

Dynamic programming procedure for searching optimal models to estimate substitution rates based on the maximum-likelihood method

Chengjun Zhang^a, Jia Wang^a, Weibo Xie^a, Gang Zhou^a, Manyuan Long^{b,1}, and Qifa Zhang^{a,1}

^aNational Key Laboratory of Crop Genetic Improvement and National Center of Plant Gene Research (Wuhan), Huazhong Agricultural University, Wuhan 430070, China; and ^bDepartment of Ecology and Evolution, University of Chicago, Chicago, IL 60637

Contributed by Qifa Zhang, March 25, 2011 (sent for review May 12, 2010)

The substitution rate in a gene can provide valuable information for understanding its functionality and evolution. A widely used method to estimate substitution rates is the maximum-likelihood method implemented in the CODEML program in the PAML package. A limited number of branch models, chosen based on a priori information or an interest in a particular lineage(s), are tested, whereas a large number of potential models are neglected. A complementary approach is also needed to test all or a large number of possible models to search for the globally optimal model(s) of maximum likelihood. However, the computational time for this search even in a small number of sequences becomes impractically long. Thus, it is desirable to explore the most probable spaces to search for the optimal models. Using dynamic programming techniques, we developed a simple computational method for searching the most probable optimal branch-specific models in a practically feasible computational time. We propose three search methods to find the optimal models, which explored $O(n)$ (method 1) to $O(n^2)$ (method 2 and method 3) models when the given phylogeny has n branches. In addition, we derived a formula to calculate the number of all possible models, revealing the complexity of finding the optimal branch-specific model. We show that in a reanalysis of over 50 previously published studies, the vast majority obtained better models with significantly higher likelihoods than the conventional hypothesis model methods.

likelihood-ratio test | natural selection | positive selection | synonymous substitution | nonsynonymous substitution

Estimating substitution rates is important in the investigation of functionality and evolution of genes. Natural selection can be also tested by comparing the substitution rates at synonymous and nonsynonymous sites, denoted usually as K_s and K_a , respectively (K_a = number of nonsynonymous substitutions per nonsynonymous site, K_s = number of synonymous substitutions per synonymous site). Such estimation is usually performed by analyzing the divergence of a protein-coding gene in a number of homologous sequences in different species.

The maximum-likelihood method is widely used for estimating the substitution rates of nucleotide sequences in protein-coding genes in molecular evolutionary analysis, although some of its techniques were recently debated (1, 2). The CODEML program in the PAML package (3) is among the most frequently used and utilizes a codon substitution model to infer evolutionary rates. Several approaches were incorporated into the program, including a site model, a clade model, a branch model, and a branch-site model. The widely used branch model allows estimation of the substitution rates with variable ratios of $\omega = K_a/K_s$ in different branches (lineages) in a phylogeny. Generally, $\omega > 1$ indicates positive selection, $\omega < 1$ indicates purifying selection with functional constraint, and $\omega \sim 1$ indicates neutral evolution (4).

The branch model was initially applied to the evolutionary analysis of the primate gene-encoding lysozyme (5). The analysis showed that the ω -parameter along the hominoid branch was significantly greater than 1, indicating that positive selection might have operated on it. This model has been widely used in molecular evolutionary studies and the functional analyses of

genes, and it is particularly valuable to detect positive selection after gene duplications (3). For example, a branch model analysis of the *Drosophila* retroposed gene *Dntf-2r* detected positive selection (6). The use of this model revealed that three young chimeric genes, *jingwei*, *Adh-Twain*, and *Adh-Finnegan*, underwent both early rapid evolution and subsequent slow evolution of protein sequences resulting from increased functional constraints (7, 8). Branch model analysis on the NOD26-like intrinsic proteins also detected strong selective pressure on highly constrained functional proteins and many positive selective events that might change the gene's functions after the duplication and speciation events in the plants (9).

In the branch model analysis, a range of ω -values can be chosen. The one-ratio model (ORM) assumes that all branches have the same one ω -parameter, whereas the free-ratio model (FRM) assigns a different ω -parameter to each branch in the tree for estimation. Between ORM and FRM are a limited number of hypothesis models, assuming that some specific branches have specific ratios based on a priori available information or interest in a possible positive selection on a branch(s) implied by FRM analysis. These models were explored and compared by likelihood-ratio tests (LRTs) (5, 10). Obviously, in this approach, it is imperative to have some good a priori reasons to restrict the estimate of spaces to explore. As Pond and Frost pointed out (11), however, this approach has a disadvantage, because it is not always possible to derive suitable hypotheses when no useful information is available or when no branch can be focused on in the model search. As a model-searching approach to complement the current approach, there is thus a need to search all possible models for the best model that has a globally maximum likelihood. Because all models, except the ORM and FRM, need to be specified with ω -parameters for certain branches, however, the analysis often becomes impractical, especially because all possible models often require an intractably large number of repeated computations of likelihoods.

To solve these technical difficulties, we proposed to search the most probable spaces to determine the optimal branch-specific models that have likelihoods equal or close to the globally maximum likelihood over all possible models with the least degrees of freedom (12). We developed a two-step method to count all possible branch models to reveal the complexity of the computation using CODEML. Then, motivated by the dynamic programming that is widely used in computation (13), we developed three simple and rapid methods in search of the optimal branch models in the most probable spaces for the maximum likelihood. Finally, the proposed methods were assessed by the lysozyme sequences of primate species (5) and reanalysis of 50 previously published

Author contributions: C.Z., M.L., and Q.Z. designed research; C.Z., J.W., W.X., and G.Z. performed research; C.Z. and M.L. analyzed data; and C.Z., M.L., and Q.Z. wrote the paper.

The authors declare no conflict of interest.

Freely available online through the PNAS open access option.

¹To whom correspondence may be addressed. E-mail: mlong@uchicago.edu or qifazhang@mail.hzau.edu.cn.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1018621108/-DCSupplemental.

studies. Through these analyses, we show that our simple methods can obtain globally better models with significantly higher likelihoods than the current approach that compares the models on the branches of particular interest. Because the current approach relies on the hypothesized branches of interest to test positive selection, we call it the “conventional hypothesis model.”

Results

Large Number of Possible Branch Models. We calculated the number of all possible branch models using a two-step strategy, which is used in a program written using the Perl script (*SI Appendix*). In the first step, we defined a model that included the number of ω 's and the branch number for each ω , recording this model in a configuration. For example, for a tree of four branches with three sequences, assuming two ω -values, ω_1 for one branch and ω_2 for the other three branches, we record this configuration as a vector (1 ω_1 , 3 ω_2), or simply (1, 3). We developed a traversing algorithm to find all the configurations of a variety of ratios. In the second step, we calculated all possible branch models with each configuration following the two formulas that we derived, as shown below.

Imagine a phylogeny of **six branches** with four sequences (*SI Appendix*, Fig. 1). The models for this tree can be divided into six groups [ranging from ORMs, to two-ratio models, up to the six-ratio model (FRM)], and in each group, the models can be divided into several configurations. For example, it has three configurations in three-ratio models: the first configuration has one branch with ω_1 , one branch with ω_2 , and the other four branches with ω_3 , expressed as (1, 1, 4); the second configuration has one branch with ω_1 , two branches with ω_2 , and the other three branches with ω_3 , expressed as (1, 2, 3); and the third configuration has three two branches with ω_1 , ω_2 , and ω_3 , respectively, expressed as (2, 2, 2).

The number of the models for the first configuration (1, 1, 4) can be calculated and expressed as K_{31} , the numbers of the models for the second and the third configurations [(1, 2, 3) and (2, 2, 2)] as K_{32} and K_{33} , respectively. In K_{32} , because the components in the configuration are not equal to each other, all possible combinations are

$$K_{32} = C_6^1 \times C_5^2 = 60$$

Because the first configuration has two different types (the numbers of branches) of components in K_{3I} and the third configuration has three components each with the same number of branches in K_{33} ,

$$K_{3l} = C_6^1 \times C_5^1 \div 2! = 15$$

Where the $2!$ is the denominator because we need only the combination, the order of arrangement does not matter. Similarly, we have

$$K_{33} = C_6^2 \times C_4^2 \div 3! = 15$$

In general, for a phylogeny with n branches, we use K_{mj} to denote the possible model numbers for the j th configuration with m ω -parameters; q_{ij} denotes the branch numbers of the i th ω -parameter of the j th configuration. By definition, we have

$$\sum_{i=1}^m q_{ij} = n, m \in (1 \text{ to } n).$$

When $q_{xj} \neq q_{yj}$ ($x \neq y, x, y \in (1 \text{ to } m), q_{0j} = 0$), the formula to calculate K_{mj} can be expressed as

$$K_{mj} = \prod_{i=1}^{m-1} C_{n-\sum q_{(i-1)j}}^{q_{ij}} \quad [1]$$

When there exist x and y variables, let $q_{xj} = q_{yj}$ [$x \neq y, x, y \in (1 \text{ to } m)$], $q_{0j} = 0$ (A_g means having g groups and A_g components

in the configuration, which have the same branch numbers), and thus we have

$$K_{mj} = \frac{\prod_{i=1}^{m-1} C_{n-\sum_{l=1}^i q_{(i-1)l}}^{q_{ij}}}{\prod_{l=1}^g A_l!} \quad [2]$$

By means of this approach, to illustrate the intractably large number of possible branch models visually, all configuration numbers and possible model numbers of phylogeny for 3, 4, 6, 8, 10, and 12 sequences are shown in Table 1 for all possible ω -values; an example of the details of the configuration and model is provided in SI Appendix.

Dynamic Programming Algorithms for Searching Optimal Branch Models. Despite present-day rapidly increasing computing powers, it is impractical to use the traversing algorithm to explore all models, as shown in Table 1. We developed three simplified methods for searching optimal models by using dynamic programming algorithms. We attempted to reduce computation to a practical workable level by exploring the most likely space that contains the maximum likelihood.

Method 1. Fig. 1A summarizes the procedure we propose. First, calculate all possible configurations for single-branch two-ratio models (SBTRMs), in which only one branch is labeled with ω_1 and all other branches are assumed to be background ratio ω_0 . Obviously, the log likelihood (lnL) values for n SBTRMs need to be calculated when the analyzed phylogeny has n branches. Second, the lnL values of all n SBTRMs are compared and sorted from maximum to minimum; the model with the maximum lnL value is considered the optimal model within two-ratio models. The branch labeled with ω_1 in the maximum lnL value model is recorded as B_1 , the branch labeled with ω_1 in the model that has the second greatest lnL value is recorded as B_2 , and so on until B_n . Then, all the optimal models of the remaining variety of ratios are generated directly. For the optimal three-ratio model, branch B_1 is labeled as ω_1 , branch B_2 is labeled as ω_2 , and all other branches are assumed to have a background ratio ω_0 and optimal models for four ratios to an “ $n - 1$ ” ratio as well. Finally, the $n - 2$ optimal models can be “predicted” in this way, and the likelihoods of these predicted models can be calculated and compared with each other to determine the final optimal model that has the maximum likelihood in the sense that the likelihood is significantly better than the likelihood of other optimal models and has the least degrees of freedom if there are more than one solutions that are not significantly different.

Method 2. This method can be described in $n - 2$ rounds with two steps in each round of iterations, as shown in Fig. 1B. The first step generates models and calculates InLs for all these models; the second step is to record the specific branch of the optimal model of this round, which is used for generating models in the next round. The models in the first round are all SBTRMs. The branch labeled with ω_1 in the maximum InL value model is recorded as B_1 . In the second round, $n - 1$ three-ratio models are generated by adding one more branch with one more ratio (ω_2) in addition to B_1 , whereas all other $n - 2$ branches have the background ratio ω_0 . The InLs for all $n - 1$ three-ratio models are calculated and compared with each other. The branch labeled ω_2 of the optimal

Table 1. Configurations and possible models

Sequence no.	Branch no.	Configuration no.	No. of possible models
3	4	3	15
4	6	9	203
6	10	40	115,975
8	14	133	190,899,322
10	18	383	6.821E + 11
12	22	1,000	4.507E + 15



Fig. 1. Sketch of the proposed methods. (A) Method 1: Searching optimal models with more than two ω -parameters directly based on the sorted results of the SBTRM (2–5 ω -parameters exemplified). The different models with the same number of ω -parameters are arranged from high-likelihood to low-likelihood values. (B) Method 2: Searching optimal models with ω -parameters until there are no free branches, based on the maximum-likelihood value model from the last round. The different models with the same number of ω -parameters are also arranged from high-likelihood to low-likelihood values. (C) Method 3: Searching optimal models by iteration. In (A–C), one color stands for one ω -parameter.

model having the maximum lnL value in all $n - 1$ three-ratio models is recorded as B_2 . This process is reiterated until all $n - 1$ ω 's are calculated. In total, $(n + 1) \cdot n/2$ models are generated and calculated; $n - 2$ optimal models of a variety of ratios are obtained and can be compared with each other, including the ORM and FRM, by LRT to determine the final optimal models.

Method 3. This is a modification of method 2 (Fig. 1C), to consider general cases of one ratio with more than one branch. First, similar to method 2, all the SBTRMs belonging to the configuration $(1, n - 1)$ are calculated in this step and the optimal model of SBTRMs (assumed to be A) is determined. This optimal model has only one branch B_1 , which is labeled as ω_1 . Then, in the second step, other $n - 1$ two-ratio models are generated, which have another branch labeled as ω_1 in addition to B_1 ; these models

belong to the configuration $(2, n - 2)$. After calculation, the optimal model is found (assumed to be B). If the difference in the lnL values between A and B is greater than k (k is a threshold that can be defined by the user to decide if one model is better than another model with the same degree of freedom when they have different branches with same ratio, the default $k = 0.5$), the models belonging to configuration $(3, n - 3)$ are generated and calculated and the optimal model C is compared with B. Such iterations continue until the difference between the two optimal models is less than k . Clearly, the optimal model obtained from the penultimate iteration will become the final optimal two-ratio model. Note that the threshold value of k will determine the number of iterations; the more iterations calculated, the more the branch would be labeled with the same ω and the fewer would be

Our finding that most final optimal models detected by our methods are significantly better than the conventional hypothesis models was further confirmed by our subsequent studies of 50 gene families. We collected the sequences from these gene families from 40 original studies (14–53), and we then applied our methods to analyze these data and to compare them with the previous results of conventional hypothesis models using the maximum-likelihood method. These analyses are summarized in *SI Appendix*, Table 5. We found that in gene families (or cases) 40 and 45, the InL value of the final optimal model our method detected and that of the conventional hypothesis model were congruent with each other; in case 38, there was no difference between the final optimal model and the current hypothesis model ($P > 0.05$). However, we were surprised to see that for the vast majority of the rest 47 cases, the InL values for the final optimal models are significantly higher than the InL values for the conventional hypothesis models ($P < 0.001$). In these cases, 22 are significant at the level $P \leq 10^{-5}$ and 8 of them even at level $P \leq 10^{-10}$. More details of the conventional hypothesis models, our optimal models, and the 50 phylogenies are provided in the data in *SI Appendix*.

Discussion

In principle, the maximum-likelihood method was proposed to find the most probable estimates, given a phylogeny of homologous sequences. It is also clear that FRM cannot guarantee a parsimonious model. It is thus expected to find the globally most probable estimate by performing an exhaustive search of the most probable model from all possible models. Such a search is often impractically time-consuming, however, because of a huge number of possible models for a tree with even a small number of sequences. The problems in calculating all possible models were raised previously (54). Our method calculated the number of all possible models for a rooted tree in full agreement with the Bell number that was used to calculate the number in an unrooted tree (54). We proposed these simplified methods to find the most probable estimates of substitution rates with the least degrees of freedom in hypothesis testing compared with the FRM. The present study highlights the finding that the optimal models obtained from the three methods described in the following text via a dynamic programming approach are extremely close to the best model obtained from the traversing algorithm. The former simple methods use a reasonably short time, whereas the latter exhaustive search is often impractical in computing time for a large dataset, such as that used in this paper.

Compared with the previous analysis of the lysozyme dataset using the conventional hypothesis models (5), our simple method 3 obtained even significantly higher likelihoods than the previous two-ratio and three-ratio hypothesis models (−842.09 vs. −844.10, $P = 0.045$; −842.09 vs. −844.10, $P = 0.045$; Table 3). The advantage of our methods is further confirmed by our large-scale case analyses of 50 previously reported gene families using the conventional hypothesis method. In these 50 cases, we found that for 47 cases (94%), our final optional models had significantly

higher likelihoods than the conventional hypothesis models and that there were only 3 cases not having significantly different likelihoods (*SI Appendix*, Table 5). The most significant differences were observed in the Chalcone Synthase Genes of *Dendranthema* (case 6: $2\Delta l = 198.91$, $df = 11$, $P < 1e-14$), the Phytochrome Gene Family in Angiosperms (case 3: $2\Delta l = 206.25$, $df = 8$, $P < 1e-14$), and the recently duplicated M_y -type MADS-box genes in *Petunia* (case 13: $2\Delta l = 175.71$, $df = 16$, $P < 1e-14$).

The compared models in the branch model should be nested, as suggested for the LRT (55). To make a more general comparison involving the models that do not meet such a condition, we also used the Akaike's information criterion (AIC) (56) method in analyses of these 50 cases, with the AIC values of the analyzed models in the data in *SI Appendix*. Again, except for 2 cases in which the final optimal model is congruent with the conventional hypothesis model, all other final optimal models have the lowest AIC value in 48 cases, even in the case (case 38) that failed in the LRT also getting a lower AIC than the conventional hypothesis model.

In additional, in the color vision gene (SWS2, case 17), in which $2\Delta l = 34.30$, $df = 6$, $P = 5.90e-006$, our optimal models suggest positive selection on the lineage *Sinocyclocheilus purpureus* (fix $\omega_{\text{purpureus}} = 1$ model vs. free $\omega_{\text{purpureus}}$ model: $2\Delta l = 5.74$, $df = 1$, $P = 0.017$), which was not detected by the previous analysis using the conventional hypothesis method. These case analyses indicate that most previous reports missed the optional models and that the conventional hypothesis method can easily miss the globally most probable model. Our methods appear to be able to detect more significant models than the conventional hypothesis method.

Although the present methods provide simplified computational procedures for the maximum-likelihood analysis, caution should be urged in using these methods. The first caveat is that, like any other phylogeny-related study, if the phylogeny tree is inaccurate or incorrect (e.g., an incorrect inference of the orthologous-paralogous relationship), the estimates of the maximum-likelihood method, which is dependent on the tree, are meaningless. The second caution is that when many models explored by our methods detected a large ω -value in some lineages, this finding may not immediately suggest positive selection, because a statistical test for its significance is needed. The model comparison as implemented by the original branch model (5) is necessary using, for example, the nested model-based LRT or AIC discussed above. Third, we note here that method 3 seems to perform better than methods 1 and 2 in detecting final optimal models using the one gene-data analysis of lysozyme. We recommended using all three methods for more genes and comparing their performance. It would be a wise practice to start from method 1 when analyzing a large dataset to gain some useful insight because of its brief computation time.

Methods

Sequence. The sequences used in calculation of all possible models to evaluate our three methods are taken from previous work (5) and can be obtained in the PAML package in the example of lysozyme. For the reanalysis of the 50 previous studies, we utilized either available sequence alignments provided

Table 3. Substitution rate values of final best model, final optimal model, and hypothesis model

	Final best model*	Final optimal model (methods 1 and 2) [†]	Final optimal model (method 3) [‡]	Hypothesis TRM [§]	Hypothesis ThreeRM [¶]
InL values	−843.25	−844.99	−842.09	−844.10	−844.10
ω_0	0.497	1.075	0.611	0.579	0.579
ω_1	4.466	0.0001	0.0001	4.224	4.333
ω_2	—	—	4.288	—	4.112
k	5.021	4.921	5.000	5.008	5.007

TRM, two-ratio model; ThreeRM, three-ratio model.

For the following phylogeny with markers for models (#1, ω_1 ; #2, ω_2):

*(((Ssc_squirrelM,Cja_marmoset),Hla_gibbon#1),Mmu_rhesus#1,(Cgu_Can_colobus,Pne_langur)#1))

[†]((Ssc_squirrelM,Cja_marmoset#1),Hla_gibbon),(Mmu_rhesus,(Cgu_Can_colobus,Pne_langur)))

[‡]((Ssc_squirrelM,Cja_marmoset#1),Hla_gibbon#2),(Mmu_rhesus#2,(Cgu_Can_colobus,Pne_langur)#2))

[§]((Ssc_squirrelM,Cja_marmoset),Hla_gibbon#1),(Mmu_rhesus,(Cgu_Can_colobus,Pne_langur)#1))

[¶]((Ssc_squirrelM,Cja_marmoset),Hla_gibbon#1),(Mmu_rhesus,(Cgu_Can_colobus,Pne_langur)#2))

in the literature or regenerated sequence realignments using MEGA 4.0 (57) when the original alignments were not available.

Calculating the Entire Range of Possible Models. We generated seven datasets of **six sequences** from these lysozyme sequences by deleting one sequence from seven. All possible models (115,975 possible models in one dataset) of these seven datasets were generated by the traversing algorithm (*SI Appendix*) and calculated. It took almost 4 d to finish all the calculations for one dataset, and according to this, it may take **160 d to calculate all possible 4,213,597 models of the seven sequences** on the server (Dawning Information Industry), which has eight AMD Opteron 2376 processors with the operation system Linux AS 5. The phylogeny used in the calculations was built by MEGA 4.0 with the neighbor-joining method (57).

