

Type I MADS-box genes have experienced faster birth-and-death evolution than type II MADS-box genes in angiosperms

Jongmin Nam^{†‡}, Joonyul Kim^{§¶}, Shinyoung Lee[§], Gynheung An[§], Hong Ma[†], and Masatoshi Nei[†]

[†]Institute of Molecular Evolutionary Genetics and Department of Biology, Pennsylvania State University, University Park, PA 16802; and [§]National Research Laboratory of Plant Functional Genomics, Division of Molecular and Life Sciences, Pohang University of Science and Technology, Pohang 790-784, Korea

Contributed by Masatoshi Nei, December 22, 2003

Plant MADS-box genes form a large gene family for transcription factors and are involved in various aspects of developmental processes, including flower development. They are known to be subject to birth-and-death evolution, but the detailed features of this mode of evolution remain unclear. To have a deeper insight into the evolutionary pattern of this gene family, we enumerated all available functional and nonfunctional (pseudogene) MADS-box genes from the *Arabidopsis* and rice genomes. Plant MADS-box genes can be classified into types I and II genes on the basis of phylogenetic analysis. Conducting extensive homology search and phylogenetic analysis, we found 64 presumed functional and 37 nonfunctional type I genes and 43 presumed functional and 4 nonfunctional type II genes in *Arabidopsis*. We also found 24 presumed functional and 6 nonfunctional type I genes and 47 presumed functional and 1 nonfunctional type II genes in rice. Our phylogenetic analysis indicated there were at least about four to eight type I genes and ≈ 15 –20 type II genes in the most recent common ancestor of *Arabidopsis* and rice. It has also been suggested that type I genes have experienced a higher rate of birth-and-death evolution than type II genes in angiosperms. Furthermore, the higher rate of birth-and-death evolution in type I genes appeared partly due to a higher frequency of segmental gene duplication and weaker purifying selection in type I than in type II genes.

Morphological/physiological evolution of organisms has been driven mainly by the evolution of genetic toolkits for developmental/physiological processes such as transcription factors and signaling pathways (1). A large proportion of genetic toolkits are highly conserved even between distantly related organisms. In flowering plants (angiosperms), MADS-box genes are among such toolkits that control various aspects of developmental processes. MADS-box genes are defined by the highly conserved 180-bp-long motif called the MADS-box and are found in animals, fungi, and plants (2). The protein region encoded by the MADS-box is called the MADS-domain (or M-domain) and is part of the DNA-binding domain. It has been proposed that there are at least two evolutionary lineages (types I and II) of MADS-box genes in animals, fungi, and plants (3) (Fig. 1).

There are ≈ 100 MADS-box genes in *Arabidopsis thaliana* (hereafter called *Arabidopsis*) and >70 MADS-box genes in *Oryza sativa* (hereafter called rice). There are ≈ 40 clearly identifiable type II MADS-box genes in each of *Arabidopsis* (4, 5) and rice (6). Most of the plant type II genes contain three additional plant-specific domains: intervening (I) domain (≈ 30 codons), keratin-like coiled-coil (K) domain (≈ 70 codons), and C-terminal (C) domain (variable length) (7) (Fig. 1). These genes are called MIKC^c-type genes. The MIKC^c-type genes can further be divided into two types based on the intron–exon structure: MIKC^c- and MIKC^{*}-type genes (8). The MIKC^c-type genes have been identified in most major evolutionary lineages of green plants such as angiosperms, gymnosperms, ferns, and mosses (9). The MIKC^{*}-type genes were originally found in mosses and

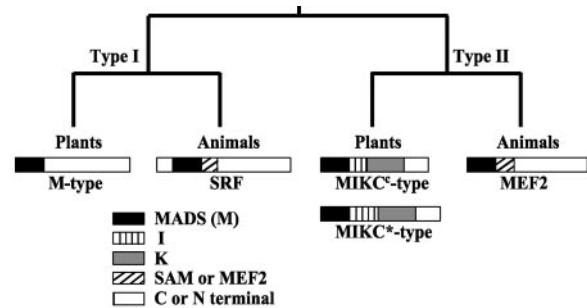


Fig. 1. Domain structures of types I and II MADS-box genes in plants and animals. Adapted from ref. 3 on the structures of types I and II genes and from ref. 8 on the structures of MIKC^c-type and MIKC^{*}-type genes.

clubmosses (8, 10), but these genes are also present in *Arabidopsis* (5). By contrast, the type I MADS-box genes in plants do not encode the K-domain and are sometimes called M-type genes (5).

