermediate taxa amongst the major groups. This would reduce the homoplasies seen from one taxon to the next and might thus reduce the overall substitution rate heterogeneity. Unequal rates of substitution amongst organisms might skew the dataset to misrepresent the positions of taxa in a phylogenetic tree (Baldauf and Palmer, 1993).

The close relationship of RPBI to RPAI rather than to RPCI shown in figures 17 and 18 is supported by neighbor-joining and maximum likelihood analyses of type A RNA polymerases by a set of complete sequences by Iwabe *et al* (1991), as well as by the parsimony and “parsimonious low sum of squares” trees inferred by Klenk *et al* (1995) from a 1168-amino acid alignment. The distance analysis by Klenk *et al* (1995) placed RPCI and RPAI as sister groups, but the analyses

in that paper were limited by a number of partial sequences, a sparse sampling of taxa and by not accounting for rate heterogeneity.

Phylogenetic analyses of domains D-G, A-G and A-H of RPB1 from 12-16 distantly related eukaryotes by Stiller and Hall (1997, 1998), Stiller *et al* (1998) and Hirt *et al* (1999) showed trees with stronger statistical support than the analyses of shorter RPB1 sequences (D-F) from a more broad selection of taxa (38 species) shown here in Figures 17 and 18. However, depending upon the selection of taxa, the alignment, and the method used to construct the trees in these studies, the relationships of animals, plants and fungi to one another were not always resolved in these analyses either, although some relationships of deeply diverging protist and algal lineages were always resolved. Stiller and Hall (1998) and Stiller *et al* (1998) showed that with consensus trees from parsimony, distance and maximum likelihood analyses of regions A-H of RPB1, red algae were distinctly separated from basal eukaryotic lineages and formed a sister group to the “crown” eukaryotic taxa. In this “crown” group, however, the relationships of animals, plants, fungi and slime molds to one another were not resolved. The relationships amongst the 3-4 most basal lineages (*Giardia*, *Trypanosoma*, *Mastigamoeba*, and *Trichomonas*) also were not well-resolved in these trees. Stiller and Hall (1997) found that phylogenetic trees of regions D-G of RPB1 inferred by both parsimony and distance methods found fungi as the sister group to a clade consisting of a plant, a green alga and a slime mold. Their maximum-likelihood analysis of regions A-H (excluding the green alga), however, placed fungi as a sister group to animals. Also, while parsimony and distance methods showed congruent relationships of basal eukaryotes (in considering regions D-G), maximum likelihood (considering regions A-H) resolved the topology of these relationships differently.

In the phylogenetic analysis of 15 sequences of regions A-G of RPB1 using a combination of parsimony and maximum likelihood methods, Hirt *et al* (1999) showed well-resolved relationships amongst the basal eukaryotic taxa, and a strongly supported relationship of Microsporidia with fungi. However, the fungi + Microsporidia clade was only placed as a sister group to animals with <50% bootstrap support. They also found that removal of fast-evolving sites from the alignment increased the bootstrap support in parsimony analyses, which are more sensitive to unequal substitution rates. For maximum-likelihood analysis, though, Hirt *et al* (1999) found that strong bootstrap support for deep relationships (amongst basal eukaryotic lineages) depended upon the inclusion of invariable sites in the alignment and was greatly reduced without them. These authors also found that by using methods which account for base composition heterogeneity between sequences and for site-by-site substitution rate heterogeneity, phylogenetic analyses of sequences for regions A-G of RPB1 yielded more conclusive results than did similar analyses of EF-1 α and EF-2 sequences, which were more strongly affected by heterogeneous base composition and substitution rates.

The work of Stiller and Hall (1997, 1998), Stiller *et al* (1998) and Hirt *et al* (1999) further demonstrates several key issues about the use of RPB1 for studying evolutionary relationships. The analyses of regions A-G and A-H gave similar results, with relationships that were better resolved than the analyses of regions D-F shown in figures 17 and 18. These papers also demonstrate that the use of phylogenetic methods which account for base composition heterogeneity between sequences and for substitution rate heterogeneity within sequences when comparing RPB1 sequences of distantly related organisms provides a more robust estimate of their evolutionary relationships than can otherwise be obtained. These key points suggest that the use of longer (A-H)

RPB1 sequences, analysed using the methods of Hirt *et al* (1999) would be more useful for inferring the evolutionary relationships of land plants.

Conclusions:

Although a single and congruent phylogeny of land plants was not be found here due to the narrow representation of taxa in this analysis and the truncated and conserved sequences, it is likely that the phylogenies presented here using the maximum likelihood methods (figures 12 and 13) are the most accurate. Analysis of longer RPB1 sequences from a more broad range of taxa, taking site-to-site as well as taxon-to-taxon substitution rate heterogeneity into account, will be needed to further establish the usefulness of RPB1 for studying the evolutionary relationships of land plants. This should include more vascular plants, additional angiosperm and gymnosperm taxa, and a more closely related outgroup or outgroups. Regions A-D of RPB1 should also be included in such an analysis, because this region is less conserved than the D-F region and should provide additional informative nonsynonymous substitutions. Although domains D-F of RPB1 are too short to provide enough information to resolve the evolutionary relationships of the major groups of land plants, the non-synonymous substitution rate of 2×10^{-10} nonsynonymous substitutions per site per year (Figure 23 and Tables 9-12) and paucity of multiple substitutions at non-synonymous sites (Table 7) indicates that RPB1 could be a good phylogenetic marker to resolve the evolutionary relationships of land plants.

REFERENCES:

- Ahearn J, Bartolmei M, West M, Cisek L, Corden J. (1987) Cloning and sequence analysis of the mouse genomic locus encoding the largest subunit of RNA polymerase II. *J Biol Chem* 262: 10695-10705.
- Allison L, Moyle M, Shales M, Ingles C. (1985) Extensive homology amongst the largest subunits of eukaryotic and prokaryotic RNA polymerases. *Cell* 42:599-610. yeast seq.
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool. *J Mol Biol* 215: 403-410.
- Altschul SF, Madden T, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucl Acids Res* 25:3389-3402
- An SS, Möpps B, Weber K, Bhattacharya D. (1999) The origin and evolution of green algal and plant actins. *Mol Biol Evol* 16: 275-285.
- Azuma Y, Yamagishi M, Ueshima R, Ishihama A (1991) Cloning and sequence determination of the *Schizosaccharomyces pombe rpb1* gene encoding the largest subunit of RNA polymerase II. *Nucl Acids Res* 19: 461-468.
- Baldauf SL, Palmer JD. (1993) Animals and fungi are each other's closest relatives: congruent evidence from multiple proteins. *Proc Natl Acad Sci USA* 90: 11558-11562.
- Baldauf SL, Palmer JD, Doolittle WF (1996) The root of the universal tree and the origin of eukaryotes based on elongation factor phylogeny. *Proc Natl Acad Sci USA* 93: 7749-7754.
- Baldauf SL, Doolittle WF. (1997) Origin and evolution of the slime molds (Mycetozoa). *Proc Natl Acad Sci USA* 94: 12007-12012.
- Baldwin JG, Frisse LM, Vida JT, Eddleman CD, Thomas WK. (1997) An evolutionary framework for the study of developmental evolution in a set of nematodes related to *Caenorhabditis elegans*. *Mol Phyl Evol* 8:249-259.
- Bateman RM, Crane PR, DiMichele WA, Kenrick PR, Rowe NP, Speck T, Stein WE. (1998) Early evolution of land plants: Phylogeny, physiology and ecology of the primary terrestrial radiation. *Annu Rev Ecol Syst* 29: 263-92.
- Bennett M.D., Smith J.B. (1976) Nuclear DNA amounts in angiosperms. 274: 227-274