- Nozawa M, Suzuki Y, Nei M (2009) Reliabilities of identifying positive selection by the branch-site and the site-prediction methods. *Proc Natl Acad Sci USA* 106:6700–6705.
- Yang Z, Nielsen R, Goldman N (2009) In defense of statistical methods for detecting positive selection. *Proc Natl Acad Sci USA* 106:E95–E96, author reply E96.
- Yang Z (2007) PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol* 24:1586–1591.
- Li W (1997) *Molecular Evolution* (Sinauer, Sunderland, MA).
- Yang Z (1998) Likelihood ratio tests for detecting positive selection and application to primate lysozyme evolution. *Mol Biol Evol* 15:568–573.
- Betrán E, Long M (2003) *Dntf-2r*, a young *Drosophila* retroposed gene with specific male expression under positive Darwinian selection. *Genetics* 164:977–988.
- Jones CD, Begun DJ (2005) Parallel evolution of chimeric fusion genes. *Proc Natl Acad Sci USA* 102:11373–11378.
- Jones CD, Custer AW, Begun DJ (2005) Origin and evolution of a chimeric fusion gene in *Drosophila subobscura*, *D. madeirensis* and *D. guanche*. *Genetics* 170:207–219.
- Liu Q, et al. (2009) Divergence in function and expression of the NOD26-like intrinsic proteins in plants. *BMC Genomics* 10:313.
- Anisimova M, Bielawski JP, Yang Z (2001) Accuracy and power of the likelihood ratio test in detecting adaptive molecular evolution. *Mol Biol Evol* 18:1585–1592.
- Pond SL, Frost SD (2005) A genetic algorithm approach to detecting lineage-specific variation in selection pressure. *Mol Biol Evol* 22:478–485.
- Goldman N (1993) Statistical tests of models of DNA substitution. *J Mol Evol* 36:182–198.
- Bellman R (1957) *Dynamic Programming* (Princeton Univ Press, Princeton). (2003) Paperback edition (Dover, New York).
- Jiggins FM, Hurst GD, Yang Z (2002) Host-symbiont conflicts: Positive selection on an outer membrane protein of parasitic but not mutualistic Rickettsiaceae. *Mol Biol Evol* 19:1341–1349.
- Marcussen T, Oxelman B, Skog A, Jakobsen KS (2010) Evolution of plant RNA polymerase IV/V genes: Evidence of subneofunctionalization of duplicated NRDP2/NRPE2-like paralogs in *Viola* (*Violaceae*). *BMC Evol Biol* 10:45.
- Yang Z, Nielsen R (2002) Codon-substitution models for detecting molecular adaptation at individual sites along specific lineages. *Mol Biol Evol* 19:908–917.
- Alba R, Kelmenson PM, Cordonnier-Pratt M-M, Pratt LH (2000) The phytochrome gene family in tomato and the rapid differential evolution of this family in angiosperms. *Mol Biol Evol* 17:362–373.
- Huttley GA, et al. (2000) Adaptive evolution of the tumour suppressor *BRCA1* in humans and chimpanzees. *Australian Breast Cancer Family Study*. *Nat Genet* 25:410–413.
- Yang J, Huang J, Gu H, Zhong Y, Yang Z (2002) Duplication and adaptive evolution of the chalcone synthase genes of *Dendranthema* (*Asteraceae*). *Mol Biol Evol* 19:1752–1759.
- Merritt TJ, Quattro JM (2001) Evidence for a period of directional selection following gene duplication in a neurally expressed locus of triosephosphate isomerase. *Genetics* 159:689–697.
- Yang Z (2002) Inference of selection from multiple species alignments. *Curr Opin Genet Dev* 12:688–694.
- Yang J, Gu H, Yang Z (2004) Likelihood analysis of the chalcone synthase genes suggests the role of positive selection in morning glories (*Ipomoea*). *J Mol Evol* 58:54–63.
- Aguileta G, Bielawski JP, Yang Z (2004) Gene conversion and functional divergence in the beta-globin gene family. *J Mol Evol* 59:177–189.
- Schein M, Yang Z, Mitchell-Olds T, Schmid KJ (2004) Rapid evolution of a pollen-specific oleosin-like gene family from *Arabidopsis thaliana* and closely related species. *Mol Biol Evol* 21:659–669.
- Narita Y, Oda S, Takenaka O, Kageyama T (2010) Lineage-specific duplication and loss of pepsinogen genes in hominoid evolution. *J Mol Evol* 70:313–324.
- Larmuseau MH, Huyse T, Vancampenhout K, Van Houdt JK, Volckaert FA (2010) High molecular diversity in the rhodopsin gene in closely related goby fishes: A role for visual pigments in adaptive speciation? *Mol Phylogenet Evol* 55:689–698.
- Arora R, et al. (2007) MADS-box gene family in rice: Genome-wide identification, organization and expression profiling during reproductive development and stress. *BMC Genomics* 8:242.
- Bemer M, Gordon J, Weterings K, Angenot GC (2010) Divergence of recently duplicated *My*-type MADS-box genes in *Petunia*. *Mol Biol Evol* 27:481–495.
- Schienman JE, Holt RA, Auerbach MR, Stewart CB (2006) Duplication and divergence of 2 distinct pancreatic ribonuclease genes in leaf-eating African and Asian colobine monkeys. *Mol Biol Evol* 23:1465–1479.
- Yu L, et al. (2010) Adaptive evolution of digestive RNASE1 genes in leaf-eating monkeys revisited: New insights from ten additional colobines. *Mol Biol Evol* 27:121–131.
- Zhang J (2006) Parallel adaptive origins of digestive RNases in Asian and African leaf monkeys. *Nat Genet* 38:819–823.
- Zhang J, Zhang YP, Rosenberg HF (2002) Adaptive evolution of a duplicated pancreatic ribonuclease gene in a leaf-eating monkey. *Nat Genet* 30:411–415.
- Viaene T, et al. (2009) *Pistillata*—Duplications as a mode for floral diversification in (Basal) asterids. *Mol Biol Evol* 26:2627–2645.
- Li Z, Gan X, He S (2009) Distinct evolutionary patterns between two duplicated color vision genes within cyprinid fishes. *J Mol Evol* 69:346–359.
- Zhou D, et al. (2009) Duplication and adaptive evolution of the *COR15* genes within the highly cold-tolerant *Draba* lineage (Brassicaceae). *Gene* 441:36–44.
- Miwa H, et al. (2009) Adaptive evolution of *rbcl* in *Conocephalum* (Hepaticeae, bryophytes). *Gene* 441:169–175.
- Zhao H, et al. (2009) The evolution of color vision in nocturnal mammals. *Proc Natl Acad Sci USA* 106:8980–8985.
- Weadick CJ, Chang BS (2009) Molecular evolution of the betagamma lens crystallin superfamily: Evidence for a retained ancestral function in gamma N crystallins? *Mol Biol Evol* 26:1127–1142.
- Wang Z, et al. (2009) Adaptive evolution of 5'HoxD genes in the origin and diversification of the cetacean flipper. *Mol Biol Evol* 26:613–622.
- Dorus S, Freeman ZN, Parker ER, Heath BD, Karr TL (2008) Recent origins of sperm genes in *Drosophila*. *Mol Biol Evol* 25:2157–2166.
- Schulenburg H, Boehnisch C (2008) Diversification and adaptive sequence evolution of *Caenorhabditis* lysozymes (Nematoda: Rhabditidae). *BMC Evol Biol* 8:114.
- Zhang L (2008) Adaptive evolution and frequent gene conversion in the brain expressed X-linked gene family in mammals. *Biochem Genet* 46:293–311.
- Zhang W, et al. (2008) Molecular evolution of *PISTILLATA*-like genes in the dogwood genus *Cornus* (Cornaceae). *Mol Phylogenet Evol* 47:175–195.
- Storz JF, Hoffmann FG, Opazo JC, Moriama H (2008) Adaptive functional divergence among triplicated alpha-globin genes in rodents. *Genetics* 178:1623–1638.
- Muggia L, Schmitt I, Grube M (2008) Purifying selection is a prevailing motif in the evolution of ketoacyl synthase domains of polyketide synthases from lichenized fungi. *Mycol Res* 112:277–288.
- Padhi A, Verghese B (2007) Evidence for positive Darwinian selection on the hepcidin gene of Perciform and Pleuronectiform fishes. *Mol Divers* 11:119–130.
- Royer B, et al. (2007) Molecular evolution of the human *SRPX2* gene that causes brain disorders of the Rolandic and Sylvian speech areas. *BMC Genet* 8:72.
- Ding K, McDonough SJ, Kullo IJ (2007) Evidence for positive selection in the C-terminal domain of the cholesterol metabolism gene *PCSK9* based on phylogenetic analysis in 14 primate species. *PLoS ONE* 2:e1098.
- Hahn Y, Jeong S, Lee B (2007) Inactivation of *MOXD2* and *S100A15A* by exon deletion during human evolution. *Mol Biol Evol* 24:2203–2212.
- Padhi A, Verghese B, Otta SK, Varghese B, Ramu K (2007) Adaptive evolution after duplication of penaeidin antimicrobial peptides. *Fish Shellfish Immunol* 23:553–566.
- Zhang Q, et al. (2007) Rapid evolution, genetic variations, and functional association of the human spermatogenesis-related gene *NYD-SP12*. *J Mol Evol* 65:154–161.
- Hou ZC, Xu GY, Su Z, Yang N (2007) Purifying selection and positive selection on the myxovirus resistance gene in mammals and chickens. *Gene* 396:188–195.
- Wang Y, et al. (2007) Isolation and characterization of a putative class E gene from *Taihangia rupestris*. *J Integr Plant Biol* 49:343–350.
- Sanderson MJ (1998) *Estimating rate and time in molecular phylogenies: beyond the molecular clock?* *Plant Molecular Systematics*, eds Soltis P, Soltis D, Doyle J (Chapman & Hall, New York), 2nd Ed, pp 242–264.
- Yang Z, Nielsen R, Hasegawa M (1998) Models of amino acid substitution and applications to mitochondrial protein evolution. *Mol Biol Evol* 15:1600–1611.
- Akaike H (1974) A new look at the statistical model identification. *IEEE Trans Automat Contr* 19:716–723.
- Tamura K, Dudley J, Nei M, Kumar S (2007) MEGA4: Molecular Evolutionary Genetics Analysis (MEGA) software version 4.0. *Mol Biol Evol* 24:1596–1599.

Supporting Information

Supporting Information 1

Table S1. Log likelihood values of a variety of ratio models of the remaining six datasets

Supporting Information 2

Table S2. Likelihood ratio statistics for testing a variety of ratio models

Supporting Information 3

Table S3. Comparison between final best model and final optimal models

Supporting Information 4

Table S4. Comparison between final optimal models and hypothesis models

Supporting Information 5

Table S5 Comparing Final Optimal Model of OBSM with Hypothesis model suggested in 50 cases

Supporting Information 6

Figure S1. An example of a phylogeny with 4 species (or genes). S1, S2, S3, S4 denote the homologous sequences of 4 species (or genes); 1~6 are six branches. For three-ratio branch-models, there're three types of configurations (1, 1, 4), (1, 2, 3), (2, 2, 2). For the first configuration (1, 1, 4), we can choose branch 1 to set as ω_1 , and branch 2 as ω_2 , and the remaining 4 branches each as ω_3 . Alternatively, we can choose branch 1 as ω_1 , branch 3 as ω_2 , and the remaining 4 branches each as ω_3 . This procedure can go on until all the topologies are considered. For the second configuration (1, 2, 3), a similar procedure assigns ω_1 , ω_2 , ω_3 to one branch, two branches, the remaining three branches respectively.

Supporting Information 7

Figure S2. The phylogenies generated from 6 dataset (A-F)s each with 6 lysozyme sequences chosen from the 7 lysozyme sequences of the original lysozyme dataset (5) (Figure 2 in the main text is the phylogeny from the last one of the all 7 possible datasets). Red color: the branch with ω_h ; Green color the branch with ω_c .

Supporting Information 8

Four Perl scripts were used in our work. The first part is used for searching the configuration and counting all the possible model numbers; the second part is used for

exploring all possible models; the third one is to do a likelihood ratio test based on the result of method 2; and the last one is to do a likelihood ratio test based on the result of method 3.

Supporting Information 9

The example of six sequences in detail with configurations and the number of possible models of each configuration.

Supporting Information 10

supplementary data for 50 cases

Table S1 Maximum Log Likelihood Values in variety ratio models of remaining six data sets of lysozymes (5)

		TwoRM	ThreeRM	FourRM	FiveRM	SixRM	SevenRM	EightRM	NineRM
Total models		511	9330	34105	42525	22827	5880	750	45
Max_lnL		-834.142582	-832.129502	-831.864334	-831.675345	-831.626547	-831.62079	-831.618484	-831.618484
Method I	lnL	-835.658129	-834.595675	-833.51955	-832.210261	-831.938366	-831.883059	-831.798407	-831.726768
		30	670	1889	757	467	408	122	21
Method II	lnL	-835.658129	-834.142582	-833.508003	-832.817495	-831.938366	-831.734028	-831.692257	-831.624213
		30	365	1870	2454	467	99	31	3
Method III	lnL	-834.142582	-833.508003	-831.942625	-831.736076	-831.675345	-831.626547	-831.62079	/
		1	90	8	6	10	3	2	
Max_lnL		-827.909168	-826.390537	-824.878946	-824.654224	-824.63444	-824.629197	-824.629197	-824.629197
Method I	lnL	-830.332163	-827.589533	-826.570755	-825.375631	-825.053261	-824.887567	-824.75887	-824.632675
		47	142	763	508	561	461	99	8
Method II	lnL	-830.332163	-827.589533	-826.535655	-825.375631	-824.881462	-824.657052	-824.637086	-824.631779
		47	142	720	508	422	125	36	6
Method III	lnL	-830.332163	-826.570755	-825.375631	-824.881462	-824.657052	-824.637086	-824.631779	-824.629197
		47	23	56	126	52	24	10	4
Max_lnL		-828.19381	-826.908267	-825.613278	-825.538777	-825.484354	-825.464357	-825.464357	-825.464357
Method I	lnL	-830.16028	-827.516759	-826.299258	-826.147571	-825.962545	-825.902357	-825.624112	-825.486661
		37	43	48	497	824	717	160	9
Method II	lnL	-830.16028	-827.516759	-826.299258	-825.813171	-825.62463	-825.54421	-825.485788	-825.464357
		37	43	48	221	277	140	13	1
Method III	lnL	-829.537486	-827.236751	-826.280292	-825.813184	-825.624662	-825.518913	-825.464357	-825.464357
		19	22	45	222	278	53	2	1
Max_lnL		-833.577481	-832.725657	-831.770264	-831.661651	-831.570233	-831.559893	-831.556326	-831.555043

Method I	lnL	-835.678272	-833.281509	-832.07163	-832.054536	-831.871882	-831.702168	-831.636489	-831.621315
		57	79	19	413	695	286	77	14
Method II	lnL	-835.678272	-833.281509	-832.07163	-831.904479	-831.754968	-831.674178	-831.613467	-831.565594
		57	79	19	215	350	184	39	4
Method III	lnL	-834.184208	-832.866043	-832.705555	-832.562737	-832.469627	-832.3931	-832.415991	-832.329388
		6	13	337	2029	2990	1611	389	34
Max_lnL		-826.875638	-826.265078	-825.669814	-825.631417	-825.598444	-825.590278	-825.589132	-825.589132
Method I	lnL	-829.704697	-828.509613	-827.351366	-826.33224	-825.816763	-825.773872	-825.637962	-825.615311
		55	1180	2774	1693	429	359	68	11
Method II	lnL	-829.704697	-827.855504	-826.682121	-826.265078	-825.897248	-825.725926	-825.630145	-825.597695
		55	448	1018	1326	775	284	47	5
Method III	lnL	-827.292062	-826.871296	-826.51281	-826.334256	-826.239186	-826.208553	-826.200039	/
		3	104	662	1706	2134	1334	339	
Max_lnL		-838.341428	-837.491141	-836.91779	-836.89203	-836.872364	-836.870807	-836.86946	-836.86946
Method I	lnL	-841.727813	-839.998874	-838.931056	-838.00642	-837.243223	-836.900313	-836.896359	-836.895105
		90	1374	3497	3838	755	148	89	22
Method II	lnL	-841.727813	-839.963719	-838.716878	-837.764604	-837.330024	-837.107644	-836.89068	-836.870807
		90	1274	2814	2259	961	443	80	4
Method III	lnL	-838.341428	-837.915019	-837.684536	-837.466694	-837.448306	-837.447477	/	
		1	42	314	868	1648	1207		

Note: The pair lnL value of the hypothesis models ($\omega_h = \omega_c$ and $\omega_h \neq \omega_c$) for seven data sets are (-844.097468, -844.096995), (-835.868563, -835.867062), (-829.92476, -828.725264), (-829.470066, -828.554522), (-834.184208, -833.480516), (-828.354787, -827.855504), (-840.41675, -839.963719). The lnL value of ORM for seven data sets are -847.329662, -838.915094, -833.338934, -833.239572, -837.801279, -832.565192, -844.104086. The lnL values of final optimal models find out by three methods are marked with red color. The final best models are marked with blue color.

Table S2 Likelihood Ratio Statistics for comparing various ratio models

p-value	ThreeRM	FourRM	FiveRM	SixRM	SevenRM	EightRM	NineRM		
TwoRM	0.082(1)	0.175(2)	0.285(3)	0.416(4)	0.559(5)	0.685(6)	0.787(7)	All results	Data set I
ThreeRM		0.496(1)	0.68(2)	0.824(3)	0.923(4)	0.969(5)	0.989(6)		
FourRM			0.578(1)	0.802(2)	0.93(3)	0.978(4)	0.994(5)		
FiveRM				0.715(1)	0.932(2)	0.986(3)	0.997(4)		
SixRM					0.928(1)	0.994(2)	1(3)		
SevenRM						0.95(1)	0.998(2)		
EightRM							1(1)		
TwoRM	0.194(1)	0.096(2)	0.157(3)	0.17(4)	0.265(5)	0.359(6)	0.461(7)		
ThreeRM		0.084(1)	0.172(2)	0.192(3)	0.313(4)	0.426(5)	0.543(6)	Method I	
FourRM			0.464(1)	0.417(2)	0.621(3)	0.749(4)	0.847(5)		
FiveRM				0.271(1)	0.539(2)	0.708(3)	0.83(4)		
SixRM					0.875(1)	0.914(2)	0.966(3)		
SevenRM						0.693(1)	0.885(2)		
EightRM							0.767(1)		
TwoRM	0.171(1)	0.096(2)	0.11(3)	0.17(4)	0.238(5)	0.296(6)	0.387(7)	Method II	
ThreeRM		0.094(1)	0.125(2)	0.207(3)	0.297(4)	0.368(5)	0.476(6)		
FourRM			0.244(1)	0.417(2)	0.553(3)	0.627(4)	0.741(5)		
FiveRM				0.533(1)	0.692(2)	0.743(3)	0.848(4)		
SixRM					0.555(1)	0.652(2)	0.804(3)		
SevenRM						0.477(1)	0.726(2)		
EightRM							0.714(1)	Method III	
TwoRM	0.016(1*)	0.041(2*)	0.089(3)	0.133(4)	0.21(5)	0.306(6)			
ThreeRM		0.441(1)	0.704(2)	0.745(3)	0.857(4)	0.93(5)			
FourRM			0.742(1)	0.726(2)	0.865(3)	0.945(4)			
FiveRM				0.466(1)	0.731(2)	0.886(3)			
SixRM					0.759(1)	0.945(2)			
SevenRM						0.888(1)		All results	
TwoRM	0.045(1*)	0.102(2)	0.177(3)	0.284(4)	0.411(5)	0.538(6)	0.654(7)		
ThreeRM		0.466(1)	0.635(2)	0.8(3)	0.907(4)	0.961(5)	0.985(6)		
FourRM			0.539(1)	0.788(2)	0.922(3)	0.974(4)	0.992(5)		
FiveRM				0.755(1)	0.947(2)	0.99(3)	0.998(4)		
SixRM					0.915(1)	0.992(2)	0.999(3)		
SevenRM						0.946(1)	0.998(2)		
EightRM							1(1)		
TwoRM	0.145(1)	0.118(2)	0.075(3)	0.114(4)	0.183(5)	0.259(6)	0.345(7)	Method I	
ThreeRM		0.142(1)	0.092(2)	0.15(3)	0.246(4)	0.348(5)	0.453(6)		
FourRM			0.106(1)	0.206(2)	0.351(3)	0.487(4)	0.61(5)		
FiveRM				0.461(1)	0.721(2)	0.844(3)	0.915(4)		
SixRM					0.739(1)	0.869(2)	0.935(3)		
SevenRM						0.681(1)	0.855(2)		
EightRM							0.705(1)		

TwoRM	0.082(1)	0.116(2)	0.128(3)	0.114(4)	0.165(5)	0.243(6)	0.327(7)	Method II	Data set III
ThreeRM		0.26(1)	0.266(2)	0.221(3)	0.307(4)	0.428(5)	0.539(6)		
FourRM			0.24(1)	0.208(2)	0.315(3)	0.458(4)	0.583(5)		
FiveRM				0.185(1)	0.338(2)	0.522(3)	0.665(4)		
SixRM					0.523(1)	0.782(2)	0.89(3)		
SevenRM						0.773(1)	0.896(2)		
EightRM							0.712(1)		
TwoRM	0.26(1)	0.111(2)	0.186(3)	0.294(4)	0.412(5)	0.538(6)		Method III	
ThreeRM		0.077(1)	0.17(2)	0.3(3)	0.439(4)	0.582(5)			
FourRM			0.52(1)	0.765(2)	0.889(3)	0.958(4)			
FiveRM				0.727(1)	0.896(2)	0.973(3)			
SixRM					0.755(1)	0.947(2)			
SevenRM						0.915(1)			
TwoRM	0.081(1)	0.048(2*)	0.089(3)	0.162(4)	0.255(5)	0.363(6)	0.476(7)	All results	
ThreeRM		0.082(1)	0.176(2)	0.319(3)	0.474(4)	0.62(5)	0.741(6)		
FourRM			0.503(1)	0.783(2)	0.919(3)	0.974(4)	0.992(5)		
FiveRM				0.842(1)	0.975(2)	0.997(3)	1(4)		
SixRM					0.918(1)	0.995(2)	1(3)		
SevenRM						1(1)	1(2)		
EightRM							1(1)		
TwoRM	0.019(1*)	0.023(2*)	0.019(3*)	0.032(4*)	0.054(5)	0.084(6)	0.122(7)	Method I	
ThreeRM		0.153(1)	0.109(2)	0.167(3)	0.248(4)	0.341(5)	0.433(6)		
FourRM			0.122(1)	0.219(2)	0.339(3)	0.459(4)	0.567(5)		
FiveRM				0.422(1)	0.614(2)	0.745(3)	0.829(4)		
SixRM					0.565(1)	0.745(2)	0.84(3)		
SevenRM						0.612(1)	0.775(2)		
EightRM							0.615(1)		
TwoRM	0.019(1*)	0.022(2*)	0.019(3*)	0.028(4*)	0.045(5*)	0.077(6)	0.122(7)	Method II	
ThreeRM		0.147(1)	0.109(2)	0.144(3)	0.209(4)	0.316(5)	0.433(6)		
FourRM			0.128(1)	0.191(2)	0.289(3)	0.434(4)	0.577(5)		
FiveRM				0.32(1)	0.487(2)	0.688(3)	0.829(4)		
SixRM					0.503(1)	0.783(2)	0.919(3)		
SevenRM						0.842(1)	0.975(2)		
EightRM							0.918(1)		
TwoRM	0.006(1*)	0.007(2*)	0.012(3*)	0.023(4*)	0.044(5*)	0.077(6)	0.122(7)	Method III	
ThreeRM		0.122(1)	0.185(2)	0.281(3)	0.424(4)	0.567(5)	0.692(6)		
FourRM			0.32(1)	0.487(2)	0.688(3)	0.829(4)	0.914(5)		
FiveRM				0.503(1)	0.783(2)	0.919(3)	0.973(4)		
SixRM					0.842(1)	0.975(2)	0.997(3)		
SevenRM						0.918(1)	0.992(2)		
EightRM							0.943(1)		

