

**The Phylogenetic Reconstruction of the Grass Family (Poaceae)
Using *matK* Gene Sequences**

by

Hongping Liang

Dissertation submitted to the Faculty of the

Virginia Polytechnic Institute and State University

in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

in Biology

APPROVED:

Khidir W. Hilu, Chairman

Brent D. Opell

Duncan M. Porter

Saghai Maroof

Brenda W. Shirley

September, 1997

Blacksburg, Virginia

Keywords: Poaceae, Grass, DNA sequence, Phylogeny, Evolution, Taxonomy

The Phylogenetic Reconstruction of the Grass Family (Poaceae) Using *matK* Gene Sequences

by
Hongping Liang

Committee Chairman: Khidir W. Hilu
Department of Biology

Abstract

Comparative DNA sequencing of *matK*, a maturase-encoding gene located within the intron of the chloroplast *trnK* gene, was evaluated for phylogenetic utility above the family level and within the grass family (Poaceae). There are three major objectives in the research. The first one is to study the utility of the *matK* gene in plant evolution. The second objective is to characterize the *matK* gene in the grass family. The last major goal is to address the phylogenetic questions in the Poaceae using the *matK* sequences from representatives of different grass groups.

In order to study the potential application of *matK* to plant systematics above the family level, eleven complete sequences from GenBank representing seed plants and liverworts and nine partial sequences generated for genera representing the monocot families Poaceae, Joinvilleaceae, Cyperaceae, and Smilacaceae were analyzed. The study underscored the following useful properties of the *matK* gene for phylogenetic reconstruction: reasonable size (1500 bp), high rate of substitution, large proportion of variation at the first and the second codon positions, low transition-transversion ratio, and the presence of mutationally-conserved sectors. The use of different sectors of the gene and the cumulative inclusion of informative sites showed that the 3' region was the most useful in resolving phylogeny, and that the topology and robustness of the tree reached a plateau after the inclusion of 100 informative sites. The presence of a relatively conserved 3' region and the less conserved 5' region provides two sets of characters that can be used at different taxonomic levels from the tribal to the division levels. It also has demonstrated the potential of partial sequencing in resolving systematic relationships from the tribe to the division level.

The *matK* gene in the Poaceae was characterized with complete sequences from 11 grass genera, representing 7 subfamilies and 11 tribes, and one outgroup (*Joinvillea plicata*, Joinvilleaceae). The alignment of 1632 base pairs from 14 species yielded a data set of 601 (36.8)% variable sites and 246 (15.1%) informative sites. The variations at nucleic and amino acid levels evenly distributed throughout the entire gene, and the 5' region appears to have more variation than the 3' region. The changes at the third codon position are very low as compared to the total of the first and second positions. This has led to a similar variation pattern at nucleic and at amino acid levels. The average tr/tv ratio

generated from 14 entire *matK* sequences is 1.29. It is intriguing to find that the tr/tv ratios were regionally related. RASA analysis of the alignment data indicated a relatively high phylogenetic signal in the data set of 14 taxa. In the two half analyses, while the tRASA of the 5' half of the *matK* gene (0.43) is not significant, the 3' of the *matK* gene showed a significant phylogenetic signal. Among the 5 sections of the 14 entire *matK* sequences, only the fourth sector contains a statistically significant phylogenetic signal. These results indicate that *matK* is a phylogenetically valuable gene and that the 3' region of the *matK* gene contains strong phylogenetic information. A single most parsimonious tree was obtained from the 246 informative sites of the 14 entire *matK* sequences. Seven major groups were well resolved on the most parsimonious tree, corresponding to the seven commonly recognized subfamilies: Aruninoideae, Bambusoideae, Centothecoideae, Chloridoideae, Panicoideae, Pooideae and Oryzoideae.

Approximately 960 base pairs of the *matK* gene were sequenced from grass species representing 48 genera, 21 tribes, and seven subfamilies to reconstruct a phylogeny for the Poaceae. *Joinvillea plicata* (Joinvilleaceae) was used as an outgroup species. The aligned sequences showed that 495 nucleotides (51%) were variable and 390 (36%) were phylogenetically informative. RASA indicated that very significant phylogenetic signals exist in this data set. The cumulative addition of informative sites starting at the internal end of the sequences revealed that at 300 sites, tree topology and bootstrap values matched those of the consensus tree based on the entire sequence. Parsimony analyses using PAUP resulted in six most parsimonious trees and a strict consensus tree showing major lineages supported by high bootstrap values. These lineages corresponded to six subfamilies: Bambusoideae, Oryzoideae, Pooideae, Chloridoideae, Panicoideae, and Arundinoideae. The Bambusoideae, including woody and herbaceous taxa, diverged as the most basal lineage, and the monophyletic oryzoid species formed a sister group. The Chloridoideae, Panicoideae, Arundinoideae, and the centothecoid *Zeugitis* (PACC group) emerged as a monophyletic assemblage with 95% bootstrap support. The Aristideae branched off as a monophyletic line basal to the chloridoid clade. Stipeae appeared as a sister taxon to the Pooideae. The *matK*-based phylogeny did not reveal a major dichotomy in the family. The *matK* gene has provided sequence information sufficient for good resolution of the major grass lineages.