- Bennett M.D., Smith J.B., Heslop-Harrison J.S. (1982) Nuclear DNA amounts in angiosperms. *Proc. Roy. Soc. Lond.* 216: 179-199.
- Bird DM, Riddle DL (1989) Molecular cloning and sequencing of *ama-I*, the gene encoding the largest subunit of *Caenorhabditis elegans* RNA polymerase II. *Mol Cell Biol* 9: 4119-4130.
- Borchiellini C, Boury-Esnault N, Vacelet J, Parco YL. (1998) Phylogenetic analysis of the Hsp70 sequences reveals the monophyly of Metazoa and the specific phylogenetic relationships between animals and fungi. *Mol Biol Evol* 15: 647-655.
- Bremer K (1985) Summary of green plant phylogeny and classification. *Cladistics* 1: 369-385.
- Brown J, Doolittle WF. (1995) Root of the universal tree of life based on ancient aminoacyl-tRNA synthetase gene duplications. *Proc Natl Acad Sci USA* 92: 2441-2445.
- Chang S, Puryear J, Cairney J. (1993) A simple method for isolating RNA from pine trees. *Plant Mol Biol Rep* 11:113-116
- Chase MW, Soltis DE, Olmstead RG, Morgan D, Les DH, Mishler BD, Duvall MR, Price RA, Hills HG, Qiu Y-L, Kron KA, Rettig JH, Conti E, Palmer JD, Manhart JR, Sytsma KJ, Michaels HJ, Kress WJ, Karol KG, Clark WD, Hedren M, Gaut BS, Jansen RK, Kim K-J, Wimpee CF, Smith JF, Furnier GF, Strauss SH, Xiang Q-Y, Plunkett GM, Soltis PS, Swensen SM, Williams SE, Gadek PA, Quinn CJ, Eguiarte LE, Golenberg E, Learn GH Jr., Graham SW, Barrett SCH, Dayanandan S, Albert VA. (1993) Phylogenetics of seed plants: An analysis of nucleic acid sequences from the plastid gene *rbcL*. *Ann Miss Botanic Gard* 80: 528-580.
- Chaw S-M, Zharkikh A, Sung H-M, Lau T-C, Li W-H. (1997) Molecular phylogeny of extant gymnosperms and seed plant evolution: analysis of nuclear 18S rRNA sequences. *Mol Biol Evol* 14: 56-68.
- Cho Y, Qiu Y-L, Kuhlman P, Palmer JD. (1998) Explosive invasion of plant mitochondria by a group I intron. *Proc Natl Acad Sci USA* 95: 14244-14249.
- Clark CG, Roger AJ. (1995) Direct evidence for secondary loss of mitochondria in *Entamoeba histolytica*. *Proc Natl Acad Sci USA* 92: 6518-6521.
- Crane PR, Friis EM, Pedersen KR. (1995) The origin and early diversification of angiosperms. *Nature* 374: 27-33.
- Croan D, Ellis J. (1996) Phylogenetic relationships between *Leishmania*, *Vianna* and *Sauvoleishmania* inferred from comparison of a variable domain within the RNA polymerase II largest subunit gene. *Mol Biochem Parasit* 79: 97-102.

- Cronquist A. (1981) *An integrated system of classification of flowering plants*. Columbia University Press, New York. p. 51.
- Darlington CD, Ammal EKJ. (1945) *Chromosome atlas of cultivated plants*. George Allen & Unwin, London. pp. 38-63.
- Denton AL, McConaughy BL, Hall BD. (1998) Usefulness of RNA polymerase II coding sequences for estimation of green plant phylogeny. *Mol Biol Evol* 15: 1082-1085.
- Dietrich MA, Prenger JP, Guilfoyle T. (1990) Analysis of the genes encoding the largest subunit of RNA polymerase II in *Arabidopsis* and soybean. *Plant Mol Biol* 15: 207-223.
- Donoghue MJ. (1994) Progress and Prospects in reconstructing plant phylogeny. *Ann Missouri Bot. Gard.* 81: 405-418.
- Donoghue MJ, Mathews S. (1998) Duplicate genes and the root of angiosperms, with an example using phytochrome sequences. *Mol Phyl Evol* 9:489-500.
- Doyle JJ, Doyle JL. (1990) Isolation of plant DNA from fresh tissue. *Focus* 12:13-15.
- Drouin G, Prat F, Ell MJ, Clarke PTG. (1999) Detecting and characterizing gene conversion events between multigene family members. *Mol Biol Evol* 16:1369-1390.
- Drouin G, Moniz de Sá M, Zuker M. (1995) The *Giardia lamblia* actin gene and the phylogeny of eukaryotes. *J Mol Evol* 41:841-9.
- Evers R, Hammer A, Kock J, Jess W, Borst P, Memet S, Cornelissen AW (1989) *Trypanosoma brucei* contains two RNA polymerase II largest subunit genes with an altered C-terminal domain. *Cell* 56(4):585-97.
- Felsenstein J. (1993) PHYLIP v3.52c, University of Washington, Seattle.
- Fitch WM, Ayala FJ. (1984) The superoxide dismutase molecular clock revisited. *Proc Natl Acad Sci USA* 91:6802-6807.
- Foster PG, Hickey DA (1999) Compositional bias may affect both DNA-based and protein-based phylogenetic reconstructions. *J Mol Evol* 48:284-90 .
- Garber AT. (1993) *Biofeedback* 14: 148-9
- Gehrig HH, Heute V, Kluge M (1998) Toward a better knowledge of the molecular evolution of phosphoenolpyruvate carboxylase by comparison of partial cDNA sequences. *J Mol Evol* 46: 107-114.

- Goremykin V, Bobrova V, Pahnke J, Troitsky A, Antonov A, Martin W. (1996) Noncoding sequences from the slowly evolving chloroplast inverted repeat in addition to *rbcL* data do not support Gnetalean affinities of angiosperms. *Mol Biol Evol* 13: 383-396.
- Gouy M, Li WH. (1989) Molecular phylogeny of the kingdoms animalia, plantae, and fungi. *Mol Biol Evol.* 6: 109-122.
- Gupta R. (1995) Phylogenetic analysis of the 90kD heat shock family of protein sequences and an examination of the relationship among animals, plants, and fungi species. *Mol Biol Evol.* 12: 1063-1073.
- Hansen A, Hansmann S, Samigullin T, Antonov A, Martin W. (1999) *Gnetum* and the angiosperms: molecular evidence that their shared morphological characters are convergent, rather than homologous. *Mol Biol Evol* 16: 1006-1009.
- Hasebe M, Kofuji R, Ito M, Kato M, Iwatsuki K, Ueda K. (1992) Phylogeny of gymnosperms inferred from *rbcL* gene sequences. *Bot Mag Tokyo* 105: 673-679.
- Hasebe M, Ito M, Kofuji R, Ueda K, Iwatsuki K. (1993) Phylogenetic relationships of ferns deduced from *rbcL* gene sequence. *J Mol Evol* 37: 476-482.
- Hasegawa, M., Hashimoto, T. (1993) rRNA trees misleading? *Nature* 361:23.
- Henikoff S, Henikoff JG. (1992) Amino acid substitution matrices from protein blocks. *Proc Natl Acad Sci USA* 89: 10915-10919.
- Higgins DG. (1994) CLUSTAL V: Multiple alignment of DNA and protein sequences. *Methods Mol Biol* 25: 307-18.
- Hirt RP, Logsdon JM Jr, Healy B, Dorey MW, Doolittle WF, Embley TM. (1999) Microsporidia are related to fungi: evidence from the largest subunit of RNA polymerase II and other proteins. *Proc Natl Acad Sci USA* 96:
- Iwabe N, Kuma KI, Kishino H, Hasegawa M. (1991) Evolution of RNA polymerases and branching patterns of the three major groups of archaebacteria. *J Mol Evol.* 32:70-78.
- Jess W, Hammer A, Cornelissen AW (1989) Complete sequence of the gene encoding the largest subunit of RNA polymerase I of *Trypanosoma brucei*. *FEBS Lett* 249:123-8.
- Jokerst RS, Weeks JR, Zehring WA, Greenleaf AL. (1989) Analysis of the gene encoding the largest subunit of RNA polymerase II in *Drosophila*. *Mol Gen Genet* 215:266-275.
- Jones DT, Taylor WR, Thornton JM (1992) The rapid generation of mutation data matrices from protein sequences. *CABIOS* 8: 275-282.