TwoRM	0.109(1)	0.076(2)	0.15(3)	0.247(4)	0.362(5)	0.486(6)	0.604(7)	All results	Data set IV
ThreeRM		0.108(1)	0.254(2)	0.416(3)	0.577(4)	0.717(5)	0.823(6)		
FourRM			0.699(1)	0.879(2)	0.96(3)	0.99(4)	0.998(5)		
FiveRM				0.741(1)	0.928(2)	0.985(3)	0.997(4)		
SixRM					0.841(1)	0.98(2)	0.998(3)		
SevenRM						1(1)	1(2)		
EightRM							1(1)		
TwoRM	0.021(1*)	0.021(2*)	0.045(3*)	0.078(4)	0.13(5)	0.17(6)	0.229(7)	Method I	
ThreeRM		0.119(1)	0.254(2)	0.375(3)	0.52(4)	0.581(5)	0.669(6)		
FourRM			0.582(1)	0.714(2)	0.851(3)	0.853(4)	0.898(5)		
FiveRM				0.543(1)	0.783(2)	0.79(3)	0.858(4)		
SixRM					0.729(1)	0.713(2)	0.813(3)		
SevenRM						0.456(1)	0.66(2)		
EightRM							0.6(1)		
TwoRM	0.021(1*)	0.021(2*)	0.034(3*)	0.059(4)	0.1(5)	0.155(6)	0.226(7)	Method II	
ThreeRM		0.119(1)	0.182(2)	0.286(3)	0.413(4)	0.541(5)	0.662(6)		
FourRM			0.324(1)	0.509(2)	0.68(3)	0.804(4)	0.893(5)		
FiveRM				0.539(1)	0.764(2)	0.884(3)	0.952(4)		
SixRM					0.688(1)	0.87(2)	0.956(3)		
SevenRM						0.732(1)	0.923(2)		
EightRM							0.836(1)		
TwoRM	0.032(1*)	0.038(2*)	0.059(3)	0.098(4)	0.154(5)	0.228(6)	0.32(7)	Method III	
ThreeRM		0.167(1)	0.241(2)	0.358(3)	0.488(4)	0.617(5)	0.738(6)		
FourRM			0.334(1)	0.519(2)	0.677(3)	0.803(4)	0.897(5)		
FiveRM				0.539(1)	0.745(2)	0.874(3)	0.952(4)		
SixRM					0.646(1)	0.852(2)	0.956(3)		
SevenRM						0.741(1)	0.947(2)		
EightRM							1(1)		
TwoRM	0.192(1)	0.164(2)	0.28(3)	0.404(4)	0.544(5)	0.671(6)	0.775(7)	All results	Data set V
ThreeRM		0.167(1)	0.345(2)	0.51(3)	0.675(4)	0.801(5)	0.886(6)		
FourRM			0.641(1)	0.819(2)	0.936(3)	0.98(4)	0.994(5)		
FiveRM				0.669(1)	0.903(2)	0.976(3)	0.995(4)		
SixRM					0.886(1)	0.986(2)	0.999(3)		
SevenRM						0.933(1)	0.995(2)		
EightRM							0.96(1)		
TwoRM	0.029(1)	0.027(2*)	0.064(3)	0.107(4)	0.159(5)	0.232(6)	0.323(7)	Method I	
ThreeRM		0.12(1)	0.293(2)	0.42(3)	0.532(4)	0.655(5)	0.768(6)		
FourRM			0.853(1)	0.819(2)	0.864(3)	0.929(4)	0.97(5)		
FiveRM				0.546(1)	0.703(2)	0.841(3)	0.929(4)		
SixRM					0.56(1)	0.79(2)	0.919(3)		
SevenRM						0.717(1)	0.922(2)		
EightRM							0.862(1)		
TwoRM	0.029(1*)	0.027(2*)	0.056(3)	0.097(4)	0.156(5)	0.229(6)	0.313(7)	0	

ThreeRM		0.12(1)	0.252(2)	0.384(3)	0.523(4)	0.648(5)	0.753(6)		Data set VI
FourRM			0.563(1)	0.729(2)	0.851(3)	0.922(4)	0.962(5)		
FiveRM				0.584(1)	0.794(2)	0.901(3)	0.954(4)		
SixRM					0.688(1)	0.868(2)	0.945(3)		
SevenRM						0.727(1)	0.897(2)		
EightRM							0.757(1)		
TwoRM	0.104(1)	0.228(2)	0.356(3)	0.489(4)	0.611(5)	0.739(6)	0.813(7)	Method III	
ThreeRM		0.571(1)	0.738(2)	0.851(3)	0.918(4)	0.97(5)	0.983(6)		
FourRM			0.593(1)	0.79(2)	0.891(3)	0.965(4)	0.98(5)		
FiveRM				0.666(1)	0.844(2)	0.961(3)	0.977(4)		
SixRM					0.696(1)	0.948(2)	0.964(3)		
SevenRM						2(1)	0.938(2)		
EightRM							0.677(1)		
TwoRM	0.269(1)	0.299(2)	0.477(3)	0.635(4)	0.766(5)	0.86(6)	0.921(7)	All results	
ThreeRM		0.275(1)	0.531(2)	0.721(3)	0.853(4)	0.93(5)	0.969(6)		
FourRM			0.782(1)	0.931(2)	0.984(3)	0.997(4)	0.999(5)		
FiveRM				0.797(1)	0.96(2)	0.994(3)	0.999(4)		
SixRM					0.898(1)	0.991(2)	0.999(3)		
SevenRM						0.962(1)	0.999(2)		
EightRM							1(1)		
TwoRM	0.122(1)	0.095(2)	0.08(3)	0.1(4)	0.164(5)	0.228(6)	0.317(7)	Method I	
ThreeRM		0.128(1)	0.113(2)	0.146(3)	0.242(4)	0.332(5)	0.447(6)		
FourRM			0.153(1)	0.216(2)	0.368(3)	0.489(4)	0.628(5)		
FiveRM				0.31(1)	0.572(2)	0.708(3)	0.838(4)		
SixRM					0.77(1)	0.836(2)	0.94(3)		
SevenRM						0.602(1)	0.853(2)		
EightRM							0.831(1)		
TwoRM	0.054(1)	0.049(2*)	0.076(3)	0.107(4)	0.159(5)	0.227(6)	0.314(7)	Method II	
ThreeRM		0.126(1)	0.204(2)	0.271(3)	0.372(4)	0.487(5)	0.607(6)		
FourRM			0.361(1)	0.456(2)	0.591(3)	0.717(4)	0.825(5)		
FiveRM				0.391(1)	0.583(2)	0.736(3)	0.855(4)		
SixRM					0.558(1)	0.766(2)	0.897(3)		
SevenRM						0.662(1)	0.88(2)		
EightRM							0.799(1)		
TwoRM	0.359(1)	0.459(2)	0.59(3)	0.716(4)	0.826(5)	0.902(6)		Method III	
ThreeRM		0.397(1)	0.584(2)	0.738(3)	0.857(4)	0.93(5)			
FourRM			0.55(1)	0.761(2)	0.894(3)	0.96(4)			
FiveRM				0.663(1)	0.882(2)	0.966(3)			
SixRM					0.805(1)	0.962(2)			
SevenRM						0.896(1)			

TwoRM	0.192(1)	0.241(2)	0.407(3)	0.568(4)	0.709(5)	0.816(6)	0.89(7)	All results	Data set VII
ThreeRM		0.284(1)	0.549(2)	0.744(3)	0.871(4)	0.941(5)	0.975(6)		
FourRM			0.82(1)	0.956(2)	0.993(3)	0.999(4)	1(5)		
FiveRM				0.843(1)	0.979(2)	0.997(3)	1(4)		
SixRM					0.955(1)	0.997(2)	1(3)		
SevenRM						0.959(1)	0.999(2)		
EightRM							1(1)		
TwoRM	0.063(1)	0.061(2)	0.059(3)	0.062(4)	0.086(5)	0.14(6)	0.208(7)	Method I	
ThreeRM		0.144(1)	0.136(2)	0.138(3)	0.185(4)	0.287(5)	0.4(6)		
FourRM			0.174(1)	0.185(2)	0.255(3)	0.397(4)	0.539(5)		
FiveRM				0.217(1)	0.331(2)	0.528(3)	0.695(4)		
SixRM					0.408(1)	0.707(2)	0.874(3)		
SevenRM						0.929(1)	0.995(2)		
EightRM							0.96(1)		
TwoRM	0.06(1)	0.049(2*)	0.048(3*)	0.066(4)	0.1(5)	0.139(6)	0.205(7)	Method II	
ThreeRM		0.114(1)	0.111(2)	0.153(3)	0.222(4)	0.292(5)	0.403(6)		
FourRM			0.168(1)	0.25(2)	0.359(3)	0.455(4)	0.595(5)		
FiveRM				0.351(1)	0.518(2)	0.626(3)	0.775(4)		
SixRM					0.505(1)	0.644(2)	0.821(3)		
SevenRM						0.51(1)	0.789(2)		
EightRM							0.842(1)		
TwoRM	0.356(1)	0.518(2)	0.626(3)	0.775(4)	0.878(5)			Method III	
ThreeRM		0.497(1)	0.639(2)	0.817(3)	0.919(4)				
FourRM			0.509(1)	0.79(2)	0.925(3)				
FiveRM				0.848(1)	0.981(2)				
SixRM					0.968(1)				

Note: The bracketed number is the degree of freedom.

Table S3. The comparing between final best model and final optimal models

Data set II		-835.658129(II)	-835.658129(II)	-834.142582(II)		
	-832.129502(III)	0.0079(1*)	0.0079(1*)	0.0448(1*)		
Data set III		-827.589533(III)	-826.570755(IV)	-825.375631(V)	-825.053261(VI)	
	-824.878946(IV)	0.0199(1*)				
		-827.589533(III)	-826.535655(IV)	-825.375631(V)	-824.881462(VI)	-824.657052(VII)
	-824.878946(IV)	0.0199(1*)				
		-826.570755(III)	-825.375631(IV)	-824.881462(V)	-824.657052(VI)	-824.637086(VII)
	-824.878946(IV)					
Data set IV		-827.516759(III)	-826.299258(IV)	-826.147571(V)		
	-828.19381(II)					
		-827.516759(III)	-826.299258(IV)	-825.813171(V)		
	-828.19381(II)					
		-827.236751(III)	-826.280292(IV)			
	-828.19381(II)					
Data set V		-833.281509(III)	-832.07163(IV)	-833.281509(III)	-832.07163(IV)	-834.184208(II)
	-833.577481(II)					
Data set VI		-829.704697(II)	-826.682121(IV)	-827.292062(II)		
	-826.875638(II)	0.0174(1*)				
Data set VII		-841.727813(II)	-839.963719(III)	-838.716878(IV)	-838.341428(II)	
	-838.341428(II)	0.0093(1*)				

The green color means final best model of each data set, the red color means the final optimal model of method I, and the orange color means the final optimal model of method II while the blue color means the final optimal model of method III. When final best model is significant better than final optimal models, it's marked by asterisk

Table S4. The comparing between final optimal models and hypothesis models

Data sets	lnL value of Hypothesis models	lnL value of optimal models by three methods				
Data set II		-835.658129(2-RM)	-835.658129(2-RM)	-834.142582(2-RM)		
	-835.868563(2-RM) -835.867062(3-RM)					
Data set III		-827.589533(3-RM)	-826.570755(4-RM)	-825.375631(5-RM)	-825.053261(6-RM)	
	-829.92476(2-RM) -828.725264(3-RM)	0.0307(1*) /	0.0349(2*) 0.0379(1*)	0.028(3*) 0.0351(2*)	0.045(4*) 0.0617(3)	
Data set IV		-827.589533(3-RM)	-826.535655(4-RM)	-825.375631(5-RM)	-824.881462(6-RM)	-824.657052(7-RM)
	-829.92476(2-RM) -828.725264(3-RM)	0.0307(1*) /	0.0349(2*) 0.0364(1*)	0.028(3*) 0.0351(2*)	0.039(4*) 0.0529(3)	0.0614(5) 0.0867(4)
		-826.570755(3-RM)	-825.375631(4-RM)	-824.881462(5-RM)	-824.657052(6-RM)	-824.637086(7-RM)
	-829.92476(2-RM) -828.725264(3-RM)	0.0096(1*) /	0.0106(2*) 0.0096(1*)	0.0178(3*) 0.0351(2*)	0.0323(4*) 0.0433(3*)	0.0605(5) 0.0853(4)
		-827.516759(3-RM)	-826.299258(4-RM)	-826.147571(5-RM)		
	-829.470066(2-RM) -828.554522(3-RM)	0.0481(1*) /	0.042(2*) 0.0337(1*)	0.0841(3) 0.0901(2)		
Data set V		-827.516759(3-RM)	-826.299258(4-RM)	-825.813171(5-RM)		
	-829.470066(2-RM) -828.554522(3-RM)	0.0481(1*) /	0.042(2*) 0.0337(1*)	0.0625(3) 0.0645(2)		
		-827.236751(3-RM)	-826.280292(4-RM)			
	-829.470066(2-RM) -828.554522(3-RM)	0.0346(1*) /	0.0412(2*) 0.0329(1*)			
		-833.281509(3-RM)	-832.07163(4-RM)	-833.281509(3-RM)	-832.07163(4-RM)	-834.184208(2-RM)
	-834.184208(2-RM) -833.480516(3-RM)					
Data set VI		-829.704697(2-RM)	-826.682121(4-RM)	-827.292062(2-RM)		

	-828.354787(2-RM) -827.855504(3-RM)					
Data set VII		-841.727813(2-RM)	-839.963719(3-RM)	-838.716878(4-RM)	-838.341428(2-RM)	
	-840.41675 (2-RM) -839.963719(3-RM)					

The green color means hypothesis models, the red color means the models of method I, orange the models of method II while blue the models of method III. Only the significant better models' p-value is shown.

Table S5 Comparing Final Optimal Model of OBSM with Hypothesis model suggested in 50 cases

Cases	Hypothesis model	Final optimal model	2Δl	df	LRTs <i>P</i> -value
1 (12)	-4106.21	-4097.60	17.23	3	6.300e-004
2 (13)	-368.96	-348.67	40.58	19	2.700e-003
3 (14, 15)	-21407.46	-21304.34	206.25	8	0
4 (14, 15)	-29760.57	-29704.59	111.97	9	0
5 (15, 16) ^A	-9341.07	-9337.50	7.13	2	2.830e-002
6 (17)	-6426.88	-6327.43	198.91	11	0
7 (18, 19)	-4086.39	-4072.63	27.52	4	1.560e-005
8 (20)	-17761.37	-17647.19	228.34	27	1.98e-014
9 (21)	-4396.93	-4366.65	60.57	12	1.780e-008
10 (22)	-2460.12	-2446.34	27.56	7	2.600e-004
11 (23)	-2657.21	-2649.35	15.72	3	1.300e-003
12 (24)	-1563.67	-1559.60	8.13	3	4.340e-002
13 (25)	-21441.24	-21353.38	175.71	16	0
14 (26-29)	-1397.58	-1387.40	20.37	4	4.200e-004
15 (30)	-2557.24	-2542.53	29.43	8	2.670e-004
16 (31)	-3019.79	-2998.27	43.04	9	2.120e-006
17 (32)	-2604.31	-2587.16	34.30	6	5.900e-006
18 (33)	-2621.56	-2600.85	41.42	9	4.200e-006
19 (34)	-2800.78	-2774.63	52.31	9	3.96e-008
20 (35)	-2043.71	-2033.04	22.34	4	1.720e-004
21 (36)	-2518.06	-2503.41	29.30	4	6.807e-006*
22 (37)	-4664.00	-4649.42	29.15	9	6.108e-004
23 (38)	-1390.02	-1386.12	7.81	2	0.0201
24 (39)	-5374.34	-5307.31	134.06	22	0
25 (39)	-2277.11	-2264.78	24.66	2	4.411e-006
26 (40)	-2451.21	-2413.22	75.98	12	2.396e-011
27 (41)	-1187.19	-1184.11	6.18	1	1.294e-002
28 (42)	-1597.54	-1585.08	24.93	2	3.861e-006
29 (43)	-1439.86	-1431.59	16.54	3	8.774e-004
30 (43)	-1345.23	-1340.19	10.08	2	6.484e-003
31 (43)	-2961.58	-2943.89	35.37	11	2.149e-004
32 (43) ^A	-1570.58	-1567.76	5.65	1	0.0175
33 (43)	-2353.63	-2342.41	22.43	4	1.643e-004
34 (43)	-1768.02	-1738.75	58.53	8	9.039e-010
35 (43)	-6917.24	-6879.80	74.89	11	1.424e-011
36 (43)	-4245.82	-4222.10	47.44	14	1.629e-005
37 (43)	-3790.66	-3774.20	32.92	8	6.380e-005
38 (43) ^A	-5350.87	-5349.17	3.38	1	6.582e-002
39 (43) ^A	-2434.20	-2430.55	7.29	1	6.948e-003
40 (43)	-2332.73	-2332.73			/

41 (44) ^A	-1310.46	-1306.08	8.76	1	3.076e-003
42 (44)	-1692.47	-1675.31	34.32	6	5.835e-006
43 (45) ^A	-2170.44	-2167.84	5.22	1	2.234e-002
44 (46)	-1374.18	-1364.61	19.14	5	1.810e-003
45 (47)	-1432.48	-1432.48			/
46 (47) ^A	-315.41	-311.94	11.02	1	8.428e-003
47 (48)	-1493.08	-1473.71	38.72	7	2.210e-006
48 (49) ^A	-3214.20	-3209.92	8.56	2	1.387e-002
49 (50)	-11734.74	-11703.17	63.14	8	1.126e-010
50 (51)	-9963.96	-9920.25	87.43	13	4.321e-013

^AThe log likelihood of final optimal models of these cases is obtained from OBSM Method III.