It has been shown that at least 11 classes of MADS-box genes are shared between *Arabidopsis* and rice/maize (9, 11). All of them are MIKC^c-type genes, and their expression patterns have been studied intensively in eudicots. Several classes of MIKC^c-type genes, called floral MADS-box genes, are concerned with the development of floral components (organs) such as petals, sepals, stamens, and carpels, as well as regulation of flowering time (12, 13). Other classes of MIKC^c-type genes play diverse roles during vegetative growth (14–16) and fruit development (17). Some of the floral MADS-box genes in monocots have functions equivalent to those of their orthologs in eudicots (18, 19), suggesting an ancient origin of the machinery of flower development. There are also a few other genes that are not shared (lineage-specific) between *Arabidopsis* and rice (9). The functions of MIKC^{*}-type and M-type (or type I) genes are poorly understood.

MADS-box genes are important regulators of development of angiosperms (and probably nonflowering plants as well), and therefore the study of evolution of MADS-box genes is expected to give important clues for understanding the morphological evolution of plants. In our previous paper (20), we indicated that the MADS-box gene family has been subject to the model of birth-and-death evolution, in which new genes are generated by gene duplication, and some duplicate genes stay in the genome as differentiated genes, whereas others are inactivated into pseudogenes or deleted from the genome (21, 22). Although it

Abbreviations: M-domain, MADS-domain; MRCA, most recent common ancestor.

[†]To whom correspondence should be addressed. E-mail: jyn101@psu.edu.

[¶]Present address: Michigan State University—Department of Energy Plant Research Laboratory and Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, MI 48824.