Acknowledgments

The author wishes to recognize several individuals for the contributions and help they have provided toward the completion of this dissertation.

First and foremost, I would like to thank Dr. Khidir W. Hilu for his encouragement, guidance and patience in his role as committee chairman of this study throughout the past years.

Thanks to Dr. Brent D. Opell, Dr. Duncan M. Porter, Dr. Saghai Maroof, and Dr. Brenda W. Shirley for serving on my committee, and for their advice and assistance. Special thanks to Dr. Opell for his help on statistical analyses.

I would like to thank Department of Biology, Graduate Student Assembly, and Sigma Xi for their support of the research.

Thanks to Nigel P. Barker, Christopher Campbell, Lynn Clark, Gary P. Fleming, Thomas Wieboldt, and Weiping Zhang for supplying DNA or plant samples. Seed material for some accessions were kindly provided by the U. S. Department of Agriculture Southern Regional and the ARS Plant Introduction Stations.

Special thanks is extended to William Speer and Weiping Zhang for the friendship we have been sharing in the Lab.

Finally, I thank my parents for their support and encouragement. Special thanks to my wife, Xiaoyun Li and my son, Jingyuan Liang

Table of Contents

Cover	i
Abstract	ii
Acknowledgments	iii
Table of Contents	iv
List of Figures	vi
List of Tables	vii
1 Introduction	1
1.1 DNA Sequence in Plant Systematics	1
1.2 The <i>matK</i> Gene	2
1.2.1 History of the <i>matK</i> Gene	2
1.2.2 Function of the <i>matK</i> Gene: Maturase MatK	3
1.2.3 Application of the <i>matK</i> Gene to Plant Systematics	3
1.3 Grass Family (Poaceae)	4
1.3.1 Size and Importance	4
1.3.2 Poaceae History	4
1.3.3 Major Subfamilies	5
1.3.4 Difficulties in Grass Systematics	7
1.4 Previous Studies on Poaceae Using DNA Sequences	8
1.5 Objectives of the Research	12
1.6 Literature Cited	13
2 The <i>matK</i> Gene: Sequence Variation and Application in Plant Systematics	20
2.1 Introduction	20
2.2 Materials and Methods	21
2.2.1 Materials	21
2.2.2 DNA Isolation, Amplification, and Sequencing	23
2.2.3 Data Analysis	24
2.3 Results and Discussions	24

	2.3.1 Transition/Transversions in <i>matK</i>	28
	2.3.2 Phylogenetic Analyses Based on the <i>matK</i> Gene	30
	2.4 Literature Cited	39
3	Preliminary Application of the <i>matK</i> Gene Sequences to Grass Systematics	44
	3.1 Introduction	44
	3.1.1 DNA Sequencing in the Poaceae	45
	3.1.2 Characteristics of the <i>matK</i> Gene	46
	3.1.3 Contrasting <i>matK</i> , <i>rbcL</i> , <i>psbA</i> , and <i>rps4</i> Genes in Grasses	47
	3.2 Materials and Methods	52
	3.2.1 DNA Extraction and Amplification	52
	3.2.2 DNA Sequencing	52
	3.2.3 Analysis of Sequence Data	53
	3.3 Results and Discussion	54
	3.3.1 Sequence Comparison	54
	3.3.2 Phylogenetic Analysis	59
	3.4 Literature Cited	65
4	Characterization of the <i>matK</i> Gene in the Poaceae	71
	4.1 Introduction	71
	4.2 Materials and Methods	73
	4.2.1 Plant Material and DNA Sequence Methods	73
	4.2.2 Data analysis	73
	4.3 Results	75
	4.3.1 Sequence Variation	75
	4.3.2 The tr/tv Ratio	77
	4.3.3 RASA Test	79
	4.3.4 Phylogenetic Analysis	80
	4.4 Discussions	83
	4.4.1 The tr/tv Ratio	83
	4.4.2 Phylogenetic Signal in <i>matK</i>	86