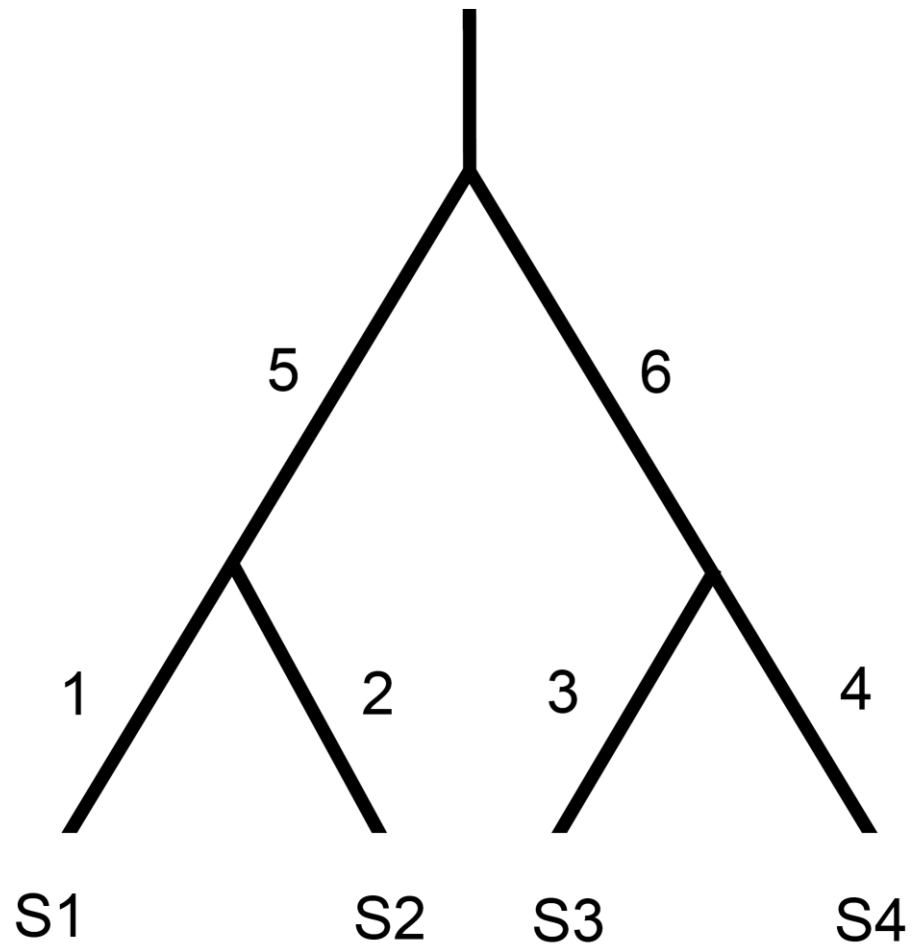


Figure S1. An example of a phylogeny with 4 species (or genes)

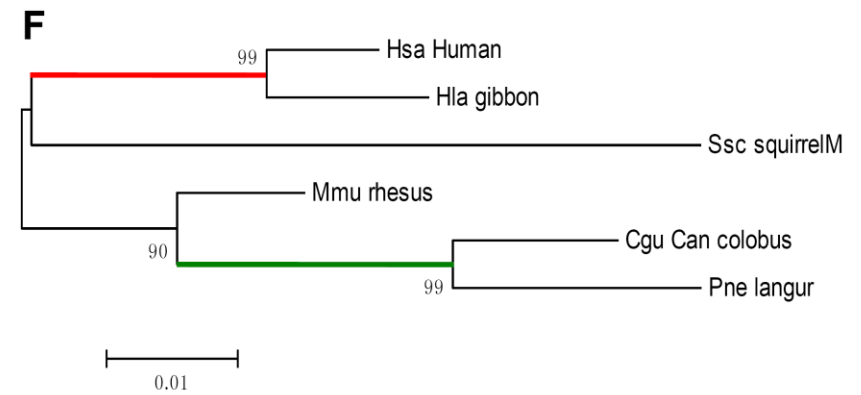
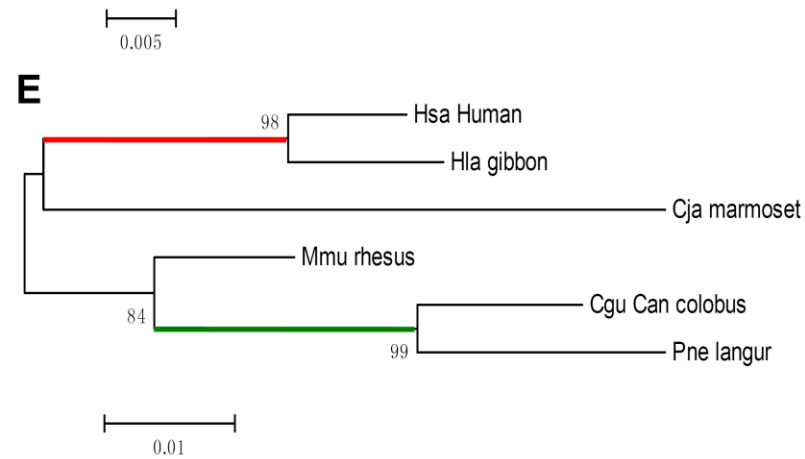
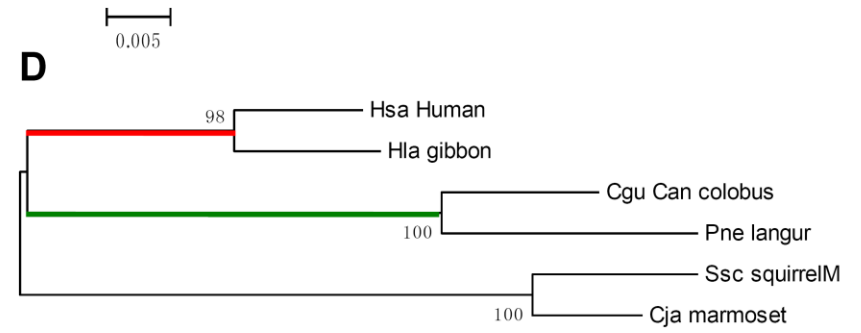
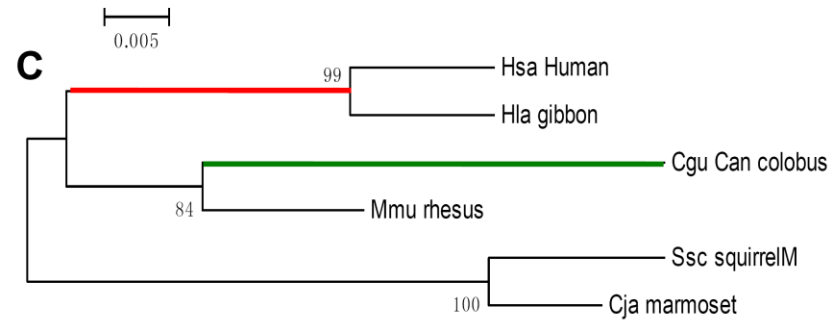
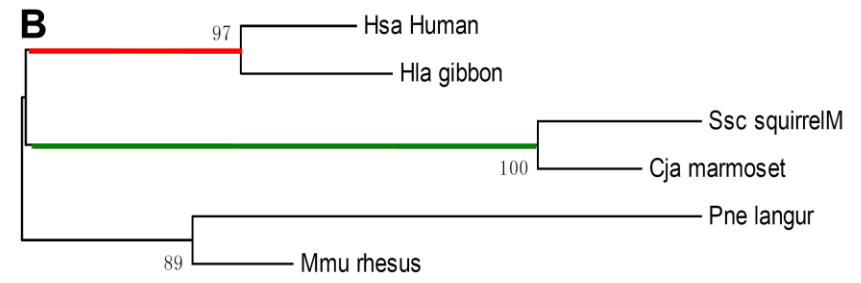
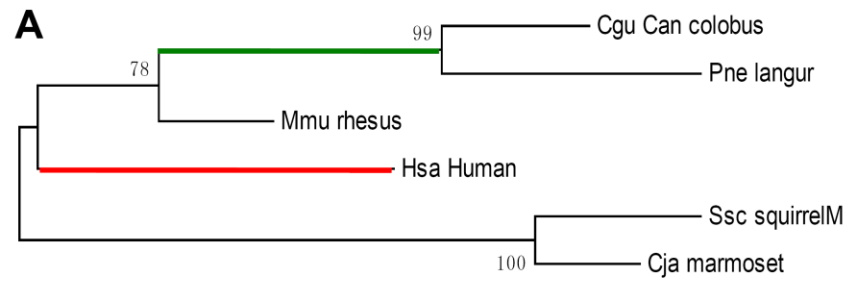


Figure S2. The phylogenies generated from 6 datasets.

Perl scripts

Perl Script 1 Start##

```
#!/usr/bin/perl
```

```
#
```

```
#after runing, two temp files will be there, and don't delete them, one temp file is needed  
by supplementary file 2.pl.
```

```
print "Please input the sequences numbers:\n";  
my $sequences_numbers=<STDIN>;  
our $branches_numbers=2*($sequences_numbers-1);  
my $temp="";  
my $j;  
my $i;  
my $temp_branches_numbers;  
my $branches=$branches_numbers;  
my @temp;  
my @temp2;  
my $flag;  
my $temp_arr;  
my %hash;  
my $temp_hash;  
my $down_number;  
my $blank_temp;  
my $temp_key;  
my $total_rel;  
my @arr;
```

```
open (TEMP, ">temp.pl");  
&maina;  
close TEMP;
```

```
#####  
#####
```

```
#first part
```

```
#####  
#####
```

```
while(&check_temp_pl) {}  
while(&check_temp_txt) {}
```

```
&get_arr;  
my $rel_file="factorization_of_". $branches_numbers. "_branch.txt";  
open (DECOMPOSITION, ">$rel_file");  
&main_2;  
close DECOMPOSITION;
```

```
sub get_arr  
{  
    undef @arr;
```

```

open (FILE, "<temp.txt");
while(<FILE>) {
    chomp;
    while($_ =~ s/\s//) {}
    next if $_ =~ /\#/;
    $arr[$i++]=$_;
}
close FILE;
}
sub check_temp_pl
{
    if (-e "temp.pl") {
        my $tail=`grep end_maker temp.pl`;
        chomp $tail;
        if ($tail eq "#end_maker#") {
            system "perl temp.pl >temp.txt";
            return 0;
        } else {
            print "noend";
            sleep (2);
            return 1;
        }
    } else {
        print "nopl";
        sleep (2);
        return 1;
    }
}
sub check_temp_txt
{
    if (-e "temp.txt") {
        my $tail=`grep decompose_maker temp.txt`;
        chomp $tail;
        if ($tail eq "#decompose_maker#") {
            return 0;
        } else {
            sleep (2);
            return 1;
        }
    } else {
        sleep (2);
        return 1;
    }
}

sub main_2
{
    foreach $temp_key (@arr) {
        chomp $temp_key;
    }
}

```

```

next if $temp_key eq "";
($blank_temp,@temp)=split/\_/, $temp_key;
undef @temp2;
$flag="F";
foreach $temp_arr (sort {$a <=> $b} @temp) {
    next if $temp_arr eq 1;
    @temp2=(@temp2, $temp_arr);
}
undef %hash;
my $k;
my $i;
for ($i=0;$i<=$#temp2;$i++) {
    next if $temp2[$i] eq "";
    $k=$i+1 if $i<$#temp2;
    if ($temp2[$i] eq $temp2[$k]) {
        $flag="T";
    }
    $temp_hash="hash"."_" . $temp2[$i];
    $hash{$temp_hash}++;
}
$down_number=$branches;
my $all_rel=1;
my $rest_rel=0;
my ($diuqi, $value);
my $blank="";
if ($flag eq "T") {
    print DECOMPOSITION join("_",@temp), "\t";
    my $key;
    foreach $key (sort keys %hash) {
        if ($hash{$key}>1) {
            ($diuqi, $value)=split/\_/, $key;
            for ($valuesi=0;$valuesi<$hash{$key};$valuesi++) {
                $rel=&combination($down_number, $value);
                $all_rel=$all_rel*$rel;
                $down_number=$down_number-$value;
            }
            $all_rel=$all_rel/(&factorial(1, $hash{$key}));
        } else {
            ($diuqi, $value)=split/\_/, $key;
            $rel=&combination($down_number, $value);
            $all_rel=$all_rel*$rel;
            $down_number=$down_number-$value;
        }
    }
    print DECOMPOSITION " ", $all_rel, "\n";
    $total_rel=$total_rel+$all_rel;
} else {
    my $kk;
    for (my $i=0;$i<$#temp;$i++) {

```

```

        next if $temp[$i] == 0;
        next if $temp[$i] eq "";
        next if $temp[$i] eq 1;
        $kk=$i+1 if $i<=$#temp;
        $rel=&combination($down_number,$temp[$i]);
        $all_rel=$all_rel*$rel;
        $rest_rel=$rest_rel+$temp[$i];
        $down_number=$branches-$rest_rel;
        last if $down_number eq $temp[$kk];
    }
    print DECOMPOSITION join("_",@temp),"\t$total_rel\n";
    $total_rel=$total_rel+$all_rel;
}
}
$total_rel=$total_rel+2;
print DECOMPOSITION $total_rel;
}

sub print_pre
{
    ($j)=@_;
    for ($i=0;$i<$j;$i++){
        $k=($i-1);
        if ($k eq -1){
            print TEMP $blank."for (\$arr_$i=0;\$arr_$i<$branches;\$arr_$i++) {\n";
        }else{
            print TEMP $blank."for
(\$arr_$i=\$arr_$k;\$arr_$i<$branches;\$arr_$i++) {\n";
        }
        $blank=$blank."\t";
    }
    for ($i=0;$i<$j;$i++){
        print TEMP "\$arr_$i";
    }
    print TEMP "\n";
    for ($i=0;$i<$j;$i++){
        $blank=~s/\t//;
        print TEMP $blank."}\n";
    }
}

sub combination
{
    my ($down_number,$up_number)=@_;
    my $numerator=&factorial(1,$down_number);
    my $denominator_1=&factorial(1,$up_number);
    my $denominator_2=&factorial(1,($down_number-$up_number));
    my $combination=$numerator/($denominator_1*$denominator_2);
    return $combination;
}

```

```

}

sub factorial()
{
    my ($f_rel, $n)=@_;
    $f_rel=$f_rel*$n;
    return $f_rel=1 if $n eq 0;
    if ($n>1)
    {
        $n--;
        &factorial($f_rel, $n)    ;
    }else{
        return $f_rel;
    }
}

sub maina
{
    for ($i=2;$i<$branches_numbers;$i++){
        print TEMP <<END;
        foreach \%key (sort keys \%hash){
            print "\$key\\n";
        }
        END
        print TEMP "undef \%hash;\\n";
        print TEMP "##### $i ratio models #####\\n";
        $temp_branches_numbers=$branches_numbers;
        for ($j=1;$j<($i-1);$j++){
            &circle_pre($j, $i, ($j-1));
            $temp=$temp."\\t";
        }
        &print_sub($j, $i, ($j-1));
        for ($j=1;$j<=($i-1);$j++){
            &circle_end;
            $temp=~s/\\t//;
        }
    }
    &print_paixu;
}

sub print_paixu
{
    print TEMP<<END;
    sub paixu
    {
        \@array=@_;
        my \%ttemp;
        foreach (sort { \%a <=> \%b } \@array){
            \%ttemp=\%ttemp."_\$_" ;
        }
    }
}

```

```

    }
    \${hash}{\${ttemp}}++;
}

foreach \${key} (sort keys \${hash}) {
    print "\${key}\n";
}
print "#decompose_maker#";
#end_maker#
END
}

sub circle_pre
{
    my (\$k, \$dd, \$k_1)=@_;
    if (\$k_1 eq 0) {
        print TEMP \$temp, "for
(\${arr}_\$k=1;\${arr}_\$k<". (int(\$temp_branches_numbers/(\$dd-\$k+1))+1). "; \${arr}_\$k++) {\n";
    } else {
        print TEMP \$temp, "for
(\${arr}_\$k=\${arr}_\$k_1;\${arr}_\$k<". (int(\$temp_branches_numbers/(\$dd-\$k+1))+1). "; \${arr}_\$k+
+)\ {\n";
    }
    \$temp_branches_numbers=\$temp_branches_numbers-1;
}

sub circle_end{
    print TEMP \$temp, "}\n";
}

sub print_sub
{
    my (\$k, \$dd, \$k_1)=@_;
    if (\$k_1 eq 0) {
        print TEMP \$temp, "for
(\${arr}_\$k=1;\${arr}_\$k<". (int(\$temp_branches_numbers/(\$dd-\$k+1))+1). "; \${arr}_\$k++) {\n";
    } else {
        print TEMP \$temp, "for
(\${arr}_\$k=\${arr}_\$k_1;\${arr}_\$k<". (int(\$temp_branches_numbers/(\$dd-\$k+1))+1). "; \${arr}_\$k+
+)\ {\n";
    }
    \$temp_branches_numbers=\$temp_branches_numbers-1;
    print TEMP \$temp, "\t\${temp}=\$branches_numbers";
    for (my \$kk=1; \$kk<=\$k; \$kk++) {
        print TEMP "-\${arr}_\$kk";
    }
    print TEMP ";\n";
    print TEMP \$temp, "\t", qq|\&paixu(|;
    for (\$kk=1; \$kk<=\$k; \$kk++) {

```



```

        print TEMP "\$arr_$kk,";
    }
    print TEMP qq|\$temp) if \$temp>0;\n|;
}

## Perl Script 1 End##

## Perl Script 2 Start##

#!/usr/bin/perl -w
#should run after supplementary file 1.pl
#
#two parameters are needed:
#
#first, the branch numbers (it should be the same with the input of supplementary file 1.pl)
#second, the phylogeny file name
#
#after runing, two files all_labled_trees.txt and all_models.txt will be there.
#
#Chengjun Zhang 2009-10-21
#

my $kk;
my @arr;
print "Please input the sequences numbers:\n";
my $sequences_numbers=<STDIN>;
my $branches_numbers=2*($sequences_numbers-1);
my $branches=$branches_numbers;
my $file="temp.txt";
chomp $file;
open (TEMP, ">temp.pl");
print TEMP "&get_arr;\n";
open (FILE, "$file");
while(<FILE>)
{
    last if $_=~/\#/;
    chomp;
    my $tempp;
    ($tempp, @line)=split/\_/, $_;
    $line=$#line+1;
    print TEMP "#####$line Ratio Model
@line#####\n";
    $last_round=$#line-1;
    &get_circles($line[0], $last_round, 1, "");
}
close FILE;

sub get_circles
{

```

```

my ($circles,$last_round,$index_number,$blank)=@_;
my $i;
my $k;
$k=1;
for ($i=0;$i<$circles;$i++){
    print TEMP $blank,"\$l=\$i_$index_number". "_$k+1;\n" if $i ne 0;
    if ($i eq 0){
        print TEMP $blank,"for
(\$i_$index_number". "_$i=1;\$i_$index_number". "_$i<=$branches;\$i_$index_number". "_$i+
+)\n";
    }else{
        print TEMP $blank,"for
(\$i_$index_number". "_$i=\$l;\$i_$index_number". "_$i<=$branches;\$i_$index_number". "_$
i++)\n";
    }
    $k=$i;
    $blank=$blank. "\t";
}

if ($index_number > 1){
    $temp_number=$index_number-1;
    print TEMP $blank,"undef \@arr_new_$index_number;\n";
    print TEMP $blank,"\@arr_new_$index_number=\@arr_new_$temp_number;\n";
}

if ($index_number eq 1){
    print TEMP $blank,"my \@arr_temp=\@arr;\n";
}

for ($i=0;$i<$circles;$i++){
    print TEMP qq($blank).qq(next
if ).qq(\$arr_new_$index_number).qq([\$i_$index_number).qq(_$k]).qq( eq "";\n) if
$index_number>1;
    $k=$i;
}

print TEMP qq($blank).qq(my \@arr_$index_number;\n);
print TEMP qq($blank).qq(undef \@arr_$index_number;\n);
for ($i=0;$i<$circles;$i++){
    print TEMP
qq($blank).qq{\@arr_$index_number=(\@arr_$index_number,\$arr_temp[\$i_$index_number}.q
q{_$k]);\n};
    $k=$i;
}

if ($last_round eq 0){
    print TEMP $blank,"print \"\"";
    for ($i=1;$i<=$index_number;$i++){
        print TEMP "\@arr_$i*";
    }
}

```

```

    }
    print TEMP "\\n\\n";\n";
}
for ($i=0;$i<$circles;$i++){
    if ($last_round>0){
        if ($index_number eq 1){
            print TEMP $blank,"delete \${arr_temp}[\${i}_$index_number"]."_$k";\n";
        }else{
            print TEMP $blank,"delete
\${arr_new}_$index_number". "[\${i}."_$index_number"]. "_$k";\n";
        }
    }
    $k=$i;
}
if ($index_number eq 1){
    print TEMP $blank,"undef \@arr_new_1;\n";
    print TEMP $blank,"\@arr_new_1=\@arr_temp;\n";
}else{
    $temp_number=$index_number+1;
    print TEMP $blank,"undef \@arr_new_$temp_number;\n";
    print TEMP $blank,"\@arr_new_$temp_number=\@arr_new_$index_number;\n";
}
if ($last_round>0){
    $last_round=$last_round-1;
    $index_number=$index_number+1;
    &get_circles($line[$index_number-1], $last_round, $index_number, $blank);
}
for ($i=0;$i<$circles;$i++){
    $blank=~s/\t//;
    print TEMP $blank,"}\n";
}
}

print TEMP <<END;
print "#decompose_maker#";
sub get_arr
{
    my \${i};
    for (\${i}=1;\${i}<=$branches;\${i}++) {
        \${arr}[\${i}]=\${i};
    }
}
END
print TEMP "#end_maker#\n";
close TEMP;
while(&check_temp_pl){}
system "perl temp.pl >templ.txt";
while(&check_temp1_txt){}
my $all_branch=&all_branch;

```

```

my $temp_five;
open (CORRECT,">all_models.txt");
open (TREES,">all_labled_trees.txt");
open (FILE,"<templ.txt");
print "Please input the tree file name(make sure in one line):\n";
my $tree_file=<STDIN>;
my $flag=0;
while(<FILE>){
    last if $_~/\#/;
    chomp;
    @line=split/\*/, $_;
    undef %hash if $flag ne $#line;
    $flag=$#line;
    $k=0;
    undef @new_arr;
    my $original_tree=&check_original_tree if $tree_file;
    my ($edit_tree,$edit_tree_lines)=&get_edit_tree($original_tree);
    $temp_five=$all_branch;
    foreach $ratio (@line){
        @members=split/ /,$ratio;
        @members=sort {$a <=> $b} @members;
        $new_arr[$k]=join(" ",@members);
        $k++;
        $maker="#". $k;
        $edit_tree=&lable_tree($maker,$edit_tree,@members);
    }
    my ($lost,@remain)=split/\_/, $temp_five;
    $new_arr[$#new_arr+1]=join(" ",@remain) if $remain[0] ne "";
    @new_arr=sort {$a cmp $b} @new_arr;
    $hash_pre=join(" ",@new_arr);
    $hash{$hash_pre}++;
    &one_line_tree($edit_tree) if $hash{$hash_pre} eq 1;
    next if substr($hash_pre,0,1) eq "*" or substr($hash_pre,0,1) eq " ";
    print CORRECT $hash_pre,"\n" if $hash{$hash_pre} eq 1;
}
close FILE;
close CORRECT;
close TREES;

sub one_line_tree
{
    my ($tree)=@_;
    my @tree=split/\n/, $tree;
    $tree=join("",@tree);
    print TREES $tree,"\n";
}

sub lable_tree
{

```

```

my ($maker, $edit_tree, @members)=@_;
my @tree=split/\n/, $edit_tree;
foreach (@members) {
    $del="_".$_;
    $temp_five=~s/$del//;
    $_=$_-1;
    my $last_letter=substr($tree[$_],-1,1);
    if ($last_letter eq ",") {
        $tree[$_]=~s/,/$maker,/;
    }
    if ($last_letter eq ")") {
        $tree[$_]=~s/\)/$maker\)/;
    }
}

$tree=join("\n", @tree);
return $tree;
}

sub check_original_tree
{
    open(TREE, "<$tree_file") or die "can't open the tree file";
    my $tree="";
    while(<TREE>)
    {
        $tree=$tree.$_;
    }
    close TREE;
    return $tree;
}

sub get_edit_tree
{
    my ($tree)=@_;
    my ($i,$j);
    while($tree=~s/,/+\\n/) {$i++;}
    while($tree=~s/\)/=\\n/) {$j++;}
    while($tree=~s/\\+/,/) {}
    while($tree=~s/\\=/)/) {}
    $i=$i+$j;
    return ($tree,$i);
}

sub check_temp_pl
{
    if (-e "temp.pl") {
        my $tail=`grep end_maker temp.pl`;
        chomp $tail;
        if ($tail eq "#end_maker#") {
            system "perl temp.pl >templ.txt";
        }
    }
}