© 2004 by The National Academy of Sciences of the USA

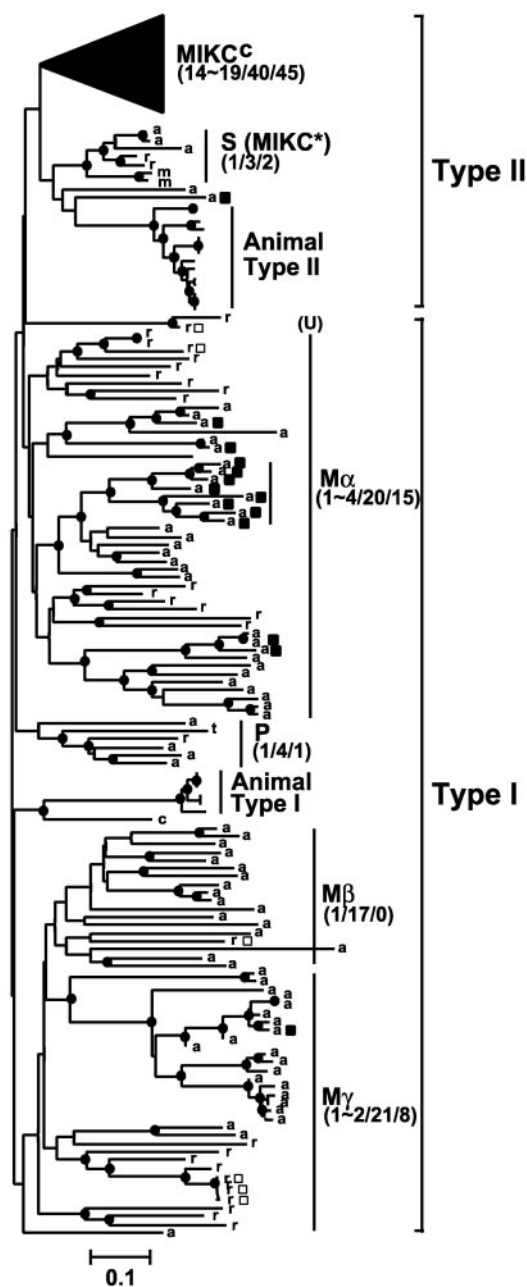


Fig. 3. Phylogenetic tree of 223 M-domain sequences from *Arabidopsis*, rice, mosses, clubmosses, and animals. This tree was constructed by the neighbor-joining method with *p*-distance and the pairwise deletion option (29) of ≈ 55 aa. *p*-distance is known to be more efficient in obtaining the correct topology when the sequence length is short (42). The genes from *Arabidopsis* and rice are labeled with "a" and "r," respectively. The reference sequences from the moss *P. patens* and the clubmoss *L. annotinum* are labeled with "m" and "c," respectively. Genes labeled with black squares (■) are pseudogenes from *Arabidopsis*, and those with open squares (□) are pseudogenes from rice. Interior branches with bootstrap values (500 bootstraps) >50% are indicated by black dots (●). The portion of the tree corresponding to the MIKCc-type genes is compressed, because it is essentially the same as that in Fig. 2. The numbers in parentheses below each class name are the numbers of ancestral MADS-box genes, MADS-box genes in *Arabidopsis*, and MADS-box genes in rice, in this order.

diverged from type II genes before the animal/plant split, although the bootstrap support is weak. These genes correspond to the type I genes proposed by Alvarez-Buylla *et al.* (3). Fig. 3

also suggests that at least one type I gene existed in the MRCA of animals and plants. This observation is consistent with that of other researchers (3, 9, 23).

Most of the shared classes observed in Fig. 2 remain unchanged in the original tree of Fig. 3 (data not shown). We also identified another shared clade (P) supported by a bootstrap value of 79%. Class P genes from *Arabidopsis* were previously classified as type I by Alvarez-Buylla *et al.* (3). However, we also observed that MIKCc*-type genes from mosses and class S and class P genes from *Arabidopsis* and rice formed a clade, when we used different sets of genes (data not shown). It is therefore possible that class P genes are also closely related to the MIKCc*-type (type II) genes from mosses as proposed by other researchers (4, 5, 23), although they are not orthologous to the latter genes. On the basis of the tree shown in Fig. 3, the remaining type I genes can further be subdivided into classes M α , M β , and M γ , in agreement with Parenicova *et al.*'s (4) classification, although bootstrap supports of these classes are very low, and class M γ genes are not monophyletic. Although our classification of type I genes is very crude, it suggests that at least about four to eight ancestral type I genes existed in the MRCA of *Arabidopsis* and rice. The numbers of functional genes and ancestral genes estimated in this way for each type of MADS-box genes in *Arabidopsis* and rice are shown in Fig. 3 (see numbers in parentheses).

Our study of ancestral MADS-box genes therefore leads to the hypothesis that there were at least $\approx 15\text{--}20$ type II genes and at least about four to eight type I genes in the MRCA of *Arabidopsis* and rice. Because there are 43 type II genes and 64 type I genes in *Arabidopsis*, the results of the present study suggest that type I genes have experienced a higher birth rate than type II genes in the *Arabidopsis* lineage. A similar pattern was also observed in rice, although it is preliminary. In addition, this pattern is quite general across most gene classes except class FLC in *Arabidopsis* and class AGL12 in rice (see numbers in parentheses in Figs. 2 and 3). One may argue that if we use more stringent criteria for estimating the number of ancestral type I genes, the number may change, and therefore the rate of gene birth would change. However, this does not affect our conclusion that type I genes have experienced a higher birth rate than type II genes. This is because many type I genes in each of classes M α , M β , and M γ from either *Arabidopsis* or rice appear to be monophyletic, suggesting that they were duplicated after the *Arabidopsis* and rice split.

Classification of MADS-Box Pseudogenes. Existence of pseudogenes means that functional genes die sometimes in the evolutionary process. To examine whether there are differences in the death rate among different types of MADS-box genes, we classified pseudogenes on the basis of sequence similarity to functional MADS-box genes. In *Arabidopsis*, four pseudogenes were most similar to the type II genes (see Table 1, which is published as supporting information on the PNAS web site), and none of these pseudogenes had the MADS-box. The remaining 37 pseudogenes were most similar to the type I genes. Fourteen of these 37 pseudogenes had the MADS-box. In the case of rice, only one of the seven pseudogenes belonged to type II, and the remainder were type I genes. These results show that the proportion of pseudogenes is significantly different between types II and I genes in both *Arabidopsis* and rice. When we applied the same criterion of pseudogenes as that of rice pseudogenes (existence of stop codons in the MADS-box), we detected nine type I pseudogenes and no type II pseudogenes in *Arabidopsis*. Our homology search and phylogenetic analysis also showed that several pseudogenes belonging to class M α are monophyletic (see Fig. 3), suggesting that the number of pseudogenes has increased recently in this lineage. Even if we exclude such lineage-specific pseudogenes, the difference in the proportion of