- Keeling PJ, Fast N, McFadden GI. (1998) Evolutionary relationship between translation initiation factor eIF-2 γ and selenocysteine-specific elongation factor SELB: change of function in translation factors. *J Mol Evol* 47: 649-655.
- Keeling PJ, Doolittle WF. (1996) Alpha-tubulin from early-diverging eukaryotic lineages and the evolution of the tubulin family. *Mol Biol Evol* 13: 1297-1305.
- Kenrick P, Crane PR. (1997) The origin and early evolution of plants on land. *Nature* 389: 33-39.
- Klenk HP, Zillig W. (1994) DNA-dependant RNA polymerase subunit B as a tool for phylogenetic reconstructions: branching topology of the archaebacterial domain. *J Mol Evol.* 38:420-432.
- Klenk HP, Zillig W, Lanzendorfer M, Grampp B, Palm P. (1995) Location of Protist lineages in a phylogenetic tree inferred from sequences of DNA-dependant RNA polymerases. *Arch Protistenkd.* 145:221-230.
- Knackmuss S, Bautz EF, Petersen G. (1997) Identification of the gene coding for the largest subunit of RNA polymerase I (A) of *Drosophila melanogaster*. *Mol Gen Genet* 253:529-34.
- Kock J, Evers R, Cornelissen AW. (1988) Structure and sequence of the gene for the largest subunit of trypanosomal RNA polymerase III. *Nucl Acids Res* 16: 8753-8772.
- Kolukisaoglu HU, Marx S, Wiegmann C, Hanelt S, Schneider-Poetsch HAW. (1995) Divergence of the phytochrome family predates angiosperm evolution and suggests that *Selaginella* and *Equisetum* arose prior to *Psilotum*. *J Mol Evol* 41: 329-337.
- Kuhner MK, Felsenstein J. (1994) A simulation comparison of phylogeny algorithms under equal and unequal evolutionary rates. *Mol Biol Evol* 11: 459-468.
- Kumar S, Rzhetsky A. (1996) Evolutionary relationships of eukaryotic kingdoms. *J Mol Evol* 42: 183-193.
- Kuzoff RK, Sweere JA, Soltis DE, Soltis PS, Zimmer EA. (1998) The phylogenetic potential of entire 26S rDNA sequences in plants. *Mol Biol Evol* 15: 251-263.
- Lam TY, Chan L, Yip P, Siu CH. (1992) The largest subunit of RNA polymerase II in *Dictyostelium*: conservation of the unique tail domain and gene expression. *Biochem Cell Biol* 70: 792-799. Carleton QP501.C3222
- Lanzendorfer M, Palm P, Grampp B, Peattie DA, Zillig W. (1992) Nucleotide sequence of the gene encoding the largest subunit of the DNA-dependant RNA polymerase III of *Giardia lamblia*. *Nucl Acids Res* 20:1145.

- Li WB, Bzik DJ, Gu HM, Tanaka M, Fox BA, Inselburg J (1989) An enlarged largest subunit of *Plasmodium falciparum* RNA polymerase II defines conserved and variable RNA polymerase domains. *Nucl Acids Res* 17: 9621-36.
- Li W-H. (1993) Unbiased estimation of the rates of synonymous and nonsynonymous substitution. *J Mol Evol* 36:96-9.
- Li W-H, Graur D. (1991) *Fundamentals of Molecular Evolution*. Chapter 4. Sinauer, Sunderland MA.
- Lloyd AT, Sharp PM. (1992) CODONS: a microcomputer program for codon usage analysis. *J Hered* 83: 239-240.
- Loconte H. (1996) "Comparison of alternative hypotheses for the origin of the angiosperms", Chapter 10 in *Flowering plant origin, evolution and phylogeny*. Taylor DW, Hickey LJ, eds. Chapman and Hall, NY. 267-286.
- Malek O, Lattig K, Hiesel R, Brennicke A, Knoop V. (1996) RNA editing in bryophytes and a molecular phylogeny of land plants. *EMBO* 15: 1403-1411.
- Mathews S, Donoghue M. (1999) The root of angiosperm phylogeny inferred from duplicate phytochrome genes. *Science* 286:947-950.
- Meacham C. (1994) Phylogenetic relationships at the basal radiation of angiosperms: further study by probability of character compatibility. *Syst. Bot.* 19: 506-522.
- Memet S, Saurin W, Sentenac A. (1988) RNA polymerases B and C are more closely related to each other than to RNA polymerase A. *J Biol Chem* 263: 10048-10051.
- Mishler BD, Churchill SP. (1985) Transition to a land flora: phylogenetic relationships of the green algae and bryophytes. *Cladistics* 1: 305-328.
- Moniz de Sá M, Drouin G. (1998) Phylogeny and substitution rates of one fern and six gymnosperm actin genes. unpublished
- Moniz de Sá M, Drouin G. (1996) Phylogeny and substitution rates of angiosperm actin genes. *Mol Biol Evol* 13: 1198-1212.
- Morris SC. (1998) Metazoan phylogenies: falling into place or falling to pieces? A paleontological perspective. *Curr Opin Genet Dev* 8:662-667.
- Nawrath C, Schell J, Koncz C. (1990) Homologous domains of the largest subunit of eukaryotic RNA polymerase II are conserved in plants. *Mol Gen Genet* 223:65-75.

- Nei M. (1996) Phylogenetic analysis in molecular evolutionary genetics. *Annu Rev Genet* 30: 371-403.
- Nickrent DL, Soltis DE. (1995) A comparison of angiosperm phylogenies from nuclear 18S rDNA and *rbcL* sequences. *Ann. Missouri Bot Gard* 82: 208-234.
- Parkinson CL, Adams KL, Palmer JD. (1999) Multigene analyses identify the three earliest lineages of extant flowering plants. *Curr Biol* 9: 1485-1488.
- Palenik B (1992) Polymerase evolution and organism evolution. *Curr Opin Genes Dev* 2:931-936.
- Ponce MR, Micol JL. (1992) PCR amplification of long DNA fragments. *Nucl Acids Res* 20:623
- Puhler G, Leffers H, Gropp F, Palm P, Klenk HP, Lottspeich F, Garrett RA, Zillig W. (1989a) Archaeobacterial DNA-dependant RNA polymerases testify to the evolution of the eukaryotic nuclear genome. *Proc Natl Acad Sci USA* 86: 4569-4573.
- Puhler G, Lottspeich F, Zillig W. (1989b) Organization and nucleotide sequence of the genes encoding the large subunits A, B and C of the DNA-dependent RNA polymerase of the archaeobacterium *Sulfolobus acidocaldarius*. *Nucl Acids Res* 17: 4517-4535.
- Purugganan MD. (1997) The MADS-box floral homeotic gene lineages predate the origin of seed plants: Phylogenetic and molecular clock estimates. *J Mol Evol* 45: 392-396.
- Qiu Y-L, Cho Y, Cox JC, Palmer JD. (1998) The gain of three mitochondrial introns identifies liverworts as the earliest land plants. *Nature* 394: 671-674.
- Qiu Y-L, Palmer JD. (1999) Phylogeny of early land plants: insights from genes and genomes. *Trends in Plant Science* 4:26-30.
- Qiu Y-L, Lee J, Bernasconi-Quadroni F, Soltis DE, Soltis PS, Zanis M, Zimmer EA, Chen Z, Savolainen V, Chase MW. (1999) The earliest angiosperms: evidence from mitochondrial, plastid and nuclear genomes. *Nature* 402: 404-407.
- Quon D, Delgadillo M, Johnson P. (1996) Transcription in the early diverging eukaryote *Trichomonas vaginalis*: an unusual RNA polymerase II and alpha-amanitin resistant transcription of protein-coding genes. *J Mol Evol* 43: 253-62.
- Regier JC, Shultz JW. (1997) Molecular phylogeny of the major arthropod groups indicates polyphyly of Crustaceans and a new hypothesis for the origin of hexapods. *Mol Biol Evol* 14: 902-913.
- Ridley M. (1996) *Evolution*. 2nd ed. Blackwell, Cambridge MA. 537-538.

Ritland K, Eckenwalder J. (1992) Chapter 17: Polymorphism, hybridization, and variable evolutionary rate in molecular phylogenies. *in Molecular Systematics of Plants*, P. Soltis, D. Soltis, J. Doyle, eds. Chapman and Hall, New York. p. 404-429.

Robinson M, Gouy M, Gautier C, Mouchiroud D (1998) Sensitivity of the relative-rate test to taxonomic sampling. *Mol Biol Evol* 15: 1091-1098.

Rodrigo AG, Bergquist PR, Bergquist PL. (1994) Inadequate support for an evolutionary link between the metazoa and the fungi. *Syst. Biol.* 43: 578-584.

Roger AJ, Svärd SG, Tovar J, Clark CG, Smith MW, Gillin FD, Sogin ML. (1998) A mitochondrial-like chaperonin 60 gene in *Giardia lamblia*: Evidence that diplomonads once harbored an endosymbiont related to the progenitor of mitochondria. *Proc Natl Acad Sci USA* 95: 229-234.

Roger AJ, Sandblom O, Doolittle WF, Philippe H. (1999) An evaluation of elongation factor 1 α as a phylogenetic marker for eukaryotes. *Mol Biol Evol* 16: 218-233.

Rose T, Schultz E, Henikoff J, Pietrokovski S, MacCallum C, Henikoff S. (1998) Consensus-degenerate hybrid oligonucleotide primers for amplification of distantly related sequences. *Nucl Acids Res* 26: 1628-1635

Sambrook J, Fritsch EF, Maniatis T. (1989) *Molecular cloning: a laboratory manual*, 2nd ed. Cold Spring Harbor Laboratory Press.

Sandhu GS, Precup JW, Kline BC. (1989) Rapid one-step characterization of recombinant vectors by direct analysis of transformed *Escherichia coli* colonies. *BioTechniques* 7: 689-690.

Seither P, Coy JF, Pouska A, Grummet I. (1997) Molecular cloning and characterization of the cDNA encoding the largest subunit of mouse RNA polymerase I. *Mol Gen Genet* 255: 180-186.