```

```

        return 0;
    }else{
        print "noend";
        sleep (2);
        return 1;
    }
}else{
    print "nopl";
    sleep (2);
    return 1;
}
}

```

```

sub all_branch
{
    my $string="";
    for (my $i=1;$i<=$branches;$i++){
        $string=$string."_".$i;
    }
    return $string;
}

```

```

sub check_tmpl_txt
{
    if (-e "templ.txt"){
        my $tail=`grep decompose_maker templ.txt`;
        chomp $tail;
        if ($tail eq "#decompose_maker#"){
            return 0;
        }else{
            sleep (2);
            return 1;
        }
    }else{
        sleep (2);
        return 1;
    }
}

```

Perl Script 2 End##

Perl Script 3 Start##

```

#!/usr/bin/perl -w
#2009-8-30 edit from YangZiheng'S PAML chi2.c
#This script is to do a likelihood ratio test based on the result of Method II.
#

```

```

print "please input the result file name of method II\n";
my $file_name=<STDIN>;
chomp $file_name;

my %hash;

open (REL, "<$file_name");
while(<REL>) {
    chomp;
    @line=split /\t/, $_;
    last if $#line <3;
    $hash{$line[1]}++;
    $temp_id=$hash{$line[1]};
    $id="1".$line[1];
    $hash{$id} {$temp_id}=$line[3];
}
for ($i=2;$i<=($line[1]>9?9:$line[1]);$i++) {
    #the number is determined to 9 here just because nine ratio is enough for normal
    analysis.
    $max=-10000;
    for ($j=1;$j<=$hash{$i};$j++) {
        $k="1".$i;
        $max=$hash{$k} {$j} if $max<$hash{$k} {$j};
    }
    @p_value=(@p_value, $max);
}
print join("\t", @p_value);
print "\n";

@input=@p_value;
for ($i=0;$i< $#input;$i++) {
    $k=$i+1;
    for ($l=0;$l<$i;$l++) {
        print "\t";
    }
    for ($j=$k;$j< $#input;$j++) {
        $df=$j-$i;
        $chi2=2*($input[$j]-$input[$i]);
        $prob=1-&CDFChi2($chi2, $df);
        $prob=int($prob*1000+0.5)/1000;
        if ($prob<0.05) {
            print "$prob($df *)\t";
        }else{
            print "$prob($df)\t";
        }
    }
    print "\n";
}

```

```

sub CDFChi2
{
    ($x, $v)=@_;
    return (&CDFGamma($x, ($v)/2.0, 0.5));
}

sub CDFGamma
{
    ($x, $alpha, $beta)=@_;
    return (&IncompleteGamma(($beta)*($x), $alpha, &LnGammaFunction($alpha)));
}

sub IncompleteGamma
{
    ($x, $alpha, $ln_gamma_alpha)=@_;
    my $i;
    $p=$alpha;
    $g=$ln_gamma_alpha;
    $accurate=1e-8;
    $overflow=1e30;
    my $factor;
    $gin=0;
    $rn=0;
    $a=0;
    $b=0;
    $an=0;
    $dif=0;
    $term=0;
    # $pn[6];

    return (0) if ($x==0) ;
    return (-1) if ($x<0 or $p<=0);

    $factor=exp($p*log($x)-$x-$g);
    goto 130 if ($x>1 and $x>=$p) ;
#    /* (1) series expansion */
    $gin=1; $term=1; $rn=$p;
120:
    $rn++;
    $term*=$x/$rn;
    $gin+=$term;
    goto 120 if ($term > $accurate) ;
    $gin*=$factor/$p;
    goto 150;
130:
#    /* (2) continued fraction */
    $a=1-$p;
    $b=$a+$x+1;

```



```

    $term=0;
    $pn[0]=1;
    $pn[1]=$x;
    $pn[2]=$x+1;
    $pn[3]=$x*$b;
    $gin=$pn[2]/$pn[3];
132:
    $a++;
    $b+=2;
    $term++;
    $an=$a*$term;
    for ($i=0; $i<2; $i++) {
        $pn[$i+4]=$b*$pn[$i+2]-$an*$pn[$i];
    }
    goto 135 if ($pn[5] == 0) ;
    $rn=$pn[4]/$pn[5];
    $dif=abs($gin-$rn);
    goto 134 if ($dif>$accurate);
    goto 142 if ($dif<=$accurate*$rn) ;
134:
    $gin=$rn;
135:
    for ($i=0; $i<4; $i++){
        $pn[$i]=$pn[$i+2];
    }
    goto 132 if (abs($pn[4]) < $overflow);
    for ($i=0; $i<4; $i++){
        $pn[$i]/=$overflow;
    }
    goto 132;
142:
    $gin=1-$factor*$gin;

150:
    return ($gin);
}
sub LnGammaFunction
{
    ($alpha)=@_;
    $x=$alpha;
    $f=0;
    my $z;

    if ($x<7) {
        $f=1;
        $z=$x-1;
        while(++$z<7){
            $f*=$z;
        }
    }

```

```

        $x=$z;
        $f=-log($f);
    }
    $z = 1/($x*$x);
    return $f + ($x-0.5)*log($x) - $x + .918938533204673+
    (((-0.000595238095238*$z+.000793650793651)*$z-.002777777777778)*$z+.083333333333333)/$x
;
}

```

Perl Script 3 End##

Perl Script 4 Start##

```

#!/usr/bin/perl -w
#2009-8-30 edit from YangZiheng'S PAML chi2.c
#This script is to do a likelihood ratio test based on the result of Method III.
#

```

```

print "please input the result file name of method III\n";
my $file_name=<STDIN>;
chomp $file_name;

```

```

my %hash;
$id_max=1;
$last_ratio=2;
$flag=1;
open (REL,"<$file_name");
while(<REL>){
    chomp;
    if ($_ =~ /www/){
        $flag++ ;
        next;
    }
    @line=split /\t/, $_ if $_ ! =~ /www/;
    last if $#line < 3 and $_ ! =~ /www/;
    if ($last_ratio ne $line[1]){
        $id_max=$id_max-$flag;
        $ratio_id="1".$last_ratio;
        for ($i=0;$i<$id_max;$i++){
            $temp_id=$hash{$last_ratio}-$i;
            #print "$temp_id\t";
            #print "$ratio_id\twhy\t";
            #print "$hash{$ratio_id} {$temp_id}\n";
            delete $hash{$ratio_id} {$temp_id};
        }
        $hash{$last_ratio}--=$id_max;
        $flag=1;
        $last_ratio = $line[1];
    }
}

```

```

    $hash{$line[1]}++;
    $temp_id=$hash{$line[1]};
    $id="1".$line[1];
    $hash{$id} {$temp_id}=$line[3];
    $id_max=$line[0] if $line[0]>$id_max;
}
for ($i=2;$i<=($line[1]>9?9:$line[1]);$i++) {
    $max=-10000;
    for ($j=1;$j<=$hash{$i};$j++) {
        $k="1".$i;
        $max=$hash{$k} {$j} if $max<$hash{$k} {$j};
    }
    @p_value=(@p_value,$max);
}
print join("\t",@p_value);
print "\n";

```

```

@input=@p_value;
for ($i=0;$i< $#input;$i++) {
    $k=$i+1;
    for ($l=0;$l<$i;$l++) {
        print "\t";
    }
    for ($j=$k;$j< $#input;$j++) {
        $df=$j-$i;
        $chi2=2*($input[$j]-$input[$i]);
        $prob=1-&CDFChi2($chi2,$df);
        $prob=int($prob*1000+0.5)/1000;
        if ($prob<0.05) {
            print "$prob($df *)\t";
        } else {
            print "$prob($df)\t";
        }
    }
    print "\n";
}

```

```

sub CDFChi2
{
    ($x,$v)=@_;
    return (&CDFGamma($x,($v)/2.0,0.5));
}

```

```

sub CDFGamma
{
    ($x,$alpha,$beta)=@_;
    return (&IncompleteGamma(($beta)*($x),$alpha,&LnGammaFunction($alpha)));
}

```

```

sub IncompleteGamma
{
    ($x, $alpha, $ln_gamma_alpha)=@_;
    my $i;
    $p=$alpha;
    $g=$ln_gamma_alpha;
    $accurate=1e-8;
    $overflow=1e30;
    my $factor;
    $gin=0;
    $rn=0;
    $a=0;
    $b=0;
    $an=0;
    $dif=0;
    $term=0;
    # $pn[6];

    return (0) if ($x==0) ;
    return (-1) if ($x<0 or $p<=0);

    $factor=exp($p*log($x)-$x-$g);
    goto 130 if ($x>1 and $x>=$p) ;
#    /* (1) series expansion */
    $gin=1; $term=1; $rn=$p;
120:
    $rn++;
    $term*=$x/$rn;
    $gin+=$term;
    goto 120 if ($term > $accurate) ;
    $gin*=$factor/$p;
    goto 150;
130:
#    /* (2) continued fraction */
    $a=1-$p;
    $b=$a+$x+1;
    $term=0;
    $pn[0]=1;
    $pn[1]=$x;
    $pn[2]=$x+1;
    $pn[3]=$x*$b;
    $gin=$pn[2]/$pn[3];
132:
    $a++;
    $b+=2;
    $term++;
    $an=$a*$term;
    for ($i=0; $i<2; $i++) {
        $pn[$i+4]=$b*$pn[$i+2]-$an*$pn[$i];
    }
}

```

```

    }
    goto 135 if ($pn[5] == 0) ;
    $rn=$pn[4]/$pn[5];
    $dif=abs($gin-$rn);
    goto 134 if ($dif>$accurate);
    goto 142 if ($dif<=$accurate*$rn) ;
134:
    $gin=$rn;
135:
    for ($i=0; $i<4; $i++){
        $pn[$i]=$pn[$i+2];
    }
    goto 132 if (abs($pn[4]) < $overflow);
    for ($i=0; $i<4; $i++){
        $pn[$i]/=$overflow;
    }
    goto 132;
142:
    $gin=1-$factor*$gin;

150:
    return ($gin);
}
sub LnGammaFunction
{
    ($alpha)=@_;
    $x=$alpha;
    $f=0;
    my $z;

    if ($x<7) {
        $f=1;
        $z=$x-1;
        while(++$z<7) {
            $f*=$z;
        }
        $x=$z;
        $f=-log($f);
    }
    $z = 1/($x*$x);
    return $f + ($x-0.5)*log($x) - $x + .918938533204673+
(((-.000595238095238*$z+.000793650793651)*$z-.0027777777777778)*$z+.083333333333333)/$x
;
}

```

Perl Script 4 End##

The example configuration of six sequences

two ratio configuration

1_9 *10*

2_8 *45*

3_7 *120*

4_6 *210*

5_5 *126*

three ratio configuration

1_1_8 *45*

1_2_7 *360*

1_3_6 *840*

1_4_5 *1260*

2_2_6 *630*

2_3_5 *2520*

2_4_4 *1575*

3_3_4 *2100*

four ratio configuration

1_1_1_7 *120*

1_1_2_6 *1260*

1_1_3_5 *2520*

1_1_4_4 *1575*

1_2_2_5 *3780*

1_2_3_4 *12600*

1_3_3_3 *2800*

2_2_2_4 *3150*

2_2_3_3 *6300*

five ratio configuration

1_1_1_1_6 *210*

1_1_1_2_5 *2520*

1_1_1_3_4 *4200*

1_1_2_2_4 *9450*

1_1_2_3_3 *12600*

1_2_2_2_3 *12600*

2_2_2_2_2 *945*

six ratio configuration

1_1_1_1_1_5 *252*

1_1_1_1_2_4 *3150*

1_1_1_1_3_3 *2100*

1_1_1_2_2_3 *12600*

1_1_2_2_2_2 *4725*

seven ratio configuration

1_1_1_1_1_1_4 *210*

1_1_1_1_1_2_3 *2520*

1_1_1_1_2_2_2 *3150*

eight ratio configuration

1_1_1_1_1_1_1_3 *120*

1_1_1_1_1_1_2_2 *630*

nine ratio configuration

1_1_1_1_1_1_1_1_2 *45*

The total number of possible models is 115975 (including one-ratio model and free-ratio model).

Supplementary data for 50 cases

CASE 1: *wsp* genes from *Wolbachia*

Sequences: ftp://ftp.ebi.ac.uk/pub/databases/embl/align/ALIGN_000201

Alignment: Not change

^aPhylogeny: Build by Clustalx1.83

^bPrevious study:

ORM:

$\ln L = -4122.011213$ $\omega_0 = 0.2616$ $\kappa = 5.84225$ $AIC = 2p - 2\ln L = 2p + 8244.022426$

TwoRM:

$\ln L = -4106.214548$ $\omega_0 = 0.3303$ $\omega_1 = 0.1099$ (clade in green color) $\kappa = 5.94271$
 $AIC = 2p + 2 + 8212.429096 = 2p + 8214.429096$

FRM:

$\ln L = -4044.313860$ $AIC = 2p + 2 \times 120 + 8088.62772 = 8328.62772$

^cOBSM Method I

TwoRM:

$\ln L = -4113.270031$ $\omega_0 = 0.25081$ $\omega_1 = 999.0000$ $\kappa = 5.84608$
 $AIC = 2p + 2 + 8226.540062 = 2p + 8228.540062$

ThreeRM:

$\ln L = -4106.101581$ $\omega_0 = 0.26773$ $\omega_1 = 999.0000$ $\omega_2 = 0.04649$ $\kappa = 5.88205$
 $AIC = 2p + 2 \times 2 + 8212.203162 = 2p + 8216.203162$

FourRM:

$\ln L = -4102.095931$ $\omega_0 = 0.28110$ $\omega_1 = 999.0000$ $\omega_2 = 0.05263$ $\omega_3 = 0.05459$ $\kappa = 5.92788$
 $AIC = 2p + 2 \times 3 + 8204.191862 = 2p + 8210.191862$

FiveRM:

$\ln L = -4097.601521$ $\omega_0 = 0.29273$ $\omega_1 = 0.05186$ $\omega_2 = 999.0000$ $\omega_3 = 0.05648$ $\omega_4 = 0.05283$ $\kappa = 5.93956$
 $AIC = 2p + 2 \times 4 + 8195.203042 = 2p + 8203.203042$

Compare:

LRT:

FourRM (OBSM) vs TwoRM (Hypothesis)

Df=2 $2\Delta l = 8.237234$ p-value=0.01627

FiveRM (OBSM) vs TwoRM (Hypothesis)

Df=3 $2\Delta l = 17.226054$ p-value=0.00063

FRM vs FiveRM

Df=116 $2\Delta l = 106.575322$ p-value=0.7231

AIC: the FiveRM is best model

a: the phylogeny may be a little different with phylogeny used in previous study

b: the TwoRM is consider to be the good model in previous study

c: only SBTRMs (single branch two ratio models) and 3-rm, 4-rm, 5-rm are calculate along the method 1 (too many branches to calculate and this three models are good enough to prove our method).

CASE 2: NRPD2/NRPE2-like Gene Family

Sequences: <http://www.biomedcentral.com/content/supplementary/1471-2148-10-45-S4.TXT>

Alignment: Not change

Phylogeny: is congruent with Fig .5 in original paper

Previous study:

ORM:

$$\ln L = -375.434947 \quad \omega_0 = 0.30837 \quad \kappa = 2.52005$$

$$AIC = 2p + 750.869894$$

TwoRM (a): the branch with red arrow is ω_1 (one branch)

$$\ln L = -372.101882 \quad \omega_0 = 0.2672 \quad \omega_1 = 999.0000 \quad \kappa = 2.52011$$

$$AIC = 2p + 2 + 744.203764 = 2p + 746.203764$$

TwoRM (b): the branches with blue arrow are all ω_1 (three branches)

$$\ln L = -368.960632 \quad \omega_0 = 0.2311 \quad \omega_1 = 999.0000 \quad \kappa = 2.53092$$

$$AIC = 2p + 2 + 737.921264 = 2p + 739.921264$$

FRM:

$$\ln L = -341.184883$$

$$AIC = 2p + 2 * 59 + 682.369766 = 2p + 800.369766$$

^aOBSM Method I:

TwoRM: same with TwoRM (a)

$$\ln L = -372.101882 \quad \omega_0 = 0.2672 \quad \omega_1 = 999.0000 \quad \kappa = 2.52011$$

$$AIC = 2p + 2 + 744.203764 = 2p + 746.203764$$

21-RM (final optimal model):

$$\ln L = -348.670597 \quad \omega_0 = 0.1226 \quad \omega_1 = 999.0000 \quad \kappa = 2.66455$$

$$\omega_{2-20} \sim (999.00000 \ 999.00000 \ 999.00000 \ 0.00010 \ 0.00010 \ 0.00010 \ 0.00010 \ 0.00010 \ 0.06952$$

$$999.00000 \ 999.00000 \ 999.00000 \ 999.00000 \ 999.00000 \ 999.00000 \ 239.96815 \ 999.00000$$

$$999.00000 \ 999.00000)$$

$$AIC = 2p + 2 * 20 + 697.341194 = 2p + 737.341194$$

Compare:

LRT:

TwoRM (b) vs 21-RM

$$Df = 19 \quad 2\Delta l = 40.58007 \quad p\text{-value} = 0.0027$$

FRM vs 21-RM

$$Df = 39 \quad 2\Delta l = 14.971428 \quad p\text{-value} = 0.9998$$

AIC: the 21-RM is the best

a: more optimal models significant better than ORM is not shown.

CASE 3: Phytochrome Gene Family BDE
Sequences: according the Genbank id in paper
^aAlignment: MEGA4 (only CDS)
Phylogeny: is congruent with Fig 3B

^bPrevious study:

ORM:

$$\ln L = -21407.461562 \quad \omega_0 = 0.0851 \quad \kappa = 1.96981$$

$$AIC = 2p + 42814.923124$$

TwoRM: the branch in bord is ω_1 (one branch)

$$\ln L = -21407.425858 \quad \omega_0 = 0.0853 \quad \omega_1 = 0.0789 \quad \kappa = 1.97135$$

$$AIC = 2p + 2 + 42814.851716 = 2p + 42816.851716$$

FRM

$$\ln L = -21302.383260$$

$$AIC = 2p + 2 \cdot 18 + 42604.76652 = 2p + 42640.76652$$

^cOBSM Method I:

TwoRM:

$$\ln L = -21376.483344 \quad \omega_0 = 0.0931 \quad \omega_1 = 0.0058 \text{ (branch in red)} \quad \kappa = 2.01938$$

$$AIC = 2p + 2 + 42752.966688 = 2p + 42754.966688$$

ThreeRM:

$$\ln L = -21349.096363 \quad \omega_0 = 0.10388 \quad \omega_1 = 0.00574 \quad \omega_2 = 0.01537 \quad \kappa = 2.01938$$

$$AIC = 2p + 2 \cdot 2 + 42698.192726 = 2p + 42702.192726$$

NineRM (final optimal model):

$$\ln L = -21304.336157 \quad \omega_0 = 0.08317 \quad \omega_1 = 0.00533 \quad \omega_2 = 0.01537 \quad \kappa = 2.05211$$

$$\omega_{3-8} = (0.04479 \ 0.02142 \ 0.19176 \ 0.27542 \ 0.11115 \ 0.14108 \ 0.24123)$$

$$AIC = 2p + 2 \cdot 8 + 42608.672314 = 2p + 42624.672314$$

Compare:

LRT:

Obviously, optimal models of OBSM method I are all significant better than ORM

Df=1	2Δl=61.956436	p-value<0.001
Df=2	2Δl=116.730398	p-value<0.001
Df=8	2Δl=206.25081	p-value<0.001

FRM vs NineRM

Df=10	2Δl=3.905794	p-value= 0.951
-------	--------------	----------------

AIC: the NineRM is the best model

b: The two-ratios model gave about the same fit to the data as the one-ratio model (sentence in original paper).

c: more optimal models significant better than ORM is not shown.

CASE 4: Phytochrome Gene Family ACF

Sequences: according the Genbank id in paper Zea Mays (AY260865)

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 3A

Previous study:

ORM:

$\ln L = -29762.994812$ $\omega_0 = 0.0883$ $\kappa = 2.02011$
 $AIC = 2p + 59525.989624$

TwoRM: the branch in bold is ω_1 (one branch)

$\ln L = -29760.573826$ $\omega_0 = 0.0889$ $\omega_1 = 0.0159$ $\kappa = 2.04276$
 $AIC = 2p + 2 + 59521.147652 = 2p + 59523.147652$

FRM:

$\ln L = -29695.459319$
 $AIC = 2p + 2 * 26 + 59390.918638 = 2p + 59442.918638$

^bOBSM Method I:

TwoRM:

$\ln L = -29741.542874$ $\omega_0 = 0.08267$ $\omega_1 = 0.28210$ (branch in red) $\kappa = 2.02372$
 $AIC = 2p + 2 + 59483.085748 = 2p + 59485.085748$

ThreeRM:

$\ln L = -29731.152707$ $\omega_0 = 0.08641$ $\omega_1 = 0.27317$ $\omega_2 = 0.02539$ $\kappa = 2.03679$
 $AIC = 2p + 2 * 2 + 59462.305414 = 2p + 59466.305414$

11-RM (final optimal model):

$\ln L = -29704.587670$ $\omega_0 = 0.09360$ $\omega_1 = 0.26076$ $\omega_2 = 0.01537$ $\kappa = 2.06294$ $\omega_{3-8} \sim (0.06238$
 $0.00782 \ 0.05597 \ 0.24491 \ 0.06556 \ 0.12880 \ 0.16551 \ 0.00385 \ 0.05778)$
 $AIC = 2p + 2 * 10 + 59409.17534 = 2p + 59429.17534$

Compare:

LRT:

Obviously, optimal models of OBSM method I are all significant better than TwoRM in original study

Df=1	2Δl=38.061904	p-value<0.001
Df=1	2Δl=58.842226	p-value<0.001
Df=9	2Δl=111.972312	p-value<0.001

FRM vs 11-RM

Df=16 2Δl=18.256702 p-value=0.309

AIC: the 11-RM is the best model

a: the alignment may have a little different

b: more optimal models significant better than ORM is not shown.