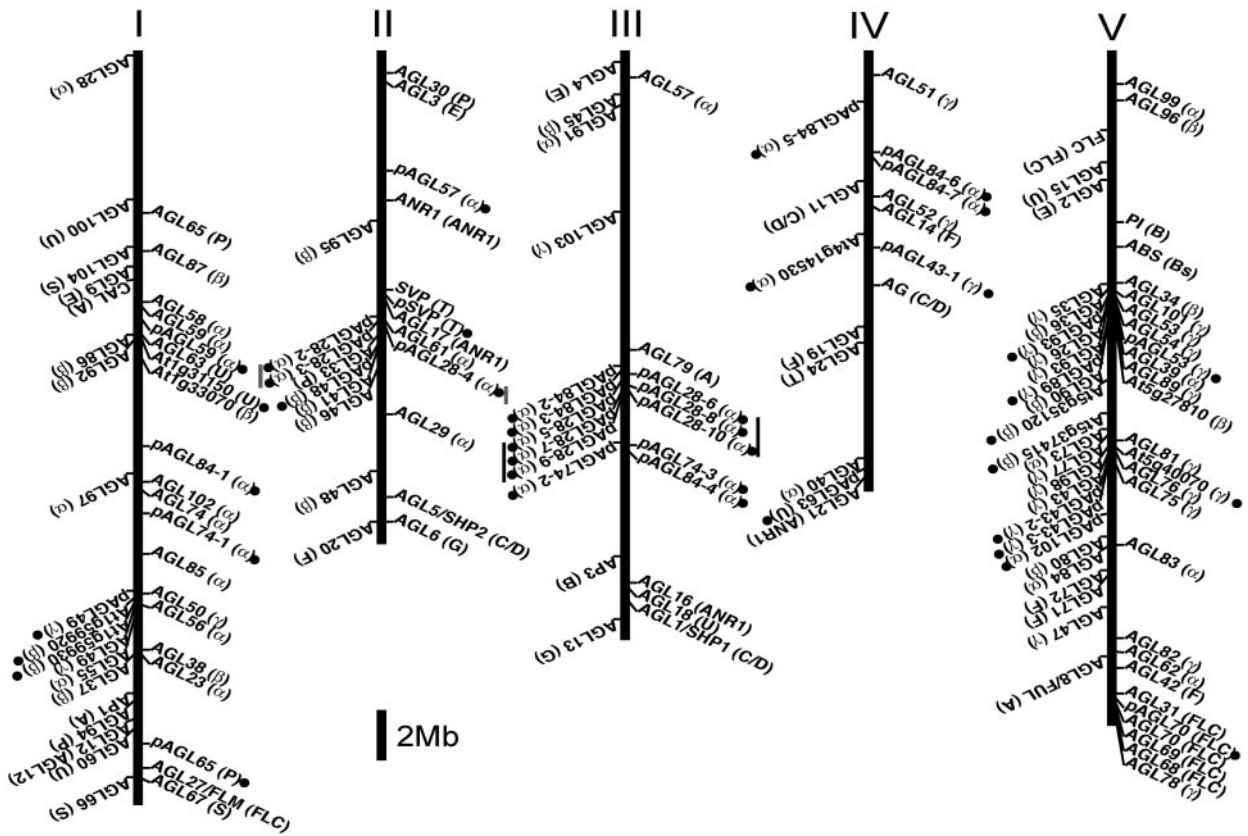


Fig. 4. Genomic organization of MADS-box genes in *Arabidopsis*. Genes with black dots (●) are pseudogenes. For seven annotated pseudogenes, we used their gene codes from GenBank. Unannotated pseudogenes are indicated by “p” in front of the name of the functional gene that is most similar to the pseudogene. For example, a pseudogene that is most similar to “SVP” is designated as “pSVP.” The class name of each gene is given in parentheses at the end of the gene name. Of these class names, “(α),” “(β),” “(γ),” and “(U)” refer to “M_α,” “M_β,” “M_γ,” and “Unassigned,” respectively. I, II, III, IV, and V represent chromosome numbers. The scale bar below chromosome II is for 2 megabase pairs (Mb).

pseudogenes between types I and II genes is still substantial. Although type I genes are expected to include more pseudogenes than type II genes because of their higher birth rate, this factor alone does not explain the difference in pseudogenes between types I and II genes. Therefore, type I genes should have had a higher death rate than type II genes.