Sidow A, Thomas WK. (1994) A molecular evolutionary framework for eukaryotic model organisms. *Curr Biol.* 4: 596-603.

Sidow A, Wilson A. (1990) Compositional statistics: an improvement of evolutionary parsimony and its application to deep branches in the tree of life. *J Mol Evol* 31: 51-68.

Smith SW, Overbeek R, Woese CR, Gilbert W, Gillevet PM (1994) The genetic data environment: an expandable GUI for multiple sequence analysis. *Comput Appl Biosci* 10:671-5.

Soltis DE, Soltis PS, Nickrent DL, Johnson LA, Hahn WJ, Hoot SB, Sweere JA, Kuzoff RK, Kron KA, Chase MW, Swensen SM, Zimmer EA, Chaw S-M, Gillespie LJ, Kress WJ, Sytsma KJ. (1997) Angiosperm phylogeny inferred from 18S ribosomal DNA sequences. *Ann Miss Bot Gard* 84: 1-49.

- Soltis PS, Soltis DE, Chase MW. (1999) Angiosperm phylogeny inferred from multiple genes as a tool for comparative biology. *Nature* 402: 402-404.
- Sparrow AH, Price HJ, Underbrink AG. (1972) A survey of DNA content per cell and per chromosome of prokaryotic and eukaryotic organisms: some evolutionary considerations. *Evolution of genetic systems*, ed by HH Smith. Gordon and Breach, NY. p.451-494.
- Stebbins GL Jr. (1950) *Variation and Evolution in Plants*. Columbia University Press. NY. pp. 298-355.
- Stewart WN, Rothwell GW. (1993) Chapter 32: Major evolutionary events and trends - in retrospect, in *Paleobotany and the evolution of plants*. Cambridge University Press, NY. pp. 505-512.
- Stiller JW, Hall BD. (1997) The origin of red algae: implications for plastid evolution. *Proc Natl Acad Sci USA* 94:4520-4525.
- Stiller JW, Duffield CS, Hall BD. (1998) Amitochondriate amoebae and the evolution of DNA-dependant RNA polymerase II. *Proc Natl Acad Sci USA* 95: 11769-11774.
- Stiller JW, Hall BD. (1998) Sequences of the largest subunit of RNA polymerase II from two red algae and their implications for rhodophyte evolution. *J Phycol* 34: 857-864.
- Strimmer K, von Haeseler A. (1996) Quartet puzzling: a quartet maximum likelihood model for reconstructing tree topologies. *Mol Biol Evol* 13: 964-969.
- Strimmer K, von Haeseler A. (1997) Likelihood-mapping: a simple method to visualize phylogenetic content of a sequence alignment. *Proc Natl Acad Sci USA* 94: 6815-6819.
- Strimmer K, von Haeseler A. (1999) Manual for PUZZLE 4.0.2, distributed with PUZZLE as: puzzle4.0.2\docs\manual.html
- Takhtadzhian AL. (1997) Introduction: Evolutionary Systematics: The Darwinian paradigm, in *Diversity and classification of flowering plants*. Columbia University Press, NY. pp. 1-7.
- Tateno Y, Takezaki N, Nei M. (1994) Relative efficiencies of the maximum likelihood, neighbor-joining, and maximum parsimony methods when substitution rate varies with site. *Mol Biol Evol* 11: 261-277.
- Taylor DW, Hickey LJ (1996) Chapter 1, "Introduction: The challenge of flowering plant history" and Chapter 9, "Evidence for and implications of an herbaceous origin for angiosperms" in *Flowering plant origin, evolution and phylogeny*. Taylor DW, Hickey LJ, eds. Chapman and Hall, NY. 1-7, 232-266.

Thompson JD, Higgins DG, Gibson TJ. (1994) CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucl Acids Res.* 22: 4673-4680.

Van de Peer Y, De Wachter R. (1997) Evolutionary relationships among the eukaryotic crown taxa taking into account site-to-site rate variation in 18S rRNA. *J Mol Evol* 45: 619-630.

Wainright PO, Hinkle G, Sogin ML, Stickel SK. (1993) Monophyletic origins of the metazoa: an evolutionary link with fungi. *Science* 260: 340-342.

Winter K-U, Becker A, Münster T, Kim JT, Saedler H, Theissen G. (1999) MADS-box genes reveal that gnetophytes are more closely related to conifers than to flowering plants. *Proc Natl Acad Sci USA* 96:7342-7347.

Wintzerith M, Acker J, Vicaire S, Vigneron M, Keding C. (1992) Complete sequence of the human RNA polymerase II largest subunit. EMBL-Genbank database, Accession number X63564. *Nucl Acids Res.* 20:910.

Yamagishi M, Nomura M. (1988) Cloning and sequence determination of the gene encoding the largest subunit of the fission yeast *Schizosaccharomyces pombe* RNA polymerase I. *Gene* 74:503-515.

Young RA (1991) RNA polymerase II. *Annu Rev Biochem* 60: 689-715.

Spirogy. II GAGAGTGGCCGAGCGCnGAA TTnTGGCAACCATTGnTGTAA TATGnTnTCACAGGTTAGAT GGATGGGAACATGCGAGGCTC ACATGAACATAAGAACATGC
Arau. het IIT.G.A.....
Cyc. rev. IIG.A.....T.
Eph. vir. IIA.CA.....
Gnk. bil. IIG.T.....
Mag. sou. IIA.A.....T.
Nym. odo. IIA.C.A.....T.
Zea. mays. 2A.T.....
Arab. tha. 2C.A.....T.

Spirogy. II AAGCAAACGAGTAAGAnTAT GAGCGCTGGAAACAnTGAGC CACCGGnGACnTATGATCTT GAAnGGTAAACAGTnTAAAA GCnGGAGAGCGGnGGCGCCCA
Arau. het IIA...T.CA...AG...C.AA.....nGA.AG.....
Cyc. rev. IIC...CAC...GC...G.....A.....A.AG.....
Eph. vir. II G.....AATAAT.CC.A...CA G.....n.....CA.....AAAG.....
Gnk. bil. IIA...CA...GC...G.....A.....A.AG.....
Mag. sou. IIA...CA...A.....G.....A.....A.AG.....
Nym. odo. IIA...CA...A.....G.....A.....A.AG.....
Zea. mays. 2A.AA..CA...G.....A.....A.AG.....
Arab. tha. 2nGCATTCA.G.G...C.G...A.nG..A.....A.AG.....

Spirogy. II TCAGnTTCGAAGAAAAGTAA GCATGTACGGGTCAAGGTC TTATAAAATTCCAATAT
Arau. het II AA.....n.....
Cyc. rev. II AA.....n.....AGC
Eph. vir. II nG.....n.....AGC
Gnk. bil. II nG.....n.....AGC
Mag. sou. II AA.....n.....AGC
Nym. odo. II nG.....n.....AGC
Zea. mays. 2 AA.....n.....AGC
Arab. tha. 2 AA.....G...C...n.....AGC

Appendix 2: 278 position alignment of 9 amino acid sequences of regions D-F of RPB1 in land plants that was used to infer the phylogenetic trees shown in figures 10-13. A green alga, *Spirogyra*, is the outgroup. Residues identical to those in the outgroup are represented by a dot.

100

| | | | | | |
|------------------|----------------------|-------------------------|-----------------------|------------------------|----------------------|
| <i>Spirogyra</i> | HVCQTFETRAETMELMMVPK | CVVSPQSNRPVMGIVQDTLL | GCRKVT"KRDTFIEKDVFMNI | LMWWEDEFEGKIPSPPTILKPR | PLWTGKQVFSLIIPRAVNLE |
| <i>Arau.het</i> | .P.S...G.VA...L... | .I.....I.....I..... | .D.D..M.T.A..... | .I.....I..... | .KQI..MT |
| <i>Cyc.rev</i> | .P.S...VL..... | .I.....I.....I..... | .D.....D..... | .N.....N..... | .QI..I |
| <i>Eph.vir</i> | .P.S...IL..... | .I.....I.....I..... | .D..V.A.A..... | .I.....I..... | .KQI..V |
| <i>Gnk.bil</i> | .P.S...VL..... | .I..K.....I.....I..... | .D.....D..... | .N.....N..... | .QI..I |
| <i>Mag.sou</i> | .P.S...VL..... | .I.....I.....I..... | .D...A..... | .N.....N..... | .KQI..I |
| <i>Nym.odo</i> | .P.C...VL..... | .I.....I.....I..... | .D...A.A.M... | .N.....N..... | .KQI..I |
| <i>Zea.mays</i> | .P.S...VL..... | .I.....I.....I..... | .Q.D...A..... | .I.....I..... | .KQI..I |
| <i>Arab.tha</i> | .P.S...VL..... | .I...A.....I.....I..... | .D..V.A.A..... | .N.....N..... | .KQI..L |