CASE 5: *BRCA1*

Sequences: according the Genbank id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: Tree is congruent with Fig 2

Previous study:

ORM:

$$\ln L = -9343.871464 \quad \omega_0 = 0.6299 \quad \kappa = 4.42745$$

$$AIC = 2p + 18687.742928$$

TwoRM: assume the human and chimpanzee lineages are under positive selection

$$\ln L = -9341.068346 \quad \omega_0 = 0.6125 \quad \omega_1 = 1.9228 \quad \kappa = 4.42646$$

$$AIC = 2p + 2 + 18682.136692 = 2p + 18684.136692$$

FRM:

$$\ln L = -9336.729565$$

$$AIC = 2p + 2 \times 13 + 18673.45913 = 2p + 18699.45913$$

OBSM Method I:

TwoRM: assume the human lineage is under positive selection

$$\ln L = -9341.818380 \quad \omega_0 = 0.6170 \quad \omega_1 = 1.9122 \quad \kappa = 4.42674$$

$$AIC = 2p + 2 + 18683.63676 = 2p + 18685.63676$$

OBSM Method II:

Same with Method I

OBSM Method III ($\kappa=0.5$):

TwoRM: same with hypothesis model

$$\ln L = -9341.068346 \quad \omega_0 = 0.6125 \quad \omega_1 = 1.9228 \quad \kappa = 4.42646$$

$$AIC = 2p + 2 + 18682.136692 = 2p + 18684.136692$$

ThreeRM:

$$\ln L = -9338.174830 \quad \omega_0 = 0.67931 \quad \omega_1 = 1.92326 \text{ (in bold)} \quad \omega_2 = 0.36342 \text{ (in red)} \quad \kappa = 4.43057$$

$$AIC = 2p + 2 \times 2 + 18676.34966 = 2p + 18680.34966$$

FourRM (final optimal model):

$$\ln L = -9337.502099 \quad \omega_0 = 0.70281 \quad \omega_1 = 1.92336 \quad \omega_2 = 0.36751 \quad \omega_3 = 0.50489 \text{ (in blue)} \quad \kappa = 4.43149$$

$$AIC = 2p + 2 \times 3 + 18675.004198 = 2p + 18681.004198$$

Compare:

LRT:

ThreeRM of Method III vs Hypothesis TwoRM:

$$Df = 1 \quad 2\Delta l = 5.787032 \quad p\text{-value} = 0.01614$$

FourRM of Method III vs Hypothesis TwoRM:

$$Df = 2 \quad 2\Delta l = 7.132494 \quad p\text{-value} = 0.02826$$

FRM vs FourRM

$$Df = 10 \quad 2\Delta l = 1.545068 \quad p\text{-value} = 0.9987$$

AIC: the ThreeRM of Method III is the best model

a: the alignment may have a little difference

CASE 6: Chalcone Synthase Genes in Dendranthema

Sequences: according the Genbank id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: according the phylogeny with Fig 1

Previous study:

^bModel A (ORM): $\omega_0=\omega_1=\omega_2=\omega_3$

$\ln L = -6445.168687$ $\omega_0 = 0.0533$ $\kappa = 1.67328$

$AIC = 2p + 12890.337374$

^bModel C: $\omega_0=\omega_1=\omega_3, \omega_2$

$\ln L = -6427.273174$ $\omega_0 = 0.0466$ $\omega_2 = 0.4666$ $\kappa = 1.67212$

$AIC = 2p + 2 + 12854.546348 = 2p + 12856.546348$

^bModel F: $\omega_0, \omega_1, \omega_2, \omega_3$

$\ln L = -6426.880750$ $\omega_0 = 0.0471$ $\omega_1 = 177.6597$ $\omega_2 = 0.5023$ $\omega_3 = 0.0220$ $\kappa = 1.67203$

$AIC = 2p + 2 * 3 + 12853.7615 = 2p + 12859.7615$

^bModel H: $\omega_0=\omega_2, \omega_1, \omega_3$

$\ln L = -6444.974866$ $\omega_0 = 0.0538$ $\omega_1 = 0.57580$ $\omega_3 = 0.03054$ $\kappa = 1.67312$

$AIC = 2p + 2 * 2 + 12889.949732 = 2p + 12893.949732$

FRM:

$\ln L = -6312.561448$

$AIC = 2p + 2 * 40 + 12625.122896 = 2p + 12705.122896$

^cOBSM Method I:

TwoRM:

$\ln L = -6424.527983$ $\omega_0 = 0.0611$ $\omega_{phchsa} = 0.0045$ (lineage marked with red arrow) $\kappa = 1.63192$

$AIC = 2p + 2 + 12849.055966 = 2p + 12851.055966$

ThreeRM:

$\ln L = -6409.685071$ $\omega_0 = 0.0538$ $\omega_{phchsa} = 0.0045$ $\omega_2 = 0.4195$ (in red) $\kappa = 1.63300$

$AIC = 2p + 2 * 2 + 12819.370142 = 2p + 12823.370142$

15-RM(final optimal model):

$\ln L = -6327.428407$ $\omega_0 = 0.08856$ $\omega_{phchsa} = 0.00381$ $\omega_2 = 0.30785$ (in red) $\kappa = 1.65631$

$\omega_{3-14} \sim (0.26757 \ 0.00610 \ 0.01269 \ 0.00200 \ 0.00863 \ 0.03584 \ 1.18773 \ 0.00978 \ 0.10479 \ 0.00755 \ 0.49616 \ 0.09588)$

$AIC = 2p + 2 * 14 + 12654.856814 = 2p + 12682.856814$

Compare:

LRT:

Obviously, optimal models of OBSM method I are all significant better than mode F in original study

FRM vs 15-RM

Df=26 $2\Delta l = 29.733918$ p-value=0.2787

AIC: the 15-RM is the best model

a: the alignment may have a little different

b: the model is congruent with Table 3 in original paper, model F is best among these models

c: more optimal models significant better than ORM is not shown.

CASE 7: Triosephosphate Isomerase

Sequences: according the Genbank id in paper, mouse NM_009415

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 1

Previous study:

ORM:

$\ln L = -4089.438559$ $\omega_0 = 0.05236$ $\kappa = 1.50940$

$AIC = 2p + 8178.877118$

ThreeRM: (vs TwoRM df=2 p-value=0.0476)

$\ln L = -4086.394347$ $\omega_0 = 0.04886$ $\omega_B = 0.51321$ $\omega_A = 0.25327$ $\kappa = 1.49552$

$AIC = 2p + 2 \times 2 + 8172.788694 = 2p + 8176.788694$

FRM: (vs ThreeRM df=14 $2\Delta l = 35.58811$ p-value=0.0012)

$\ln L = -4068.600292$ $\omega_0 = 0.08235$ $\omega_B = 0.09533$ $\omega_A = 383.10392$ $\kappa = 1.50190$

$\omega_{3-17} = (0.00010 \ 0.07190 \ 0.11167 \ 0.12394 \ 0.05770 \ 0.01807 \ 999.00000 \ 0.09457 \ 0.03877 \ 0.05280 \ 0.01263 \ 0.02210 \ 0.04044 \ 0.05025)$

$AIC = 2p + 2 \times 16 + 8137.200584 = 2p + 8169.200584$

^bOBSM Method I:

TwoRM:

$\ln L = -4083.894080$ $\omega_0 = 0.0495$ $\omega_1 = 999.0000$ (in red) $\kappa = 1.50150$

$AIC = 2p + 2 + 8167.78816 = 2p + 8169.78816$

ThreeRM:

$\ln L = -4081.261109$ $\omega_0 = 0.05455$ $\omega_1 = 999.00000$ $\omega_2 = 0.02006$ $\kappa = 1.51481$

$AIC = 2p + 2 \times 2 + 8162.522218 = 2p + 8166.522218$

7-RM (final optimal model):

$\ln L = -4072.634101$ $\omega_0 = 0.05390$ $\omega_1 = 999.00000$ $\omega_2 = 0.02058$ $\kappa = 1.50605$ $\omega_{3-6} = 0.02192$

$1.83280 \ 0.12169 \ 0.00010$

$AIC = 2p + 2 \times 6 + 8145.268202 = 2p + 8157.268202$

Compare:

LRT:

Obviously, optimal model TwoRM ThreeRM and 7-RM of OBSM method I are all significant better than ThreeRM in original study, and FRM is not significant better than final optimal model 7-RM of method I (df=10, $2\Delta l = 8.067618$, p-value=0.6222)

AIC: the 7-RM is the best model

a: the alignment may have a little different

b: more optimal models significant better than ORM is not shown.

CASE 8: Chalcone Synthase Genes in Morning Glories (Ipomoea)

Sequences: according the Genbank id paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 1

Previous study:

ORM:

$\ln L = -17880.400682$ $\omega_0 = 0.06646$ $\kappa = 1.56773$
 $AIC = 2p + 35760.801364$

TwoRM:

$\ln L = -17878.564824$ $\omega_0 = 0.06572$ $\omega_a = 0.15751$ $\kappa = 1.56350$
 $AIC = 2p + 2 + 35757.129648 = 35759.129648$

7-RM:

$\ln L = -17761.365327$ $\omega_0 = 0.03899$ $\omega_a = 0.19258$ $\omega_b = 0.07423$ $\omega_c = 0.09128$ $\omega_A = 0.22714$
 $\omega_B = 0.03954$ $\omega_C = 0.12176$ $\kappa = 1.58659$
 $AIC = 2p + 2 \times 6 + 35522.730654 = 2p + 35534.730654$

FRM:

$\ln L = -17574.394267$
 $AIC = 2p + 2 \times 87 + 35148.788534 = 2p + 35322.788534$

^bOBSM Method I:

11-RM:

$\ln L = -17750.913944$ $\omega_0 = 0.05864$ $\omega_1 = 0.31287$ $\omega_2 = 0.00063$ $\kappa = 1.55469$
 $\omega_{3-10} \sim (0.42278 \ 0.787158 \ 0.35731 \ 0.00296 \ 0.01615 \ 0.01569 \ 0.21626 \ 0.16033)$
 $AIC = 2p + 2 \times 10 + 35501.827888 = 2p + 35521.827888$

17-RM:

$\ln L = -17713.043181$ $\omega_0 = 0.06284$ $\omega_1 = 0.31279$ $\omega_2 = 0.00079$ $\kappa = 1.56033$
 $\omega_{3-16} \sim (0.18931 \ 0.78795 \ 0.35979 \ 0.00279 \ 0.01517 \ 0.01599 \ 0.21449 \ 0.16052 \ 0.20426 \ 0.01161$
 $0.03188 \ 0.01282 \ 0.99054 \ 1.12049)$
 $AIC = 2p + 2 \times 16 + 35426.086362 = 2p + 35458.086362$

23-RM:

$\ln L = -17682.501932$ $\omega_0 = 0.06663$ $\omega_1 = 0.31268$ $\omega_2 = 0.00011$ $\kappa = 1.55667$
 $\omega_{3-22} \sim (0.19474 \ 0.00010 \ 0.35606 \ 0.00271 \ 0.01530 \ 0.02032 \ 0.21387 \ 0.16031 \ 0.18937 \ 0.01379$
 $0.02993 \ 0.01296 \ 0.96491 \ 1.11444 \ 80.63478 \ 0.02020 \ 0.01631 \ 44.84182 \ 0.03889 \ 0.31906)$
 $AIC = 2p + 2 \times 22 + 35365.003864 = 35409.003864$

30-RM:

$\ln L = -17647.194774$ $\omega_0 = 0.07302$ $\omega_1 = 0.31061$ $\omega_2 = 0.02090$ $\kappa = 1.54646$
 $\omega_{3-29} \sim (0.18670 \ 0.00010 \ 0.33513 \ 0.00250 \ 0.00893 \ 0.03025 \ 0.20192 \ 0.16088 \ 0.20693 \ 0.01832$
 $0.02034 \ 0.01717 \ 0.95975 \ 1.12263 \ 999.00000 \ 0.02882 \ 0.01367 \ 999.00000 \ 0.02850 \ 0.31543$
 $0.16603 \ 0.82103 \ 0.02126 \ 0.26883 \ 0.01184 \ 0.23291 \ 0.01319)$
 $AIC = 2p + 2 \times 29 + 35294.389548 = 2p + 35352.389548$

Compare:

LRT:

Obviously, The Optimal models of OBSM method I are all significant better than 7-RM in original study

Df=4 $2\Delta l = 20.902766$ p-value=3.31E-04

Df=10 $2\Delta l = 96.644292$ p-value=5.42E-07

Df=16 $2\Delta l = 157.72679$ p-value=2.67E-10

Df=23 $2\Delta l = 228.341106$ p-value=1.98E-14

FRM vs 30-RM (too many calculation, didn't finish searching)

Df=58 $2\Delta l = 145.601014$ p-value=1.754e-009

AIC: the FRM is the best model

a: the alignment may have a little different

b: more optimal models significant better than ORM is not shown

CASE 9: Globin Gene Family

^aSequences: according the Genbank id in paper

^bAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 3

Previous study:

ORM:

$$\ln L = -4401.552483 \quad \omega_0 = 0.23436 \quad \kappa = 2.15341$$

$$\text{AIC} = 2p + 8803.104966$$

^cThreeRM:

$$\ln L = -4396.934657 \quad \omega_0 = 0.45619 \quad \omega_\beta = 0.17242 \quad \omega_\gamma = 0.21530 \quad \omega_\epsilon = 0.27402 \quad \kappa = 2.14474$$

$$\text{AIC} = 2p + 2 \times 2 + 8793.869314 = 2p + 8797.869314$$

FRM:

$$\ln L = -4356.459281$$

$$\text{AIC} = 2p + 2 \times 47 + 8712.918562 = 2p + 8806.918562$$

^dOBSM Method I:

TwoRM:

$$\ln L = -4392.037887 \quad \omega_0 = 0.2241 \quad \omega_1 = 999.0000 \quad \kappa = 2.14232$$

$$\text{AIC} = 2p + 2 + 8784.075774 = 2p + 8786.075774$$

ThreeRM:

$$\ln L = -4385.218567 \quad \omega_0 = 0.21393 \quad \omega_1 = 999.00000 \quad \omega_2 = 1.44661 \quad \kappa = 2.14455$$

$$\text{AIC} = 2p + 2 \times 2 + 8770.437134 = 2p + 8774.437134$$

15-RM(Final optimal model):

$$\ln L = -4366.651880 \quad \omega_0 = 0.21374 \quad \omega_1 = 999.00000 \quad \omega_2 = 1.44624 \quad \kappa = 2.14788$$

$$\omega_{3-14} = (0.05223 \ 0.06145 \ 0.59826 \ 999.00000 \ 0.08814 \ 0.78402 \ 0.13002 \ 0.40781 \ 0.07391 \ 3.57236 \ 999.00000 \ 0.00010)$$

$$\text{AIC} = 2p + 2 \times 14 + 8733.30376 = 2p + 8761.30376$$

Compare:

LRT:

Obviously, optimal models of OBSM method I are all significant better than FourRM in original study

FRM vs 15-RM

$$Df = 33 \quad 2\Delta l = 20.385198 \quad p\text{-value} = 0.9578$$

AIC: the 15-RM is the best model

a: few sequences (human beta and mouse beta) may not congruent with the original paper

b: the alignment may have a little different

c: we try 4 models, and chose the model of best lnL value among these four

d: more optimal models significant better than ORM is not shown.

CASE 10: Pollen-Specific Oleosin-Like Gene Family

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 4

^bPrevious study:

ORM:

$$\ln L = -2460.115266 \quad \omega_0 = 0.67749 \quad \kappa = 1.64825 \\ AIC = 2p + 4920.230532$$

FRM:

$$\ln L = -2432.734251 \\ AIC = 2p + 2 \times 45 + 4865.468502 = 2p + 4955.468502$$

^cOBSM Method I:

TwoRM:

$$\ln L = -2457.443865 \quad \omega_0 = 0.6643 \quad \omega_1 = 999.0000 \quad \kappa = 1.65263 \\ AIC = 2p + 2 + 4914.88773 = 2p + 4916.88773$$

ThreeRM:

$$\ln L = -2454.763631 \quad \omega_0 = 0.6484 \quad \omega_1 = 999.00000 \quad \omega_2 = 999.00000 \quad \kappa = 1.65944 \\ AIC = 2p + 2 \times 2 + 4909.527262 = 2p + 4913.527262$$

8-RM(Final optimal model):

$$\ln L = -2446.335287 \quad \omega_0 = 0.72178 \quad \omega_1 = 999.00000 \quad \omega_2 = 999.00000 \quad \kappa = 1.65607 \quad \omega_{3-7} \sim (0.05638 \\ 0.00010 \ 0.17969 \ 0.00010 \ 0.27021 \ 0.06238) \\ AIC = 2p + 2 \times 7 + 4892.670574 = 2p + 4906.670574$$

Compare:

LRT:

Obviously, optimal models of OBSM method I are all significant better than ORM

Df=1	2Δl=5.342802	p-value=0.0208
Df=2	2Δl=10.70327	p-value=0.0047
Df=7	2Δl=27.559958	p-value=0.00026

FRM vs 8-RM

$$Df=38 \quad 2\Delta l=27.202072 \quad p\text{-value}=0.903$$

AIC: the 8-RM is the best model

a: the alignment may have a little different

b: Comparisons of branch models revealed no significant differences between different lineages in the gene phylogeny (sentence in original paper).

c: more optimal models significant better than ORM is not shown.

CASE 11: Pollen-Specific Oleosin-Like Gene Family

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 4

Previous study:

ORM:

$$\ln L = -2662.201006 \quad \omega_0 = 0.1945 \quad \kappa = 2.24635$$

$$AIC = 2p + 5324.402012$$

Hypothesis TwoRM:

$$\ln L = -2657.211303 \quad \omega_0 = 0.1706 \quad \omega_1 = 2.5434 \quad \kappa = 2.24962$$

$$AIC = 2p + 2 + 5314.422606 = 2p + 5316.422606$$

FRM:

$$\ln L = -2641.379441$$

$$AIC = 2p + 2 \times 25 + 5282.758882 = 2p + 5332.758882$$

OBSM Method I:

TwoRM: same with hypothesis TwoRM

$$\ln L = -2657.211303 \quad \omega_0 = 0.1706 \quad \omega_1 = 2.5434 \quad \kappa = 2.24962$$

$$AIC = 2p + 2 + 5314.422606 = 2p + 5316.422606$$

ThreeRM:

$$\ln L = -2652.975199 \quad \omega_0 = 0.18837 \quad \omega_1 = 2.40186 \quad \omega_2 = 0.00010 \quad \kappa = 2.24783$$

$$AIC = 2p + 2 \times 2 + 5305.950398 = 2p + 5309.950398$$

5-RM(Final optimal model, ω ratio is labeled above the branch):

$$\ln L = -2649.352240 \quad \omega_0 = 0.21385 \quad \omega_1 = 2.35867 \quad \omega_2 = 0.00010 \quad \kappa = 2.24539 \quad \omega_{3,4} \sim (0.05109$$

0.00010)

$$AIC = 2p + 2 \times 4 + 5298.70448 = 2p + 5306.70448$$

Compare:

LRT:

Obviously, optimal models (ThreeRM and 5-RM) of OBSM method I are significant better than Hypothesis

TwoRM

$$Df=1 \quad 2\Delta l=8.472208 \quad p\text{-value}=0.0036$$

$$Df=3 \quad 2\Delta l=15.718126 \quad p\text{-value}=0.0013$$

FRM vs 5-RM

$$Df=21 \quad 2\Delta l=15.945598 \quad p\text{-value}=0.7726$$

AIC: the 5-RM is the best model

a: the alignment may have a little different

CASE 12: Rhodopsin Gene

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 2(not include outgroup)

^bPrevious study:

ORM:

$$\ln L = -1563.668037 \quad \omega_0 = 0.2655 \quad \kappa = 2.81281$$

$$AIC = 2p + 3127.336074$$

FRM:

$$\ln L = -1549.634709$$

$$AIC = 2p + 2 \cdot 31 + 3099.269418 = 2p + 3161.269418$$

OBSM Method I:

FourRM:

$$\ln L = -1559.603182 \quad \omega_0 = 0.26495 \quad \omega_1 = 999.00000 \quad \omega_{2,3} = 0.00010 \quad \kappa = 2.80964$$

$$AIC = 2p + 2 \cdot 3 + 3119.206364 = 2p + 3125.206364$$

Compare:

LRT:

FRM vs FourRM

$$Df = 3 \quad 2\Delta l = 8.12971 \quad p\text{-value} = 0.0434$$

AIC: the FourRM is the best model .

a: the alignment may have a little different

b: no likelihood ratio test was significant in the comparison between the ω ratio of model M0 and the two-ratio model that estimates two different ω ratios (sentence in original paper).