	4.4.3 Grass Phylogeny	87
	4.5 Literature Cited	88
5	Phylogenetic Construction of Poaceae Based on <i>matK</i> Sequences	91
	5.1 Introduction	91
	5.2 Materials and Methods	93
	5.2.1 Plant Material and Nucleic Acid Methods	93
	5.2.2 Data Analysis	95
	5.3 Results	97
	5.3.1 Sequence Comparison	97
	5.3.2 Phylogenetic Analysis	98
	5.4 Discussion	100
	5.4.1 Grass Phylogeny	101
	5.5 Literature Cited	110
6	Conclusions and Suggestions	116
	6.1 Application of <i>matK</i> above the Family Level	116
	6.2 Preliminary Application of <i>matK</i> to Poaceae	117
	6.3 Characterization of <i>matK</i> in Poaceae	117
	6.4 Phylogeny of the Grass Family	118
	Vita	121

List of Figures

		page
Fig. 2.1	Relative position of the PCR amplification primers	23
Fig. 2.2	Comparative sequence variation among taxa representing different taxonomic hierarchies using the GenBank sequences of the <i>matK</i> coding region	26
Fig. 2.3	Five cladograms based on sequence data from different part of the <i>matK</i> coding region	32
Fig. 2.4	Cladograms based on the incremental addition of informative sites from the 3' end of the <i>matK</i> gene	34
Fig. 2.5	The single most-parsimonious tree rooted with <i>Marchantia</i> for all 20 taxa based on 306 bp from the 3' <i>matK</i> coding region	37
Fig. 3.1	Relative position of the PCR amplification primers	47
Fig. 3.2	Variability within <i>matK</i> and other chloroplast genes between <i>Oryza Sativa</i> L. and <i>Hordeum Vulgare</i> L.	50
Fig. 3.3	Variability of the <i>matK</i> gene in the grass family	56
Fig. 3.4	Dendrogram produced by the Neighbor-Joining method with Jukes-Cantor distance	60
Fig. 3.5	The strict consensus tree derived from <i>matK</i> sequence analysis for Poaceae	61
Fig. 4.1	Variability of the <i>matK</i> gene in the grass family at nucleotide and amino acid levels	76
Fig. 4.2	The tr/tv ratio, A+T content and variable sites at different sections of the entire <i>matK</i> gene	77
Fig. 4.3	Mapping of the tr/tv ratios on the most parsimonious tree	78
Fig. 4.4	The Power and Effect Tests starting from both directions of the entire <i>matK</i> gene data set	80

Fig. 4.5	The single most parsimonious tree generated from <i>matK</i> sequence analysis for 13 grass species and the outgroup <i>Joinvillea</i>	81
Fig. 4.6	Consensus trees generated from the two halves of the 13 grass species and the outgroup <i>Joinvillea</i>	82
Fig. 4.7	The strict consensus trees generated from each of the five sectors of the entire <i>matK</i> gene	84
Fig. 5.1	A tRASA Power Plot based on <i>matK</i> sequences from the 49 taxa showing the amount of phylogenetic signal (tRASA) with the accumulative addition of informative sites	98
Fig. 5.2	The strict consensus tree of six most parsimonious trees derived from <i>matK</i> sequence analysis for 48 grass species and the outgroup <i>Joinvillea</i>	99

List of Tables

	Page
TABLE 1.1 Tribes and genera in the previous studies with DNA sequence	11
TABLE 2.1 Twenty taxa, their families, and the sequence length used in this study	22
TABLE 2.2 Transition/transversion ratios of the 11 taxa from Genbank	29
TABLE 3.1 Comparison of DNA sequences of four chloroplast genes from rice and barley for nucleotide variation, length, G+C content and transition/transversion ratios (tr/tv)	48
TABLE 3.2 Eighteen taxa and their respective tribes and subfamilies used in the sequence analysis	51
TABLE 3.3 The ratios of transition/transversion among six subfamilies and the outgroup	56
TABLE 4.1 The fourteen entire sequences of <i>matK</i>	74
TABLE 4.2 RASA test of the 5 Sections of the 14 entire <i>matK</i> genes	79
TABLE 5.1 Summary of the phylogenetic studies of Poaceae using DNA sequences	92
TABLE 5.2 Forty-nine taxa and their respective tribes and subfamilies used in sequence analysis	94
TABLE 5.3 Structure and location of primers used in sequencing	96