200

| | | | | | |
|------------------|---------------------|-------------------------|----------------------|-----------------------|----------------------|
| <i>Spirogyra</i> | RYSAWHPDSEKDFSPGDTQ | VRVEKGELLAGLLCKKSLGT | SGGSLVHIIWEEVGPDAARK | FLGHTQWLNVNYWLLQQGFSI | GIGDTIADASTMDTINETIA |
| <i>Arau.het</i> | .TA..N...ST.VT...S | .I.....T.T...T...T... | .I.V.....I..... | .H.....H..... | .SA..EK..... |
| <i>Cyc.rev</i> | .SE..T.FIT...C | .I...V.S.T...T...T... | .S...I.V..... | .A.E.....A.E..... | .S |
| <i>Eph.vir</i> | .NE.D..HMTL...V | .I...VIT.T...T...T.A | .S...I.V..... | .M.....M..... | .T.....Q |
| <i>Gnk.bil</i> | .SE..T.FIT...C | .I...V.S.T...T...T... | .S...I.V..... | .A.E.....A.E..... | .S |
| <i>Mag.sou</i> | .T...SEA.T.FIT... | .I.R...T...T...T...T... | .T...I.V..... | .N.....N..... | .EK.....S |
| <i>Nym.odo</i> | .SE..T.FIT...C | .I.R...S.T...TM... | .S...I.V..... | .NA.....NA..... | .EK.....S |
| <i>Zea.mays</i> | .F...SEE.K.FIT...M | .I...S.T...T...T...T... | .GS..I.V..... | .N.....N..... | .E...D..S |
| <i>Arab.tha</i> | .A.T.T.FIT... | .I.R...T...T...T...T... | .N...V.....N..T... | .S.....S..... | .EK.....S |

278

| | | | | |
|------------------|----------------------|----------------------|-----------------------|--------------------|
| <i>Spirogyra</i> | KARTEVKDLIEAACEKQLEA | QPGRTLMESFENRVNQVLNK | ARDDAGRAAQSSLSSESNVVK | AMVTAGSKGSFINISQMI |
| <i>Arau.het</i> | ...QD..E..K.SQ..S..P | ...K.....R | ...E..SS..K..... | |
| <i>Cyc.rev</i> | ...A..NQ..QL.HQ.A... | E...M..... | ...SS..K.....L. |T |
| <i>Eph.vir</i> | D..IK.QE...KYMAHK..Q | E...L.....Q..... | ...NS..K.....L. |T |
| <i>Gnk.bil</i> | ...N..Q..K..Q..A... | E...M..... | ...SS..R.....L. |T |
| <i>Mag.sou</i> | ...N..E..K..Q..... | E...M..... | ...SS..K.....L. |T |
| <i>Nym.odo</i> | ...N..E..K..Q..... | E...M..... | ...SS..R.....L. |T |
| <i>Zea.mays</i> | ...NA..E..KK.H..... | E...M..... | ...SS..N.....L. |T |
| <i>Arab.tha</i> | N...A.....RQFQG.E.DP | E.....MRDT | ...SS..K..A.T..L. |T |

Appendix 3: 282 position alignment of 19 amino acid sequences of regions D-F of RPB1 in animals, plants and fungi that was used to infer the phylogenetic trees shown in figures 14-16. A red alga, *Porphyra yezoensis*, is the outgroup. Residues identical to those in the outgroup are represented by dots while gaps in the alignment are shown as dashes.

| | | | | | |
|------------------|------------------------|----------------------|---------------------|----------------------|---|
| <i>Por.yez</i> | HVPQTHATRAEVMELMMVPR | CIVSPQGNKPMGIVQDITLL | GCMLFTYRDTFRRDVTMSL | LLHVEGWDGVI PPPAIKPE | PLWTGKQLFSLLLP-DVNLV |
| <i>Sac.cer</i> | ...SEE...LSQ.CA..L.Q. | ...S...C..... | ...C...IRKL.L.... | ...IEL.QVLNM | ...Y.W.PD.....T.....K...S...IL.VAI.NGIH.Q |
| <i>Sch.pom</i> | ...SEE...IQ.IT...K.Q. | ...S..... | ...A...VRK.SL.N.. | ...T.NAV.NI | ...M.W.PD...IL...V.L..K.V.....IL..II.KGI..I |
| <i>Mag.sou</i> | ...SFE...L...K | ...S.R..... | ...R.KI.K.... | ...IEK..F.NI | ...M.WW.DF..K.A.T.L.R |
| <i>Nym.odo</i> | ...CFE...L...K | ...S.R..... | ...R.KI.K.... | ...IEK..F.NI | ...M.WW.DF..K.A..M.R |
| <i>Zea.mays</i> | ...SFE...L...K | ...S.R..... | ...R.KI.K.... | ...LIEK..F.NI | ...M.WWQDF..K.A.T.L.R |
| <i>Cyc.rev</i> | ...SFE...L...K | ...S.R..... | ...R.KI.K.... | ...IEK..F.NI | ...M.WW.DF..K.S.T.L.R |
| <i>Gnk.bil</i> | ...SFE...L...K | ...K.S.R..... | ...R.KV.K.... | ...IEK..F.NI | ...M.WW.DF..K.S.T.L.R |
| <i>Arab.tha</i> | ...SFE...L...K | ...A.R..... | ...R.KI.K.... | ...IEK..F.NT | ...M.WW.DF..KV.A..L.R |
| <i>Eph.vir</i> | ...SFE...IL...K | ...S.R..... | ...R.KI.K.... | ...IEK..F.NI | ...M.WW.DF..KV.A..L.R |
| <i>Spirogyra</i> | ...C.FE...T...K.V. | ...S.R..... | ...R.KV.K.... | ...IEK..F.NI | ...M.WW.DFE.K.S.T.L.R |
| <i>Arau.het</i> | ...SFE...G.A..L..K | ...S.R..I..... | ...R.KV.K.... | ...IEK..F.NI | ...M.WWDDF..KM.T..L.R |
| <i>Dro.mel</i> | ...SME...ENIHIT.. | Q.IT..A..... | ...T.AVRKM.K..V. | IT.EQV.N. | ...M.FLPT..AKM.Q.C.L.R |
| <i>Art.sal</i> | ...A.SLE...LENIHIT.. | Q.IT..A.R..... | ...C.AVRKM.K..V.. | EKEQM.T. | ...M.YLPT..KL.Q..L.K |
| <i>Cra.gig</i> | ...L..SLE.K..ISN.AL... | M.IT..A.R..... | ...T.AVRKM.K..V..D. | GQX.N. | ...M.FLPR...HV.Q..L.K |
| <i>Ily.obs</i> | ...L..SLE...I.N.CA... | M.IT..A.R..... | ...T.AVRKM.K..V..T. | AQM.H. | ...M.FLPT...RM.Q.....I...VI.GR..VI |
| <i>Mus.mus</i> | ...L..SLE...IQ..A... | M.T..S.R..... | ...T.AVRK..K..V..E. | GEV.N. | ...M.FLST..KV.Q..L.R |
| <i>Cae.ele</i> | ...L..SLE...IE.IA... | QLIT..A..... | ...C.AVRMM.K..V. | IDWPFM.D. | ...M.YLPT..KV.Q..L.K |
| <i>Hel.sta</i> | ...LA.SLE...ISQ.AS.K. | M.IT..A.R..... | ...S.T.AVNKM.R... | ITK.EI.NI | ...M.YLP..A.KL.Q..L.R |