CASE 13: M_gamma_type_MADS_Box_Genes

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 2 (original paper)

Previous study:

ORM:

$$\ln L = -21447.371961 \quad \omega_0 = 0.51565 \quad \kappa = 1.65870$$

$$\text{AIC} = 2p + 42894.743922$$

TwoRM:

$$\ln L = -21441.236742 \quad \omega_0 = 0.47932 \quad \omega_1 = 0.76610 \quad \kappa = 1.66651$$

$$\text{AIC} = 2p + 2 + 42882.473484 = 2p + 42884.473484$$

FRM:

$$\ln L = -21330.684136$$

$$\text{AIC} = 2p + 2 \times 67 + 42661.368272 = 2p + 42795.368272$$

^bOBSM Method I:

TwoRM:

$$\ln L = -21412.775657 \quad \omega_0 = 0.54361 \quad \omega_1 = 0.00351 \quad \kappa = 1.67500$$

$$\text{AIC} = 2p + 2 + 42825.551314 = 2p + 42827.551314$$

18-RM:

$$\ln L = -21353.379969 \quad \omega_0 = 0.56269 \quad \omega_1 = 0.00010 \quad \omega_2 = 0.00404 \quad \kappa = 1.71444$$

$$\omega_{3-18} \sim (0.04736 \ 0.08443 \ 0.08913 \ 214.44949 \ 0.00668 \ 1.12431 \ 0.27295 \ 0.16102 \ 1.40938 \ 0.93848 \\ 0.10296 \ 1.73498 \ 0.15243 \ 999.00000 \ 999.00000 \ 1.93258)$$

$$\text{AIC} = 2p + 2 \times 17 + 42706.759938 = 2p + 42740.759938$$

Compare:

LRT:

Obviously, optimal models of OBSM method I are all significant better than hypothesis model.

FRM vs 18-RM

$$Df = 50 \quad 2\Delta l = 45.391666 \quad p\text{-value} = 0.6585$$

AIC: 18-RM is the best model

a: the alignment result show that it may not suitable for such research (too distant to deduce its ancestral sequences).

b: more optimal models significant better than ORM is not shown.

CASE 14: Digestive_RNASE1_Genes

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 6

^bPrevious study:

ORM:

$$\ln L = -1409.406777 \quad \omega_0 = 0.66410 \quad \kappa = 2.95678$$

$$AIC = 2p + 2818.813554$$

SixRM:

$$\ln L = -1397.582193 \quad \omega_0 = 0.37396 \quad \omega_1 = 1.07877 \quad \omega_5 = 1.08460 \quad \omega_{2,3,4} = 999.0000 \quad \kappa = 2.96104$$

$$AIC = 2p + 2 \cdot 5 + 2795.164386 = 2p + 2805.164386$$

FRM:

$$\ln L = -1379.257362$$

$$AIC = 2p + 2 \cdot 57 + 2758.514724 = 2p + 2872.514724$$

^cOBSM Method I:

SixRM: (marked with red arrow and number)

$$\ln L = -1394.828338 \quad \omega_0 = 0.62326 \quad \omega_{2,3,4} = 0.00010 \quad \omega_{1,5} = 999.0000 \quad \kappa = 2.96773$$

$$AIC = 2p + 2 \cdot 5 + 2789.656676 = 2p + 2799.656676$$

TenRM:

$$\ln L = -1387.396695 \quad \omega_0 = 0.75823 \quad \omega_1 = 0.12173 \quad \omega_2 = 0.08177 \quad \kappa = 2.96773$$

$$\omega_{3,6} = 0.00010 \quad \omega_{7,9} = 999.0000$$

$$AIC = 2p + 2 \cdot 9 + 2774.79339 = 2p + 2792.79339$$

Compare:

LRT:

Obviously, optimal SixRM of OBSM method I is significant better than Hypothesis SixRM, and final optimal model TenRM is also significant better than Hypothesis SixRM:

$$Df = 4 \quad 2\Delta l = 20.370996 \quad p\text{-value} = 0.00042$$

FRM vs TenRM

$$Df = 48 \quad 2\Delta l = 16.278666 \quad p\text{-value} = 0.9999$$

AIC: the TenRM is the best model

a: the alignment may have a little different

b: Only Hypothesis SixRM (Model D in original paper) is calculate while other two models (Model B and C) is confused described.

c: more optimal models significant better than ORM is not shown.

CASE 15: Pistillata PI genes

Sequences: according the id in paper, Matrix A

^aAlignment: MEGA4 (only CDS)

Phylogeny: is manual edit according Fig 1 (20)

Previous study:

ORM:

$\ln L = -2562.634189$ $\omega_0 = 0.29816$ $\kappa = 1.35424$

$AIC = 2p + 5125.268378$

^bHypothesis_TwoRM:

$\ln L = -2557.238920$ $\omega_0 = 0.3360$ $\omega_1 = 0.0750$ (red branches) $\kappa = 1.35376$

$AIC = 2p + 2 + 5114.47784$

FRM:

$\ln L = -2536.246228$

$AIC = 2p + 2 * 29 + 5072.492456 = 2p + 5130.492456$

^cOBSM Method I:

TwoRM:

$\ln L = -2558.902183$ $\omega_0 = 0.32180$ $\omega_1 = 0.01833$ $\kappa = 1.35176$

$AIC = 2p + 2 + 5117.804366 = 2p + 5119.804366$

FourRM:

$\ln L = -2555.807051$ $\omega_0 = 0.30531$ $\omega_1 = 999.00000$ $\omega_2 = 0.01985$ $\omega_3 = 2.56347$ $\kappa = 1.36062$

$AIC = 2p + 2 * 3 + 5111.614102 = 2p + 5117.614102$

TenRM:

$\ln L = -2542.526246$ $\omega_0 = 0.36921$ $\omega_1 = 0.04096$ $\omega_2 = 273.72251$ $\kappa = 1.37044$

$\omega_{3-9} = (2.33347 \ 0.80354 \ 0.11693 \ 0.04545 \ 999.00000 \ 0.14308 \ 0.13157)$

$AIC = 2p + 2 * 9 + 5085.052492 = 2p + 5103.052492$

Compare:

LRT:

Final optimal model is significant better than Hypothesis model.

Df=2 $2\Delta l = 2.863738$ p-value= /

Df=8 $2\Delta l = 29.425348$ p-value=2.670e-004

FRM vs TenRM

Df=20 $2\Delta l = 12.560036$ p-value=0.89545

AIC: the TenRM is the best model

a: the alignment may have a little different

b: Three hypothesis models are tried and the best one is show here.

c: more optimal models significant better than ORM is not shown.

CASE 16: SWS1_HeShunping

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 2a (not include outgroup Danio_rerio)

Previous study:

ORM:

$\ln L = -3022.465018$ $\omega_0 = 0.43018$ $\kappa = 2.31164$
 $AIC = 2p + 6044.930036$

TwoRM:

$\ln L = -3019.786644$ $\omega_0 = 0.5369$ $\omega_1 = 0.3157$ $\kappa = 2.31152$
 $AIC = 2p + 2 + 6039.573288 = 2p + 6041.573288$

FRM:

$\ln L = -2989.152966$
 $AIC = 2p + 2 * 63 + 5978.305932 = 2p + 6104.305932$

^bOBSM Method I:

TwoRM:

$\ln L = -3017.843433$ $\omega_0 = 0.36791$ $\omega_1 = 0.93827$ (red arrow) $\kappa = 2.31570$
 $AIC = 2p + 2 + 6035.686866 = 2p + 6037.686866$

ThreeRM:

$\ln L = -3014.933243$ $\omega_0 = 0.40355$ $\omega_1 = 0.93886$ $\omega_2 = 0.13203$ $\kappa = 2.32292$
 $AIC = 2p + 2 * 2 + 6029.866486 = 2p + 6033.866486$

11-RM:

$\ln L = -2998.267273$ $\omega_0 = 0.42293$ $\omega_1 = 0.93878$ $\omega_2 = 0.13461$ $\kappa = 2.32208$
 $\omega_{3-11} \sim (0.00010 \ 0.88031 \ 0.17243 \ 2.74113 \ 0.06632 \ 0.00010 \ 0.00010 \ 0.00010)$
 $AIC = 2p + 2 * 10 + 5996.534546 = 2p + 6016.534546$

Compare:

LRT:

Obviously, optimal models of OBSM method I are significant better than Hypothesis TwoRM

Df=1 $2\Delta l = 9.706802$ p-value=0.00184
Df=9 $2\Delta l = 43.038742$ p-value=2.121e-006

FRM vs 11-RM

Df=53 $2\Delta l = 18.228614$ p-value= 0.9999

AIC: the 11-RM is the best model

a: the alignment may have a little different and some CDS region marked in GenBank is confused and re-predict by FgenesH-M

b: more optimal models significant better than ORM is not shown.

CASE 17: SWS2_HeShunping

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 2b (not include outgroup Danio_rerio)

Previous study:

ORM:

$$\ln L = -2606.878710 \quad \omega_0 = 0.35308 \quad \kappa = 2.65541 \\ AIC = 2p + 5213.75742$$

TwoRM:

$$\ln L = -2604.305376 \quad \omega_A = 0.4883 \quad \omega_B = 0.2517 \quad \kappa = 2.65793 \\ AIC = 2p + 2 + 5208.610752 = 2p + 5210.610752$$

FRM:

$$\ln L = -2576.232247 \\ AIC = 2p + 2 * 39 + 5152.464494 = 2p + 5230.464494$$

^bOBSM Method I:

^cTwoRM:

$$\ln L = -2600.116139 \quad \omega_0 = 0.31892 \quad \omega_1 = 999.00000 \text{ (red)} \quad \kappa = 2.65897 \\ AIC = 2p + 2 + 5200.232278 = 2p + 5202.232278$$

ThreeRM:

$$\ln L = -2595.542652 \quad \omega_0 = 0.28778 \quad \omega_1 = 999.00000 \quad \omega_2 = 3.43585 \quad \kappa = 2.66214 \\ AIC = 2p + 2 * 2 + 5191.085304 = 2p + 5195.085304$$

EightRM:

$$\ln L = -2587.157466 \quad \omega_0 = 0.38075 \quad \omega_1 = 999.00000 \quad \omega_2 = 3.37710 \quad \kappa = 2.66500 \\ \omega_{3-7} \sim (0.00010 \ 999.00000 \ 0.14775 \ 0.00010 \ 0.17214) \\ AIC = 2p + 2 * 7 + 5174.314932 = 2p + 5188.314932$$

Compare:

LRT:

Obviously, TwoRM ThreeRM and EightRM of OBSM method I are all significant better than TwoRM in original study

FRM vs 8-RM

$$Df = 32 \quad 2\Delta l = 21.850438 \text{ p-value} = 0.9113$$

AIC: the EightRM is the best model

a: the alignment may have a little different

b: more optimal models significant better than ORM is not shown.

c: when we fix the $\omega_1 = 1$, the likelihood ratio test show the branch is significant larger than 1 (p-value=0.0166)

CASE 18: COR15 gene family

^aSequences: according the id in paper

^bAlignment: MEGA4 (only CDS)

Phylogeny: is built by ClustalX v1.83

^cPrevious study:

ORM:

$$\ln L = -2625.334652 \quad \omega_0 = 0.50507 \quad \kappa = 1.68299$$

$$AIC = 2p + 5250.669304$$

TwoRM A (clade model):

$$\ln L = -2625.315516 \quad \omega_0 = 0.51190 \quad \omega_1 = 0.48956 \text{ (clade a and b)} \quad \kappa = 1.68277$$

$$AIC = 2p + 2 + 5250.631032 = 2p + 5252.631032$$

TwoRM B:

$$\ln L = -2621.561591 \quad \omega_0 = 0.48746 \quad \omega_1 = 999.00000 \quad \kappa = 1.68444$$

$$AIC = 2p + 2 + 5243.123182 = 2p + 5245.123182$$

ThreeRM (clade model):

$$\ln L = -2622.053852 \quad \omega_0 = 0.48832 \quad \omega_1 = 0.95778 \text{ (clade a)} \quad \omega_2 = 0.28658 \text{ (clade b)} \quad \kappa = 1.68368$$

$$AIC = 2p + 2 \times 2 + 5244.107704 = 2p + 5248.107704$$

ThreeRM:

$$\ln L = -2623.293414 \quad \omega_0 = 0.52395 \quad \omega_a = 0.69154 \quad \omega_b = 0.20211 \quad \kappa = 1.68367$$

$$AIC = 2p + 2 \times 2 + 5246.586828 = 2p + 5250.586828$$

FRM:

$$\ln L = -2588.725430$$

$$AIC = 2p + 2 \times 50 + 5177.45086 = 2p + 5277.45086$$

^dOBSM Method I:

TwoRM:

$$\ln L = -2621.186909 \quad \omega_0 = 0.54334 \quad \omega_1 = 0.13860 \quad \kappa = 1.68683$$

$$AIC = 2p + 2 + 5242.373818 = 2p + 5244.373818$$

ThreeRM:

$$\ln L = -2617.640992 \quad \omega_0 = 0.52420 \quad \omega_1 = 0.13933 \quad \omega_2 = 999.00000 \quad \kappa = 1.68850$$

$$AIC = 2p + 2 \times 2 + 5235.281984 = 2p + 5239.281984$$

11-RM:

$$\ln L = -2600.853632 \quad \omega_0 = 0.63718 \quad \omega_1 = 0.15782 \quad \omega_2 = 999.00000 \quad \kappa = 1.67854$$

$$\omega_{3-10} \sim (0.04557 \ 999.00000 \ 999.00000 \ 2.64062 \ 0.20135 \ 0.22141 \ 0.15076 \ 0.00010)$$

$$AIC = 2p + 2 \times 10 + 5201.707264 = 2p + 5221.707264$$

Compare:

LRT:

Obviously, ThreeRM and 11-RM of OBSM method I are all significant better than TwoRM B in original study.

$$Df=1 \quad 2\Delta l=7.841198 \quad p\text{-value}=0.0051$$

$$Df=9 \quad 2\Delta l=41.415918 \quad p\text{-value}=4.201e-006$$

FRM vs 11-RM

$$Df=40 \quad 2\Delta l=24.256404 \quad p\text{-value}=0.9765$$

AIC: the 11-RM is the best model

a: the Genbank id showed in the original paper is confused with fig 2 in the original paper.

b: the alignment may have a little different

c: since the variety of sequences and phylogeny, the values of Hypothesis models suggested in original research can't reappear closely. We calculate several hypothesis models according the original paper implied.

d: more optimal models significant better than ORM is not shown.

CASE 19: rbcL genes in Conocephalum

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 2

^bPrevious study:

ORM:

$\ln L = -2809.960994$ $\omega_0 = 0.04252$ $\kappa = 2.26661$
 $AIC = 2p + 5619.921988$

TwoRM A:

$\ln L = -2800.778895$ $\omega_0 = 0.0381$ $\omega_1 = 999.0000$ (red branch) $\kappa = 2.26984$
 $AIC = 2p + 2 + 5601.55779 = 2p + 5603.55779$

TwoRM B (clade model):

$\ln L = -2803.416120$ $\omega_0 = 0.03830$ $\omega_1 = 0.85378$ (clade F) $\kappa = 2.26967$
 $AIC = 2p + 2 + 5606.83224 = 2p + 5608.83224$

FRM:

$\ln L = -2769.503035$
 $AIC = 2p + 2 * 27 + 5539.00607 = 2p + 5593.00607$

^cOBSM Method I:

TwoRM:

$\ln L = -2798.002656$ $\omega_0 = 0.08036$ $\omega_1 = 0.01317$ (blue branch) $\kappa = 2.27228$
 $AIC = 2p + 2 + 5596.005312 = 2p + 5598.005312$

11-RM:

$\ln L = -2774.625477$ $\omega_0 = 0.04157$ $\omega_1 = 0.00010$ $\omega_2 = 0.02645$ $\kappa = 2.28439$
 $\omega_{3-10} \sim (999.00000 \ 999.00000 \ 999.00000 \ 0.35504 \ 999.00000 \ 999.00000 \ 999.00000 \ 999.00000)$
 $AIC = 2p + 2 * 10 + 5549.250954 = 2p + 5569.250954$

Compare:

LRT:

Obviously, TwoRM and 11-RM of OBSM method I are all better than TwoRM A in original study.

Df=1 $2\Delta l = 7.841198$ p-value=0.0051
Df=9 $2\Delta l = 52.306836$ p-value=3.955e-008

FRM vs 11-RM

Df=17 $2\Delta l = 10.244884$ p-value=0.8930

AIC: the 11-RM is the best model

a: the alignment may have a little different

b: We calculate several hypothesis models according the original paper implied and only two optimal model is shown.

d: more optimal models significant better than ORM is not shown.

CASE 20: M_LWS_gene

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 2

^bPrevious study:

ORM:

$$\ln L = -2044.437533 \quad \omega_0 = 0.07246 \quad \kappa = 5.62795$$

$$AIC = 2p + 4088.875066$$

TwoRM A: (fruit bats)

$$\ln L = -2044.376785 \quad \omega_0 = 0.07171 \quad \omega_1 = 0.10955 \quad \kappa = 5.63077$$

$$AIC = 2p + 2 + 4088.75357 = 2p + 4090.75357$$

TwoRM B: (Yinpterochiroptera):

$$\ln L = -2043.712372 \quad \omega_0 = 0.07429 \quad \omega_1 = 0.00010 \quad \kappa = 5.62842$$

$$AIC = 2p + 2 + 4087.424744 = 2p + 4089.424744$$

TwoRM C: (Yangochiroptera):

$$\ln L = -2044.209877 \quad \omega_0 = 0.07124 \quad \omega_1 = 0.14036 \quad \kappa = 5.62884$$

$$AIC = 2p + 2 + 4088.419754 = 2p + 4090.419754$$

FRM:

$$\ln L = -2022.899266$$

$$AIC = 2p + 2 \times 31 + 4045.798532 = 2p + 4107.798532$$

OBSM Method I:

TwoRM:

$$\ln L = -2040.314104 \quad \omega_0 = 0.06467 \quad \omega_1 = 0.35723 \quad \kappa = 5.63068$$

$$AIC = 2p + 2 + 4080.628208 = 2p + 4082.628208$$

ThreeRM:

$$\ln L = -2037.013979 \quad \omega_0 = 0.06209 \quad \omega_1 = 0.35761 \quad \omega_2 = 999.00000 \quad \kappa = 5.63350$$

$$AIC = 2p + 2 \times 2 + 4074.027958 = 2p + 4078.027958$$

SixRM:

$$\ln L = -2033.039125 \quad \omega_0 = 0.07333 \quad \omega_1 = 0.34682 \quad \omega_2 = 999.00000 \quad \kappa = 5.64782$$

$$\omega_{3-5} \sim (0.02076 \ 0.00010 \ 0.00010)$$

$$AIC = 2p + 2 \times 5 + 4066.07825 = 2p + 4076.07825$$

Compare:

LRT:

Obviously, TwoRM and SixRM of OBSM method I are all significant better than ORM and TwoRM B in original study.

$$Df=1 \quad 2\Delta l=6.796536 \quad p\text{-value}=0.0091$$

$$Df=1 \quad 2\Delta l=13.396786 \quad p\text{-value}=2.521e-004$$

$$Df=4 \quad 2\Delta l=22.341504 \quad p\text{-value}=1.714e-004$$

FRM vs SixRM

$$Df=26 \quad 2\Delta l=20.279718 \quad p\text{-value}=0.77811$$

AIC: SixRM is the best model

a: the alignment may have a little different

b: We calculate several hypothesis models according the original paper implied

CASE 21: Gama N Crystallin Superfamily

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 1b

^bPrevious study:

ORM:

$$\ln L = -2520.154970 \quad \omega_0 = 0.05602 \quad \kappa = 2.16197$$

$$AIC = 2p + 5040.30994$$

ThreeRM:

$$\ln L = -2518.060965 \quad \omega_0 = 0.05518 \quad \omega_1 = 0.03996 \quad \omega_2 = 0.08011 \quad \kappa = 2.17468$$

$$AIC = 2p + 2 \times 2 + 5036.12193 = 2p + 5040.12193$$

FRM:

$$\ln L = -2490.212643$$

$$AIC = 2p + 2 \times 32 + 4980.425286 = 2p + 5044.425286$$

^cOBSM Method I:

TwoRM:

$$\ln L = -2514.560860 \quad \omega_0 = 0.06322 \quad \omega_1 = 0.00368 \quad \kappa = 2.13310$$

$$AIC = 2p + 2 + 5029.12172 = 2p + 5031.12172$$

ThreeRM:

$$\ln L = -2512.520642 \quad \omega_0 = 0.06098 \quad \omega_1 = 0.00389 \quad \omega_2 = 0.48310 \quad \kappa = 2.12556$$

$$AIC = 2p + 2 \times 2 + 5025.041284 = 2p + 5029.041284$$

SevenRM:

$$\ln L = -2503.413056 \quad \omega_0 = 0.06976 \quad \omega_1 = 0.00386 \quad \omega_2 = 0.60896 \quad \kappa = 2.09140$$

$$\omega_{3-6} = (0.00010 \ 0.01776 \ 0.00010 \ 999.00000)$$

$$AIC = 2p + 2 \times 6 + 5006.826112 = 2p + 5018.826112$$

Compare:

LRT:

Obviously, TwoRM, ThreeRM and 11-RM of OBSM method I are all significant better than ThreeRM (clade model) in original study.

$$Df=4 \quad 2\Delta l=29.295818 \quad p\text{-value}=6.807e-006$$

FRM vs SevenRM

$$Df=26 \quad 2\Delta l=26.400826 \quad p\text{-value}=0.4412$$

AIC: the SevenRM is the best model

a: the alignment may have a little different. This is a **typical case** that the sequences used in previous study may have some change (the CDS region) in the later version in GenBank.

b: We calculate several hypothesis models according the original paper implied and only ThreeRM (clade model) is shown.

c: more optimal models significant better than ORM is not shown.