It is not easy to have an unambiguous definition of pseudogenes, because even a fragmentary gene can be functional (35, 36), and young pseudogenes may not be distinguishable from functional genes. Therefore, different criteria for pseudogenes may change our conclusion about the death rates of types I and II genes. As mentioned above, however, our conclusions about the death rates based on two different criteria in *Arabidopsis* are essentially the same. Note also that our searches for pseudogenes are apparently biased for pseudogenes similar to more conserved functional genes (type II genes in this study) than for less conserved functional genes. Therefore, our conclusion about the difference in death rate between types I and II genes is conservative.

Genomic Organization of MADS-Box Genes in *Arabidopsis*. The genomic locations of all MADS-box genes in *Arabidopsis* are shown in Fig. 4. In general, MADS-box genes are scattered all over the chromosomes. However, we also observed a number of clusters of closely located MADS-box genes in *Arabidopsis*. Most of these genes belonged to type I genes, and in general the genes in each cluster are evolutionarily closely related. These closely related MADS-box genes were probably generated by recent segmental duplication. The genomic locations of pseudogenes

are also shown in Fig. 4. Most pseudogenes are closely located to each other as well as to their closely related functional MADS-box gene, although there are several exceptions. We also found a genomic cluster of type I pseudogenes without any functional MADS-box genes (but there are other genes) on chromosome 3 (genes with vertical bars in Fig. 4). The genomic locations and the phylogenetic tree of M-domain sequences (Fig. 3; see also Fig. 5, which is published as supporting information on the PNAS web site) suggest that this gene cluster was formed by segmental duplication of an ancestral pseudogene cluster, which was in turn duplicated from another pseudogene cluster on chromosome 2 (genes with gray bars in Fig. 4).

Rice MADS-box genes are also scattered all over the chromosomes, and more clusters of type I genes were found than those of type II genes (J.N., unpublished work).

Discussion

We have seen that type I genes have experienced faster birth-and-death evolution than type II genes in the *Arabidopsis* and rice lineages. The higher birth rate of type I genes is apparently caused by a higher rate of gene duplication, because duplicate genes generally do not cause harmful effects. In fact, type I genes are associated with a higher frequency of segmental duplications than type II genes in *Arabidopsis* (see Fig. 4). (We do not think the genome duplication is responsible for the different birth rates of types I and II genes, because in this case the birth rate should be the same for all genes. Therefore, we will not discuss this factor.) By contrast, the death of functional genes may have harmful effects, and therefore the death rate may be influenced

by functional requirements of duplicate genes as well as genomic events and fixation by genetic drift. Our estimates of the numbers of nonsynonymous nucleotide substitutions per nonsynonymous site (d_N) and synonymous nucleotide substitutions per synonymous site (d_S) suggested that type I genes have been under weaker purifying selection than type II genes (see Fig. 6, which is published as supporting information on the PNAS web site). This observation may explain why type I genes have experienced a higher death rate than type II genes, because the death of type I genes could be less harmful than that of type II genes. It is possible that, after duplication, type II genes became functionally differentiated in a relatively short time and therefore have been maintained as functional genes in the genome. This might be related to the extensive morphological diversification of angiosperms.

Although type I genes are apparently under weaker purifying selections than type II genes, they still might have played some important roles, because most of the recently duplicated type I gene pairs show significantly lower d_N and d_S (Fig. 6). Recently, it has been proposed that the expression of a type I MADS-box gene, *PHERESI*, in *Arabidopsis* is associated with seed abortion in a certain mutant background (37). However, the functions of type I genes are not well understood. If there are functionally redundant duplicate genes, it would be difficult to study their functions by mutagenesis experiments. Moreover, if type I genes are involved in a short period of developmental processes, it may also be difficult to study their functions. At the present time, gaining insights into the functional constraints of type I genes by evolutionary analysis may be of some help for future experimentation. Our study suggests that type I genes may be more variable among different angiosperm species than type II genes

because of faster birth-and-death evolution than that of type II genes. In addition, type I genes are generally less conserved than type II genes.