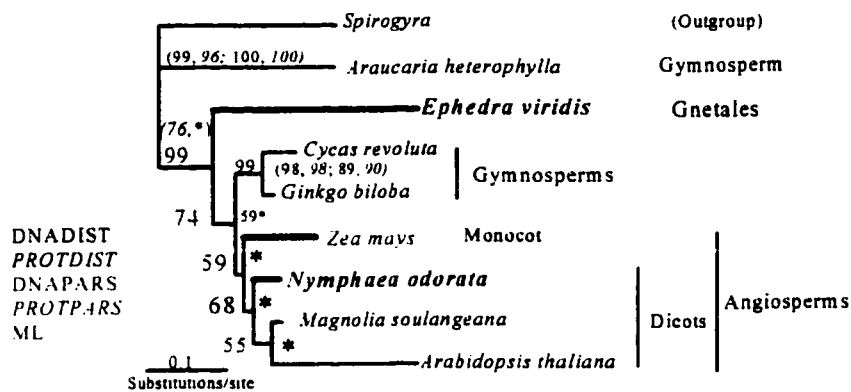
Por.yez RYCNTHPDDE--S--TDISP GDRVLIIVGGELITGIVDKR TVGSAANGLIHVTWKEKGPE RTCVLISAIQVLVNHVYIMR GQSIGIGDTIADAHTDANVR
Sac.cer .FDEGTTLLS----- K.NGM..ID.QI.F.V.E.K...SNG...VTR...Q VCAK.FGN.KV..FWLLHN .F.T.....GP.MREIT
Sch.pom .DDDKQSLSN----- T.SGM..EN..I.Y.V...K...ASQG.V.TI..... ICKGFFNG.RV..YWLLHN .F.....D.MKE.T
Mag.sou .TSAW.SEA.-TG---F.T. ...Q.R.ER..LA.TLC.K.L.TSTGS...I.E.V..D AARKFLGHT.W...YWLLQN .F.....S.MEKIN
Nym.odo .SAW.SES.-TG---F.T. ...C.R.ER..LS.TLC.K.M.TSSGS...I.E.V..D AARKFLGHT.W...YWLLQN AF.....S.MEKIN
Zea mays .FSAW.SEE.-KG---F.T. ...M.R.EK..LS.TLC.K.SL.TSSGS...I.E.V..D AARKFLGHT.W...YWLLQN .F.....S.METIN
Cyc.rev .SAW.SES.-TG---F.T. ...C.R.EK..VLS.TLC.K.L.TSSGS...I.E.V..D AARKFLGHT.W...YWLLQQ .F.....A.METIN
Gnk.bil .SAW.SES.-TG---F.T. ...C.R.EK..VLS.TLC.K.L.TSSGS...I.E.V..D AARKFLGHT.W...YWLLQQ .F.....A.METIN
Arab.tha .SAW.A.T.-TG---F.T. ...Q.R.ER..LA.TLC.K.L.TSNGS.V.I.E.V..D AARKFLGHT.W...YWLLQN .FT.....SS.MEKIN
Eph.vir .SAW.NESD-RG---HMTL ...V.R.EK..V...TLC.K.L.ASSGS...I.E.V..D AARKFLGHT.W...YWLLQQ .F.M.....T.MDTIN
Spirogyra .SAW...S.-RG---F...Q.RVEK...LA.LC.K.SL.TSSGS.V.II.E.V..D AARKFLGHT.W...YWLLQQ .F.....S.MDTIN
Arau.het .TAAW.N.S.-ST---VT. ...S.R.EK..L...TLC.K.L.TSSGS...I.E.V..D AARKFLGHT.W...YWLLQH .F.....SA.MEKIN
Dro.mel .THS...E.DEGPYKW...K.MVEH...M.LC.K.SL.TS.GS.L.ICFL.L.HD IAGRFYGN.TVI.NWLLFE.H.....PQ.YNEIQ
Art.sal .KTHS...E.DGPKYK...K.MVEH...M.LC.K.SL.TS.GS.L.ICFL.L.HD IAGRFYGN.TVI.NWLLFE.H.....PQ.YNEIQ
Cra.gig .THS...E.DGPKYK...K...ED.M.S.LC.K.L.TSSGS.A.VEM.Y.WV JAGEMYGH.T...NWLLYE.H.....PQ.YNSIQ
Ily.obs .THS...G.DSGPYKW...K...ED...S.LC.K.L.TS.GS.V.IVFL.M.F.VAGE.YGN.TV.NWLLLE.H.....PQ.YIDIQ
Mus.mus .THS...D.SGPKYK...K.VVEN...M.LC.K.SL.TS.GS.V.ISYL.M.HD I.RLFY.N.TVI.NWLLIE.H.....S.SK.YQDIQ
Cae.ete .THS...S.DSGPKYK...K.I.EH..LS...CSK...KS.GN.L.VTL.L.Y.IAANFY.H.TVI.AWL.RE.HT.....QA.YLDIQ
Hel.sta .THS...DRGPHKW...K.K.VED.K.LS.LC.K.SL.ASSGS.Q.IIHH.L.SD.A.ADFYAY.MVT.WLLVT.HT...A.....K.YSDIQ

Por.yez ATITGAQAEVTQLERRAQEG ELTLLPGKSMMEFEVEVKN VLNWARDTSGSSAQLSILKS NNIKRMVSAGSKGSFINISO IC
Sac.cer E..AE.KKK.LDVTKE..AN L..AKH.MTLR...DN.VR F..E...KA.RL.EVN.KDL .V.Q..M.....A. MS
Sch.pom R.VKE.RRQ.AECIQD..HN R.KPE..MTLR...AK.SR I..Q...NA.R.EH..KD. .V.Q..A.....MS
Mag.sou E..SK.KN..KE.IKA..K.Q.EAE..RT...NR..Q...K...DA...K..SE. .L.A..T.....MT
Nym.odo E..SK.KN..KE.IKA..K.Q.EAE..RT...NR..Q...K...DA...K..SE. .L.A..T.....MT
Zea mays D..SK.KNA.KE.IKK.H.K.Q.EAE..RT...NR..Q...K...DA...K..SE. .L.A..T.....MT
Cyc.rev E..SK.K..N..IQL.HQK.A.EAE..RT...NR..Q...K...DA...K..SE. .L.A..T.....MT
Gnk.bil E..SK.KN..K..IKA..K.A.EAE..RT...NR..Q...K...DA...K..SE. .L.A..T.....MT
Arab.tha E..SN.KTA.KD.I.QF.GK .DPE..KT.RDT..NR..Q...K...DA...K..AET .L.A..T.....MT
Eph.vir E..QD.KIK.QE.IEKYMAH K.EQE..RTL...NQ..Q...K...DA.N..R..SE. .L.A..T.....MT
Spirogyra E..AK.KT..KD.IEA.C.K.Q.EAQ..RTL...NR..Q...K...DA.RA.S..SE. .V.A..T.....MI
Arau.het E..AK.KQD.KE.IKAS..K.S.EPQ..RTL...NK..Q...R..EA...K..SE. .V.A..T.....MI
Dro.mel QA.KK.KDD.INVIQK.HNM .EPT.NTLRQT.NK..R I.D.H.KT.G.KK.TEY .L.A.VS.....VI
Art.sal T..KK.KED.IDVITK.HNN .EPT.NTLRQT.NQ..R I.D..KT.G.KN.TEY .L.A.VS.....VI
Cra.gig D..KK.KHD.IEVIEK.HND D.EPT.NTLRQT.NM..R I.D..KT.K.K.SDY .F.A.V.....VI
Ily.obs E..RK.K.D.IEVIEK.HND .EPT.NTLRQT.NQ..R I.D..KT.K.K.SEF .F.A.V.....VI
Mus.mus N..KK.KQD.IEVIEK.HNN .EPT.NTLRQT.NQ..R I.D..KT...K.SEY .F.S.VS.A...K.....VI
Cae.ete N..RK.KQD.VDVIEK.HND D.EPT.NTLRQT.NK..Q I.D..RT...K.SEF .F.S.VS...K.....VI
Hel.sta TA.KK.KSD.VEVIEK.HND .EPM.NTLRQT.NQ..R I.D..KT.L.K.SEF .F.S.V.....NK.....VI

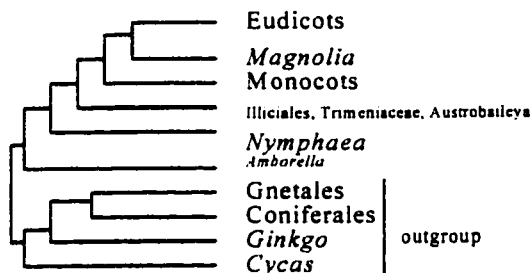
Sulf. acido L L L E G V F D K K A I G N Q P E S N L H W S I R E Y G T E Y G K W L M D N V F K M F I R F L E H R G F T M T L E D I T I P D E A Q N E I T T K I K E G Y S Q V D E Y I R K F N E G Q L E P I P G R T I
Por. yez. II E . I T . I V . . R T V . S A A N G L I . . V T W K . K . P . R T C V . I S A I Q V L V N H Y V I M . . Q S I G I G . T I A D A H T D A N V R A T . T G A Q A E . T Q L E . R A Q . . E . T L L . . K S M
Bon. ham. II Q . I C . I V . . R T V . S S A N G L I . . I T W K . F . P K I T D T . I S Q I Q V L V N H Y I L Q . R Q S I G I G . T I A D . A T M R N V I D T . Q Q A K E E . K L L V . A Q . K E . V L L . . K G M
Dro. mel. II E . I M . I L C . . S L . T S A G . L . . I C F L . L . H D I A G R F Y G . I Q T V I N N W . I . F E . H S I G I G . T I A D P Q T Y . . . Q Q A . . K A K D D . I N V . Q . A H N M E . . . T . . N . L
Art. sal. II E . I M . I L C . . T L . A X A G . I . . I I F L . L . H D I C G K F X G . I Q T V V N N W . L Y E . H S I G I G . T I A D P Q T Y . S . Q . T . . K A K E D . I D V . T . A H N N E . . . T . . N . L
Cra. gig. II M . I S . I L C . . T L . T S S G . L A . . V V F M . . . W V I A G E M Y G H I Q T L V N N W . L L E . H S I G I G . T I A D P Q T Y I D . Q D T . . K A K H D . I . V . E . A H N D D . . . T . . N . L
Ily. obs. II E . I S . I L C . . T L . T S A G . L V . . I V F L . M . F . V A G E . Y G . I Q T V V N N W . L I E . H S I G I G . T I A D Q Q T Y H . . Q E T . R K A K A D . I . V . E . A H N D E . . . T . . N . L
Mus. mus. II E . I M . I L C . . S L . T S A G . L V . . I . Y L . M . H D I T R L F Y S . I Q T V I N N W . L I E . H . I G I G . S I A D S K T Y Q D . Q N T . . K A K Q D . I . V . E . A H N N E . . . T . . N . L
Cae. ele. II E . . S . I V C S . T V . K S A G N L . . V V T L . L . Y . I A A N F Y S H I Q T V I N A W . I R E . H . I G I G . T I A D Q A T Y L D . Q N T . R K A K Q D . V D V . E . A H N D D . . . T . . N . L
Hel. sta. II K . . S . I L C . . S L . A S S G . L Q . . I I H H . L . S D A T A D F Y A Y I Q M V T N H W . L V T . H . I G I A . T I A D A K T Y S D . Q . A . . K A K . D . V . V . E . A H N D E . . . M . . N . L
Spirogyra. II E . . A . I L C . . S L . T S G G . L V . . I I W E . V . P D A A R K F L G H T Q W L V N Y W . L Q Q . . S I G I G . T I A D A S T M D T . N E T . A K A K T E . K D L . E A A C . K . . A Q L
Arau. het. II E . . T . T L C . . T L . T S G G . L I . . V I W E . V . P D A A R K F L G H T Q W L V N Y W . L Q H . . S I G I G . T I A D S A T M E K . N E T . A K A K Q D . K . L . K A S Q . K S . . . Q L
Cyc. rev. II E V . S . T L C . . T L . T S S G . L I . . V I W E . V . P D A A R K F L G H T Q W L V N Y W . L Q Q . . S I G I G . T I A D A A T M E T . N E T . S K A K A E . N Q L . Q L A H Q K A . A E M
Eph. vir. II E V I T . T L C . . T L . A S S G . L I . . V I W E . V . P D A A R K F L G H T Q W L V N Y W . L Q Q . . S I G I G . T I A D A A T M E T . N E T . S K A K A E . Q . L . E . Y M A H K . Q E L
Gnk. bil. II E V . S . T L C . . T L . T S S G . L I . . V I W E . V . P D A A R K F L G H T Q W L V N Y W . L Q Q . . S I G I G . T I A D A A T M D T . N E T . S K A K N E . K Q L . K A A Q . K A . A E M
Mag. sou. II E . . A . T L C . . T L . T S T G . L I . . V I W E . V . P D A A R K F L G H T Q W L V N Y W . L Q N . . S I G I G . T I A D A S T M E K . N E T . S K A K N E . K . L . K A A Q . K . . A E M
Nym. odo. II E . . S . T L C . . T M . T S S G . L I . . V I W E . V . P D A A R K F L G H T Q W L V N Y W . L Q N A . S I G I G . T I A D A S T M E K . N E T . S K A K N E . K . L . K A A Q . K . . A E M
Zea. may. 2 E . . S . T L C . . S L . T G S G . L I . . V I W E . V . P D A A R K F L G H T Q W L V N Y W . L Q N . . S I G I G . T I A D A S T M E T . N D T . S K A K N A . K . L . K . A H . K . . A E M
Arab. tha. 2 E . . A . T L C . . T L . T S N G . L V . . V I W E . V . P D A A R K F L G H T Q W L V N Y W . L Q N . . I G I G . T I A D S S T M E K . N E T . S N A K T A . K D L . . Q . Q G K E . D . E M
Dic. dis. II E . . A . I L C . R S L . A A N G . I I . . V V M N . H . H D T C R L F I . Q T Q T V V N H W . I N G I G . T I A D S A T M A K V . L T . S S A K N . . K . L . I . A Q N K . F . C Q K S V
Ac. cas. II E . V S . I L N . . T L . T S H K . L V . . V I W N . H . S . V C . H F L N Q . Q H V V N Y W . L . H . . S V G V G . T I A D E . T L A K . . Q T . R K A K D E . K . R Q L E A Q Q R Q M
Sac. cer. II Q I I F . . V E . . T V . S S N G G L I . . V V T . . K . P Q V C A K . F G . I Q . . V V N F W . L . N . . S T G I G . T I A D G P T M R . . . E T . A . A K K K . L D V T K E A Q A N L . T A K H . M . L
Sch. pom. II E I I Y . . V . . T V . A S Q G G L V . . T I W K . K . P . I C . G F F N G I Q R V V N Y W . L . N . . S I G I G . T I A D A D T M K . V . R T V . . A R R . . A . C . Q D A Q H N R . K . E . . M . L
Nos. loc. II V . . N . . I C . . S V . T A Q G . L I . . I I A N D . H . E I T R F I . S L Q . . L I S T Y . T L I T . S V G I G . T I S S P . T M A H . S R A . G D A K A E . K Q L . Y D S R R R S . . K L . . M N L
Vai. nec. II E I . S . I I V . S T Q G G L I . . I I A N D F . P D R V T C F F . D A Q . . M N I Y F A T I A . S I G I G . A I A D K . T M S Q V Q R S . E T A K E . . N . I . V . A Q N K . . R L . . M S M
Pla. fal. II E . . S . I I C . R T V . S S S G . L I . . V L W H . M . P D K T . D F L S A L Q . . V T N N W . . Y V . . V S C S . I A S N K V I G K V R E I L D K S K . E . S K L V E . A Q K . E . . C Q . . K S L
Mas. inv. II Q . H T . . I T N . . V . K S Q G . I I . . I L W K D Q . P M A A R D E L S R . Q L L T N A Y I L T . . . S V G T . . T L A D P D T L Q A V K A E . E G A K R N . K I H . D D A R A . R . K V Q A . . S L
Tri. vag. II N . . A . I L G . . T V A R S E G . L I . . V V . N S . N . N I A . D F L N Q T Q L I V N N W S I G . I . C V V . . F V L Q . V K H E . D D V E . K . Q A T . I . A Q D R M . . M S Y
Try. bru. II Q . . C . P I T . S I V . A A . G . L I . . V I F N . H . S D E V A R F I N G . Q R V T T F . . L N F . . S V G V Q . T V A D S D T L R Q M N D V L V K T P R N . E K I G A A A . N R T . N R K A . M . L
Gia. lam. II S V V G T E S A L I . . I L F . D . . I . P C R A F I . . C Q R V V C . M L D H . . S V G M G . M V S S E H T E R K V A E I Q T K L S K D I S . I L F K L S I H Y K I K L A . . Q S R
Sac. cer. III Q I . S . M . . S V L . D G K H . V F Y T I L . D . . P Q E A A N A . N R M A . L C A . . . G N S I G I N . V P A . D L K Q K K E E L V E I A . H K Q . . K . . T L . . K . E . . T Q . . C N E
Gia. lam. III E H . V . R L T . T F L A S S . N C I F Y F L V Q N . . P V S A A R I . L R F A . V A A . . . M N Y . . . I G I D . V M P S Q R V L G K K E V I V E K A Q . K . . D Y E S . K . A S . V
Try. bru. III C F I S . R L . . L L . G K D G L F A R L H T I A . G G . T A R V . S R I A Q F T S . Y . T N Y . . S L G . G . V A P T P . L N K Q K A A V L A R S V E V C . G L . K S A K T . R M I . L . . L . V
Mus. mus. I E . . C . L . . A H Y . S S A Y G L V . C C Y E I . . G . T S G R V L T C L A R L . T A Y . Q L Y . . L G V L V K P N . D V V R Q R I . E . A V K A A L S L P E T A S C D S K D Q R D F N M .
Dro. mel. I E . . V . L . . Q Q Y . A T T Y G L I . C M Y E L . G D V S T L . V T A F T . V T F . Q L E . . L G V K . . L V T . V . D R K R R K I . R . A V A A A L . L E D E P P H D E -
Try. bru. I E . I T . F M C . . Q L . A S N M . A P . H V Y E L . P H R T G Q . F A A F G R V L L L A . R K E . L S L A M D . M F L V E R R R C D L L R K L . D I A L D V P D E E -
Sac. cer. I A . . C . I L . . S Q Y . A S K Y G I V . S L H E V . . P . V A A K V L S V L G R I . T N Y I T A . A . . C G M D . L R L T A . G N K W R . D I L . T A A A E . T N L D K D T P A D D N N K S G I L
Sch. pom. I E . . C . I L . . S S F . A S A F G L V . S V H E L . . P D I A G R . L S V L S R L . T A Y A Q M C R M D . L R L D E Q G D . W R R Q L L E N A A . E Y V G L S T D S P I A D D E K L Q G L