There are a substantial number of type II duplicate genes, although the birth rate of type II genes is lower than that of type I genes. Therefore, some extent of functional redundancy or differentiation is expected to be observed among highly similar type II genes. For example, three class E genes (*AGL2/4/9* or *SEP1/2/3*) in *Arabidopsis* are known to be functionally redundant, because single gene mutations showed only subtle phenotypic changes, whereas triple mutants showed significant phenotypic changes in flowers of *Arabidopsis* (38). Nevertheless, our d_N and d_S analysis suggests that these genes are generally subject to strong purifying selection (see File 3, which is published as supporting information on the PNAS web site). Therefore, more careful study of single gene mutations may reveal some unrecognized phenotypic effects in plants. Moreover, there is substantial conservation or differentiation in gene expressions (4, 5) and in protein coding region (39–41) among paralogous MADS-box genes. By combining experimental studies with evolutionary analyses, we may be able to have a better insight into gene functions.

We thank Hongzhi Kong, Yoshi Niimura, Nikos Nikolaidis, Li Hao, Jim Leebens-Mack, Claude dePamphilis, Kerstin Kaufmann, Mitsuyasu Hasebe, Guenter Theissen, Doug Soltis, Mike Purugganan, and Lucia Colombo for useful comments. This work was supported, in part, by National Institutes of Health Grant GM20293 (to M.N.) and grants from the Crop Functional Genomic Center, the 21st Century Frontier Program, Korea (CG1111), and the Biogreen 21 Program, Rural Development Administration, Korea (to G.A.). J.N. was partially supported by a scholarship from the Rotary Foundation.