Sulf. acido EESLESYILDITDKLRKVAG EIATKYLDPFNNVYIMAITG ARGSELNITQMT
Por. yez. II M..F.VEVNKV.NGA.DTS. SS.QLS.LKS..IKR.VSA. SK..FI..S.IC
Bon. ham. II M..F.TEVNKV.NGA.DKS. AS.QRS.LKS..IKR.VSA. SK..FI..S.IC
Dro. mel. II RQTF.NKVNRI.NDAHDKT. GS.K.S.TEY..LKA.VVS. SK..NI..S.VI
Art. sal. II RQTF.NQVNRI.NDA.DKT. GS.KNS.TEY..LKA.VVS. SK..NI..S.VI
Cra. gig. II RQTF.NMVNRI.NDA.DKT. SK.Q.S.SDY..FKA.VVA. SK..KI..S.VI
Ily. obs. II RQTF.NQVNRI.NDA.DKT. SK.Q.S.SE..FKA.VVA. SK..KI..S.VI
Mus. mus. II RQTF.NQVNRI.NDA.DKT. SS.Q.S.SEY..FKS.VVS. .K..KI..S.VI
Cae. ele. II RQTF.NKVNQI.NDA.DRT. SS.Q.S.SE..FKS.VVS. SK..KI..S.VI
Hel. sta. II RQTF.NQVNRI.NDA.DKT. SL.Q.S.SE..FKS.VVA. SK.NKI..S.VI
Spirogyra. II M..F.NRVNQV.N.A.DD.. RA.QSS.SES..KA.VTA. SK..FI..S..I
Arau. het. II M..F.NKVNQV.NRA.DE.. SS.Q.S.SES..KA.VTA. SK..FI..S..I
Cyc. rev. II M..F.NRVNQV.N.A.DD.. SS.Q.S.SES..LKA.VTA. SK..FI..S..I
Eph. vir. II L..F.NQVNQV.N.A.DD.. NS.QRS.SES..LKA.VTA. SK..FI..S..I
Gnk. bil. II M..F.NRVNQV.N.A.DD.. SS.QRS.SES..LKA.VTA. SK..FI..S..I
Mag. sou. II M..F.NRVNQV.N.A.DD.. SS.Q.S.SES..LKA.VTA. SK..FI..S..I
Nym. odo. II M..F.NRVNQV.N.A.DD.. SS.QRS.SES..LKA.VTA. SK..FI..S..I
Zea. may. 2 M..F.NRVNQV.N.A.DD.. SS.QNS.SES..LKA.VTA. SK..FI..S..I
Arab. tha. 2 RDTF.NRVNQV.N.A.DD.. SS.Q.S.AET..LKA.VTA. SK..FI..S..I
Dic. dis. II I.TF.QKVNQV.N.A.DT.. SS.QDS.SED..LKA.VTA. SK..FI..S..M
Aca. cas. II M..F.FV.NQI.N.A.DD.. NS.Q.S.RRS..FKA.V.A. SK..AI..S.VL
Sac. cer. II R..F.DNVRF.NEA.DK.. RL.EVN.KDL..KQ.VMA. SK..FI..A..S
Sch. pom. II R..F.AKVSRI.NQA.DN.. RS.EHS.KDS..KQ.VAA. SK..FI..S..S
Nos. loc. II Q.T..K.NLA.N.A.DIS. TR.VES.NHL.GLKQ.LKA. SK..YI..S.I.
Vai. nec. II R..F.QVNYI.N.P.DIS. AS.S.S.SFC..MRT.VLA. SK..FI..S.V.
Pla. fal. II Y..F.TRVNNE.NCA.EM.. KV.SES..ER..IFS.VAS. SK..II..S.II
Mas. inv. II V..F.AKTNS.QDALSN.. KKSLS.S.KYD..FKL.IES. SK..M..C.I.
Tri. vag. II MQ.F.TEVNNTNEILSKTY KVINAKIRGD.SLSE.LSA. SK.ADT.MS.II
Try. bru. II LQ.F.ADVNSA.N.C.EE.A KK.I.SNVRRT.SPKV.IEA. SK.TD..C.IA
Gia. lam. II NDAF.QEVISKVSGLALE KVI.DAAPHR.ALLV.INA. SK.KKF.MM.IS
Sac. cer. III .QT..AK.GGL.S.V.EEV. DVCINE..NW.APL..TC. SK..T..VS..V
Gia. lam. III Q.T..ATLNQI.SNV.ESCA Q..L.E.HFT.KPL..SLC. SK..PI..A..I
Try. bru. III KQ..ARLNTE.S.V.DEC. TA.VQT.SIH..PL..VQS. SK..A..A..M
Mus. mus. I DMKFKEEVNHYSNEIN.ACM PLGLHRQF.E..LQM.VQS. .K..TV.TM.IS
Dro. mel. I DRKYK.LLDGYTNDINSTCL PRGLITKF.S..LQL.VLS. .K..MV.TM.IS
Try. bru. I -ATAAPM.A.YAT.IQOEFV PQRMLVPF.K.HLLL.T.S. .K..N..A..IS
Sac. cer. I DAVTS.KVNAITSQVSKCV PDG.MKKF.C.SMQA..LS. .K..NV.VS.IM
Sch. pom. I DAAMKGMNGLTSSIINKCI PDGLLTKE.Y.HMQT.TVS. .K..NV.VS.IS

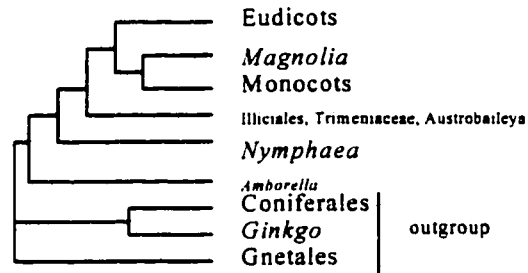
Appendix 5: Results of this study of regions D-F of RPBI in plants compared with the most recently published studies, which were based upon larger datasets and found that *Nymphaea*, a herbaceous dicot, was amongst the earliest extant angiosperms.



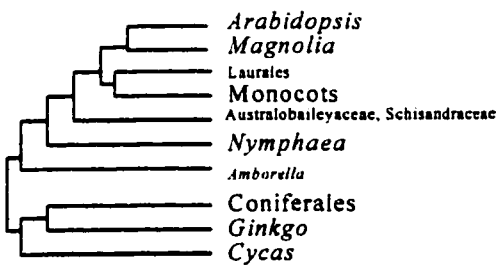
A maximum-likelihood phylogenetic tree of the 278-position amino acid sequence alignment of plant RPBI sequences, inferred using the JTT model of substitution. The γ parameter $\alpha = 1.01$, and the fraction of invariable sites = 0.41.



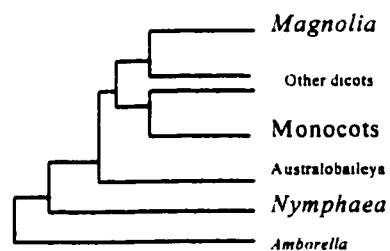
Adapted from: Qiu *et al* (1999) *Nature* 402:404-407. Parsimony (PAUP*) analysis inferred from a 8733 bp alignment of nuclear 18S rRNA, mt *atp1* & *matR*, cp *atpB* & *rbcL* from 105 taxa.



Adapted from: Soltis *et al* (1999) *Nature* 402:402-404. Parsimony (PAUP*) analysis inferred from a 4733 bp alignment of nuclear 18S rRNA and cp *atpB* & *rbcL* from 567 taxa.



Adapted from Parkinson *et al* (1999) *Current Biology* 9:1485-8. ML analysis of 6564 bp alignment of mt (SSU rDNA, *cox1*, *rps2*), cp *rbcL* and nuclear SSU rDNA from 51 taxa inferred with fastDNAm1, PUZZLE & PAUP* with 4 γ -distributed rates, accounting for rate heterogeneity



Adapted from Mathews & Donoghue (1999) *Science* 286:947-9. Parsimony (PAUP*) analysis of 2208 bp alignment of phytochrome A & C from 26 angiosperms. *PHYA* & *PHYC* used to root one another.