- Carroll, S. B., Grenier, J. K. & Weatherbee, S. D. (2001) *From DNA to Diversity* (Blackwell Science, Malden, MA).
- Shore, P. & Sharrocks, A. D. (1995) *Eur. J. Biochem.* **229**, 1–13.
- Alvarez-Buylla, E. R., Pelaz, S., Liljegren, S. J., Gold, S. E., Burgeff, C., Ditta, G. S., Ribas de Pouplana, L., Martinez-Castilla, L. & Yanofsky, M. F. (2000) *Proc. Natl. Acad. Sci. USA* **97**, 5328–5333.
- Parenicova, L., De Folter, S., Kieffer, M., Horner, D. S., Favalli, C., Busscher, J., Cook, H. E., Ingram, R. M., Kater, M. M., Davies, B., et al. (2003) *Plant Cell* **15**, 1538–1551.
- Kofuji, R., Sumikawa, N., Yamasaki, M., Kondo, K., Ueda, K., Ito, M. & Hasebe, M. (2003) *Mol. Biol. Evol.* **20**, 1963–1977.
- Lee, S., Kim, J., Son, J. S., Nam, J., Jeong, D. H., Jang, S., Lee, J., Lee, D. Y., Kang, H. G. & An, G. (2003) *Plant Cell Physiol.* **44**, 1403–1411.
- Ma, H., Yanofsky, M. F. & Meyerowitz, E. M. (1991) *Genes* **5**, 484–495.
- Henschel, K., Kofuji, R., Hasebe, M., Saedler, H., Munster, T. & Theissen, G. (2002) *Mol. Biol. Evol.* **19**, 801–814.
- Becker, A. & Theissen, G. (2003) *Mol. Phylogenet. Evol.* **29**, 464–489.
- Svensson, M. E., Johannesson, H. & Engstrom, P. (2000) *Gene* **253**, 31–43.
- Munster, T., Delcu, W., Wingen, L. U., Ouzunova, M., Cacharron, J., Faigl, W., Werth, S., Kim, J. T. T., Saedler, H. & Theissen, G. (2002) *Maydica* **47**, 287–301.
- Weigel, D. & Meyerowitz, E. M. (1994) *Cell* **78**, 203–209.
- Theissen, G. (2001) *Curr. Opin. Plant Biol.* **4**, 75–85.
- Zhang, H. & Forde, B. G. (1998) *Science* **279**, 407–409.
- Michaels, S. D. & Amasino, R. M. (1999) *Plant Cell* **11**, 949–956.
- Alvarez-Buylla, E. R., Liljegren, S. J., Pelaz, S., Gold, S. E., Burgeff, C., Ditta, G. S., Vergara-Silva, F. & Yanofsky, M. F. (2000) *Plant J.* **24**, 457–466.
- Liljegren, S. J., Ditta, G. S., Eshed, Y., Savidge, B., Bowman, J. L. & Yanofsky, M. F. (2000) *Nature* **404**, 766–770.
- Ma, H. & dePamphilis, C. (2000) *Cell* **101**, 5–8.
- Kang, H. G., Jeon, J. S., Lee, S. & An, G. (1998) *Plant Mol. Biol.* **38**, 1021–1029.
- Nam, J., dePamphilis, C. W., Ma, H. & Nei, M. (2003) *Mol. Biol. Evol.* **20**, 1435–1447.
- Nei, M. (1969) *Nature* **221**, 40–42.
- Nei, M., Gu, X. & Sitnikova, T. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 7799–7806.
- De Bodt, S., Raes, J., Florquin, K., Rombauts, S., Rouze, P., Theissen, G. & Van de Peer, Y. (2003) *J. Mol. Evol.* **56**, 573–586.
- Altschul, S. F., Madden, T. L., Schaffer, A. A., Zhang, J., Zhang, Z., Miller, W. & Lipman, D. J. (1997) *Nucleic Acids Res.* **25**, 3389–3402.
- Yu, J., Hu, S., Wang, J., Wong, G. K., Li, S., Liu, B., Deng, Y., Dai, L., Zhou, Y., Zhang, X., et al. (2002) *Science* **296**, 79–92.
- Eddy, S. R. (2001) HMMER: Profile Hidden Markov Models for Biological Sequence Analysis (Ver. 2.3.2), <http://hmmer.wustl.edu>.
- Katoh, K., Misawa, K., Kuma, K. & Miyata, T. (2002) *Nucleic Acids Res.* **15**, 3059–3066.
- Saitou, N. & Nei, M. (1987) *Mol. Biol. Evol.* **4**, 406–425.
- Kumar, S., Tamura, K., Jakobsen, I. B. & Nei, M. (2001) *Bioinformatics* **17**, 1244–1245.
- Swofford, D. L. (1998) PAUP*: Phylogenetic Analysis Using Parsimony (*and Other Methods) (Sinauer, Sunderland, MA), Ver. 4.
- Suzuki, Y., Glazko, G. V. & Nei, M. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 16138–16143.
- Cummings, M. P., Handley, S. A., Myers, D. S., Reed, D. L., Rokas, A. & Winka, K. (2003) *Syst. Biol.* **52**, 477–487.
- Misawa, K. & Nei, M. (2003) *J. Mol. Evol.* **57**, S290–S296.
- Becker, A., Kaufmann, K., Freialdenhoven, A., Vincent, C., Li, M. A., Saedler, H. & Theissen, G. (2002) *Mol. Genet. Genomics* **266**, 942–950.
- Chen, F., Kook, H., Milewski, R., Gitler, A. D., Lu, M. M., Li, J., Nazarian, R., Schnepf, R., Jen, K., Biben, C., et al. (2002) *Cell* **110**, 713–723.
- Shin, C. H., Liu, Z. P., Passier, R., Zhang, C. L., Wang, D. Z., Harris, T. M., Yamagishi, H., Richardson, J. A., Childs, G. & Olson, E. (2002) *Cell* **110**, 725–735.
- Kohler, C., Hennig, L., Spillane, C., Pien, S., Gruissem, W. & Grossniklaus, U. (2003) *Genes Dev.* **12**, 1540–1553.
- Pelaz, S., Ditta, G. S., Baumann, E., Wisman, E. & Yanofsky, M. F. (2000) *Nature* **405**, 200–203.
- Lamb, R. S. & Irish, V. F. (2003) *Proc. Natl. Acad. Sci. USA* **100**, 6558–6563.
- Immink, R. G., Ferrario, S., Busscher-Lange, J., Kooiker, M., Busscher, M. & Angenent, G. C. (2003) *Mol. Genet. Genomics* **268**, 598–606.
- Vandenbussche, M., Theissen, G., Van de Peer, Y. & Gerats, T. (2003) *Nucleic Acids Res.* **31**, 4401–4409.
- Takahashi, K. & Nei, M. (2000) *Mol. Biol. Evol.* **17**, 1251–1258.