CASE 22: HoxD Genes

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 2

^bPrevious study:

ORM:

$$\ln L = -4670.009057 \quad \omega_0 = 0.10712 \quad \kappa = 4.57654 \\ AIC = 2p + 9340.018114$$

Hypothesis model D:

$$\ln L = -4664.528087 \quad \omega_0 = 0.10328 \quad \omega_c = \omega_{tm} = 1.10715 \quad \kappa = 4.58076 \\ AIC = 2p + 2 + 9329.056174 = 2p + 9331.056174$$

Hypothesis model E:

$$\ln L = -4663.998915 \quad \omega_0 = 0.10332 \quad \omega_c = 0.77848 \quad \omega_{tm} = 347.95623 \quad \kappa = 4.58077 \\ AIC = 2p + 2 * 2 + 9327.99783 = 2p + 9331.99783$$

FRM:

$$\ln L = -4633.711096 \\ AIC = 2p + 2 * 67 + 9267.422192 = 2p + 9401.422192$$

^cOBSM Method I:

TwoRM:

$$\ln L = -4666.625301 \quad \omega_0 = 0.10437 \quad \omega_c = 0.77759 \quad \kappa = 4.57886 \\ AIC = 2p + 2 + 9333.250602 = 2p + 9335.250602$$

ThreeRM:

$$\ln L = -4663.476101 \quad \omega_0 = 0.10078 \quad \omega_c = 0.78253 \quad \omega_2 = 999.00000 \text{ (red arrow)} \quad \kappa = 4.56168 \\ AIC = 2p + 2 * 2 + 9326.952202 = 9330.952202$$

12-RM:

$$\ln L = -4649.422772 \quad \omega_0 = 0.09676 \quad \omega_1 = 0.78520 \quad \omega_2 = 999.00000 \quad \kappa = 4.56555 \\ \omega_{3-11} \sim (0.00010 \ 1.11937 \ 999.00000 \ 0.04225 \ 0.22779 \ 999.00000 \ 999.00000 \ 0.03805 \ 0.56985) \\ AIC = 2p + 2 * 11 + 9298.845544 = 2p + 9320.845544$$

Compare:

LRT:

12-RM of OBSM method I is significant better than two hypothesis models suggested in original study.

$$\begin{array}{lll} Df=10 & 2\Delta l=30.21063 & p\text{-value}=7.912e-004 \\ Df=9 & 2\Delta l=29.152286 & p\text{-value}=6.108e-004 \end{array}$$

FRM VS 12-RM

$$Df=56 \quad 2\Delta l=31.423352 \quad p\text{-value}=0.9967$$

AIC: the 12-RM is the best model

a: the alignment may have a little different.

b: We calculate several hypothesis models according the original paper implied and only optimal two models are shown.

c: more optimal models significant better than ORM is not shown.

CASE 23: Sperm Genes

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 2a

Previous study:

ORM:

$$\ln L = -1393.366296 \quad \omega_0 = 0.24026 \quad \kappa = 1.42403$$

$$AIC = 2p + 2786.732592$$

TwoRM:

$$\ln L = -1390.023095 \quad \omega_0 = 0.21576 \quad \omega_1 = 2.01228 \quad \kappa = 1.42746$$

$$AIC = 2p + 2 + 2780.04619 = 2p + 2782.04619$$

FRM:

$$\ln L = -1385.132916$$

$$AIC = 2p + 2 \cdot 9 + 2770.265832 = 2p + 2788.265832$$

OBSM Method I:

TwoRM:

It's same with hypothesis TwoRM.

FourRM: (branches are labeled by arrows)

$$\ln L = -1386.115852 \quad \omega_0 = 0.35322 \quad \omega_1 = 2.00656 \quad \omega_2 = 0.13203 \quad \omega_3 = 0.09261 \quad \kappa = 1.45027$$

$$AIC = 2p + 2 \cdot 3 + 2772.231704 = 2p + 2778.231704$$

Compare:

LRT:

FourRM of OBSM method I is significant better than hypothesis TwoRM in original study.

$$Df = 2 \quad 2\Delta l = 7.814486 \quad p\text{-value} = 0.0201$$

FRM vs FourRM

$$Df = 6 \quad 2\Delta l = 1.965872 \quad p\text{-value} = 0.9228$$

AIC: the FourRM is the best model

a: the alignment may have a little different.

CASE 24: Celegent_Lysozymes_lys

Sequences: according the id in paper

Alignment: MEGA4 and edit by manual to adjust with original paper

Phylogeny: is congruent with Fig 6B

^aPrevious study:

ORM:

$$\ln L = -5385.853900 \quad \omega_0 = 0.09985 \quad \kappa = 1.62418$$

$$\text{AIC} = 2p + 10771.7078$$

TwoRM A:

$$\ln L = -5374.337116 \quad \omega_0 = 0.09365 \quad \omega_j = 999.00000 \quad \kappa = 1.63526$$

$$\text{AIC} = 2p + 2 + 10748.674232 = 2p + 10750.674232$$

TwoRM B:

$$\ln L = -5376.391808 \quad \omega_0 = 0.10530 \quad \omega_{AJ} = 0.00010 \quad \kappa = 1.62222$$

$$\text{AIC} = 2p + 2 + 10752.783616 = 2p + 10754.783616$$

FRM:

$$\ln L = -5297.096110$$

$$\text{AIC} = 2p + 2 \cdot 43 + 10594.19222 = 2p + 10680.19222$$

^bOBSM Method I:

ThreeRM:

$$\ln L = -5365.417863 \quad \omega_0 = 0.09867 \quad \omega_j = 999.00000 \quad \omega_{AJ} = 0.00010 \quad \kappa = 1.63309$$

$$\text{AIC} = 2p + 2 \cdot 2 + 10730.835726 = 2p + 10734.835726$$

24-RM:

$$\ln L = -5307.308567 \quad \omega_0 = 0.08185 \quad \omega_j = 999.00000 \quad \omega_{AJ} = 0.00010 \quad \kappa = 1.60670$$

$$\omega_{3-23} \sim (0.01039 \ 0.46570 \ 999.00000 \ 0.24787 \ 1.84484 \ 0.10786 \ 999.00000 \ 0.26308 \ 0.22000 \\ 0.60653 \ 0.02760 \ 0.03657 \ 0.02596 \ 0.04064 \ 0.39669 \ 0.03612 \ 0.52866 \ 0.22819 \ 0.42078 \ 1.11020 \\ 0.36133)$$

$$\text{AIC} = 2p + 2 \cdot 23 + 10614.617134 = 2p + 10660.617134$$

Compare:

LRT:

Obviously, ThreeRM and 24-RM of OBSM method I are all significant better than hypothesis models in original study.

$$\text{Df} = 1 \quad 2\Delta l = 17.838506 \quad p\text{-value} = 2.405e-005$$

$$\text{Df} = 22 \quad 2\Delta l = 134.057098 \quad p\text{-value} = 0$$

FRM VS 24-RM

$$\text{Df} = 20 \quad 2\Delta l = 20.424914 \quad p\text{-value} = 0.4316$$

AIC: the 24-RM is the best model

a: We calculate several hypothesis models according the original paper implied and only optimal two are shown.

b: more optimal models significant better than ORM is not shown.

CASE 25: Celegent_Lysozymes_lys

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 4B

Previous study:

ORM:

$$\ln L = -2285.487402 \quad \omega_0 = 0.08775 \quad \kappa = 1.46754$$

$$AIC = 2p + 4570.974804$$

TwoRM A: (Branch J)

$$\ln L = -2277.110133 \quad \omega_0 = 0.07059 \quad \omega_1 = 999.00000 \quad \kappa = 1.47874$$

$$AIC = 2p + 2 + 4554.220266 = 2p + 4556.220266$$

TwoRM B: (Branch N)

$$\ln L = -2277.411278 \quad \omega_0 = 0.07774 \quad \omega_1 = 2.74348 \quad \kappa = 1.49366$$

$$AIC = 2p + 2 + 4554.822556 = 2p + 4556.822556$$

FRM:

$$\ln L = -2255.780988$$

$$AIC = 2p + 2 \times 17 + 4511.561976 = 2p + 4545.561976$$

OBSM Method I:

ThreeRM:

$$\ln L = -2267.532752 \quad \omega_0 = 0.06000 \quad \omega_1 = 2.84147 \quad \omega_2 = 999.00000 \quad \kappa = 1.50711$$

$$AIC = 2p + 2 \times 2 + 4535.065504 = 2p + 4539.065504$$

FourRM:

$$\ln L = -2264.778762 \quad \omega_0 = 0.07076 \quad \omega_1 = 3.30886 \quad \omega_2 = 999.00000 \quad \omega_3 = 0.01957 \quad \kappa = 1.52597$$

$$AIC = 2p + 2 \times 3 + 4529.557524 = 2p + 4535.557524$$

Compare:

LRT:

ThreeRM and FourRM of OBSM method I are all significant better than TwoRM A and TwoRM B in original study.

$$\text{ThreeRM vs A Df}=1 \quad 2\Delta l=19.154762 \quad p\text{-value}=1.205e-005$$

$$\text{FourRM vs A Df}=2 \quad 2\Delta l=24.662742 \quad p\text{-value}=4.411e-006$$

FRM VS FourRM

$$\text{Df}=14 \quad 2\Delta l=17.995548 \quad p\text{-value}=0.20698$$

AIC: the FourRM is the best model

a: the alignment may have a little different

b: We calculate several hypothesis models according the original paper implied and only two of them are shown

CASE 26: X_Linked_Gene_Family

^aSequences: according the id in paper

^bAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 2

^cPrevious study:

ORM:

$$\ln L = -2456.267806 \quad \omega_0 = 0.28262 \quad \kappa = 2.70747 \\ AIC = 2p + 4912.535612$$

TwoRM A:

$$\ln L = -2454.333403 \quad \omega_0 = 0.27450 \quad \omega_1 = 999.00000 \quad \kappa = 2.71278 \\ AIC = 2p + 2 + 4908.666806 = 2p + 4910.666806$$

ThreeRM:

$$\ln L = -2453.178179 \quad \omega_0 = 0.26442 \quad \omega_1 = 999.00000 \quad \omega_2 = 0.57204 \quad \kappa = 2.65807 \\ AIC = 2p + 2 * 2 + 4906.356358 = 2p + 4910.356358$$

FourRM:

$$\ln L = -2451.212063 \quad \omega_0 = 0.25053 \quad \omega_1 = 0.56358 \quad \omega_2 = 999.00000 \quad \omega_3 = 0.84446 \quad \kappa = 2.70908 \\ AIC = 2p + 2 * 3 + 4902.424126 = 2p + 4908.424126$$

FRM:

$$\ln L = -2396.264656 \\ AIC = 2p + 2 * 66 + 4792.529312 = 2p + 4924.529312$$

^dOBSM Method I:

TwoRM:

$$\ln L = -2448.331000 \quad \omega_0 = 0.30898 \quad \omega_1 = 0.01533 \quad \kappa = 2.66214 \\ AIC = 2p + 2 + 4896.662 = 2p + 4898.662$$

ThreeRM:

$$\ln L = -2442.218704 \quad \omega_0 = 0.33152 \quad \omega_1 = 0.01505 \quad \omega_2 = 0.03336 \quad \kappa = 2.65807 \\ AIC = 2p + 2 * 2 + 4884.437408 = 2p + 4888.437408$$

16-RM:

$$\ln L = -2413.221621 \quad \omega_0 = 0.25372 \quad \omega_1 = 0.01310 \quad \omega_2 = 0.03481 \quad \kappa = 2.63900 \\ \omega_{3-15} \sim (1.06662 \ 0.03730 \ 0.98333 \ 1.91290 \ 999.00000 \ 0.00010 \ 0.00010 \ 154.56875 \ 0.88120 \ 0.65897 \\ 99.08026 \ 0.00010 \ 0.59155) \\ AIC = 2p + 2 * 15 + 4826.443242 = 2p + 4856.443242$$

Compare:

LRT:

ThreeRM and FourRM of OBSM method I are all significant better than TwoRM A and TwoRM B in original study.

ThreeRM vs A Df=1 $2\Delta l = 19.154762$ p-value=1.205e-005

FourRM vs 16-RM Df=12 $2\Delta l = 75.980884$ p-value=2.396e-011

FRM vs 16-RM

Df=51 $2\Delta l = 33.91393$ p-value= 0.96863

AIC: the 16-RM is the best model

a: some id given by original paper have improved and we select the longest CDS.

b: the alignment may have a little different

c: We calculate several hypothesis models according the original paper implied and only three of them are shown

d: more optimal models significant better than ORM is not shown.

CASE 27: PISTILLATA_like_genes

Sequences: according the id in paper

^aAlignment: MEGA4 (only CDS)

Phylogeny: is congruent with Fig 4B

Previous study:

ORM:

$$\ln L = -1187.751149 \quad \omega_0 = 0.13369 \quad \kappa = 1.53112$$

$$AIC = 2p + 2375.502298$$

TwoRM A:

$$\ln L = -1187.259645 \quad \omega_0 = 0.11798 \quad \omega_1 = 0.16143 \quad \kappa = 1.54170$$

$$AIC = 2p + 2 + 2374.51929 = 2p + 2376.51929$$

ThreeRM B:

$$\ln L = -1187.194399 \quad \omega_0 = 0.11385 \quad \omega_1 = 0.16165 \quad \omega_2 = 0.13723 \quad \kappa = 1.49366$$

$$AIC = 2p + 2 * 2 + 2374.388798 = 2p + 2378.388798$$

FRM:

$$\ln L = -1153.603318$$

$$AIC = 2p + 2 * 111 + 2307.206636 = 2p + 2529.206636$$

OBSM Method I:

TwoRM A:

$$\ln L = -1184.105906 \quad \omega_0 = 0.12624 \quad \omega_1 = 999.00000 \quad \kappa = 1.53211$$

$$AIC = 2p + 2 + 2368.211812 = 2p + 2370.211812$$

Compare:

LRT:

TwoRM of OBSM method I is significant better than hypothesis models in original study.

$$Df=1 \quad 2\Delta l=6.307478 \quad p\text{-value}=1.202e-002$$

$$Df=1 \quad 2\Delta l=6.176986 \quad p\text{-value}=1.294e-002$$

FRM vs TwoRM

$$Df=110 \quad 2\Delta l=61.005176 \quad p\text{-value}=0.99995$$

AIC: the TwoRM is the best model

a: the alignment may have a little different

b: We calculate several hypothesis models according the original paper implied and best two of them are shown

CASE 28: triplicated_alpha_globin_genes

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 4

^bPrevious study:

ORM:

$$\ln L = -1612.072779 \quad \omega_0 = 0.31077 \quad \kappa = 2.31811 \\ AIC = 2p + 3224.145558$$

TwoRM:

$$\ln L = -1598.147350 \quad \omega_0 = 0.1913 \quad \omega_1 = 1.1151 \quad \kappa = 2.36482 \\ AIC = 2p + 2 + 3196.2947 = 2p + 3198.2947$$

ThreeRM:

$$\ln L = -1597.543353 \quad \omega_0 = 0.1902 \quad \omega_1 = 0.8270 \quad \omega_2 = 1.6875 \quad \kappa = 2.37466 \\ AIC = 2p + 2 * 2 + 3195.086706 = 2p + 3199.086706$$

FRM:

$$\ln L = -1580.053550 \\ AIC = 2p + 2 * 15 + 3160.1071 = 2p + 3190.1071$$

OBSM Method I:

TwoRM:

$$\ln L = -1603.272372 \quad \omega_0 = 0.27191 \quad \omega_1 = 999.00000 \quad \kappa = 2.34649 \\ AIC = 2p + 2 + 3206.544744 = 2p + 3208.544744$$

ThreeRM:

$$\ln L = -1593.631244 \quad \omega_0 = 0.20754 \quad \omega_1 = 1.45908 \quad \omega_2 = 999.00000 \quad \kappa = 2.40139 \\ AIC = 2p + 2 * 2 + 3187.262488 = 2p + 3191.262488$$

FourRM:

$$\ln L = -1591.690155 \quad \omega_0 = 0.26197 \quad \omega_1 = 1.45196 \quad \omega_2 = 999.00000 \quad \omega_3 = 0.13052 \quad \kappa = 2.45056 \\ AIC = 2p + 2 * 3 + 3183.38031 = 2p + 3189.38031$$

SixRM:

$$\ln L = -1585.078668 \quad \omega_0 = 0.16188 \quad \omega_1 = 1.47085 \quad \omega_2 = 999.00000 \quad \kappa = 2.43738 \\ \omega_3 = 0.13318 \quad \omega_4 = 0.82982 \\ AIC = 2p + 2 * 5 + 3170.157336 = 2p + 3180.157336$$

Compare:

LRT:

ThreeRM, FourRM and FiveRM of OBSM method I are all significant better than TwoRM and ThreeRM in original study.

$$Df=1 \quad 2\Delta l=11.706396 \quad p\text{-value}=6.229e-004 \\ Df=2 \quad 2\Delta l=24.92937 \quad p\text{-value}=3.861e-006$$

FRM vs SixRM

$$Df=10 \quad 2\Delta l=10.050236 \quad p\text{-value}=0.43609$$

AIC: the SixRM is the best model

a: the sequences may have a little change

b: it seems there're some difference between our results with original paper, this difference may be caused by a parameter change in control file.

CASE 29: Ketoacyl synthase domains Clade I

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -1442.533873 \quad \omega_0 = 0.00516 \quad \kappa = 3.51222$$
$$\text{AIC} = 2p + 2885.067746$$

TwoRM:

$$\ln L = -1439.863723 \quad \omega_0 = 0.00010 \quad \omega_1 = 0.00718 \quad \kappa = 3.47665$$
$$\text{AIC} = 2p + 2 + 2879.727446 = 2p + 2881.727446$$

FRM:

$$\ln L = -1427.795754$$
$$\text{AIC} = 2p + 2 \cdot 29 + 2855.591508 = 2p + 2913.591508$$

^cOBSM Method I:

TwoRM:

$$\ln L = -1437.170953 \quad \omega_0 = 0.00392 \quad \omega_1 = 0.23541 \quad \kappa = 3.58453$$
$$\text{AIC} = 2p + 2 + 2874.341906 = 2p + 2876.341906$$

ThreeRM:

$$\ln L = -1434.465677 \quad \omega_0 = 0.00333 \quad \omega_1 = 0.23891 \quad \omega_2 = 0.26447 \quad \kappa = 3.61716$$
$$\text{AIC} = 2p + 2 \cdot 2 + 2868.931354 = 2p + 2872.931354$$

FiveRM:

$$\ln L = -1431.592282 \quad \omega_0 = 0.00243 \quad \omega_1 = 0.23772 \quad \omega_2 = 0.27729 \quad \kappa = 3.49458$$
$$\omega_3 = 0.00010 \quad \omega_4 = 0.01551$$
$$\text{AIC} = 2p + 2 \cdot 4 + 2863.184564 = 2p + 2871.184564$$

Compare:

LRT:

ThreeRM and FiveRM of OBSM method I are all significant better than TwoRM in original study.

$$\text{Df} = 1 \quad 2\Delta l = 10.796092 \quad \text{p-value} = 1.017\text{e-}003$$

$$\text{Df} = 3 \quad 2\Delta l = 16.542882 \quad \text{p-value} = 8.774\text{e-}004$$

FRM vs FiveRM

$$\text{Df} = 25 \quad 2\Delta l = 7.593056 \quad \text{p-value} = 0.99968$$

AIC: FiveRM is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we try several models according the original paper, but it seems there're some difference

c: the FourRM of OBSM method I have smaller log likelihood -1440.119958.

CASE 30: Ketoacyl synthase domains Clade II

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

Previous study:

ORM:

$$\ln L = -1347.126347 \quad \omega_0 = 0.01994 \quad \kappa = 3.44062 \\ AIC = 2p + 2694.252694$$

TwoRM:

$$\ln L = -1345.232158 \quad \omega_0 = 0.01450 \quad \omega_1 = 0.04716 \quad \kappa = 3.42820 \\ AIC = 2p + 2 + 2690.464316$$

FRM:

$$\ln L = -1336.208517 \\ AIC = 2p + 2 \cdot 14 + 2672.417034 = 2p + 2700.417034$$

^bOBSM Method I:

TwoRM:

$$\ln L = -1344.601473 \quad \omega_0 = 0.01561 \quad \omega_1 = 0.07851 \quad \kappa = 3.45083 \\ AIC = 2p + 2 + 2689.202946 = 2p + 2691.202946$$

ThreeRM:

$$\ln L = -1342.629160 \quad \omega_0 = 0.01886 \quad \omega_1 = 0.07786 \quad \omega_2 = 0.00010 \quad \kappa = 3.45195 \\ AIC = 2p + 2 \cdot 2 + 2685.25832 = 2p + 2689.25832$$

FourRM:

$$\ln L = -1340.193687 \quad \omega_0 = 0.01102 \quad \omega_1 = 0.07873 \quad \omega_2 = 0.00010 \quad \omega_3 = 0.07279 \quad \kappa = 3.43848 \\ AIC = 2p + 2 \cdot 3 + 2680.387374 = 2p + 2686.387374$$

Compare:

LRT:

ThreeRM and FiveRM of OBSM method I are all significant better than TwoRM in original study.

$$\begin{array}{ll} Df=1 & 2\Delta l=5.205996 \quad p\text{-value}=0.0225 \\ Df=2 & 2\Delta l=10.076942 \quad p\text{-value}=6.484e-003 \end{array}$$

FRM vs FourRM

$$Df=1 \quad 2\Delta l=7.97034 \quad p\text{-value}=0.7159$$

AIC: the FourRM is the best model

a: the sequences may have a little change (may be shorter than original paper)

CASE 31: Ketoacyl synthase domains Clade III

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$\ln L = -2969.405876$ $\omega_0 = 0.01015$ $\kappa = 1.87659$
 $AIC = 2p + 5938.811752$

TwoRM: (branch b)

$\ln L = -2966.588079$ $\omega_0 = 0.01106$ $\omega_1 = 0.00032$ $\kappa = 1.94088$
 $AIC = 2p + 2 + 5933.176158 = 2p + 5935.176158$

TwoRM: (clade a)

$\ln L = -2961.581034$ $\omega_0 = 0.00338$ $\omega_1 = 0.02002$ $\kappa = 2.00419$
 $AIC = 2p + 2 + 5923.162068 = 2p + 5925.162068$

FRM:

$\ln L = -2940.835126$
 $AIC = 2p + 2 * 25 + 5881.670252 = 2p + 5931.670252$

^cOBSM Method I:

TwoRM:

$\ln L = -2963.922594$ $\omega_0 = 0.01210$ $\omega_1 = 0.00010$ $\kappa = 1.85404$
 $AIC = 2p + 2 + 5927.845188 = 2p + 5929.845188$

ThreeRM:

$\ln L = -2957.525489$ $\omega_0 = 0.01474$ $\omega_1 = 0.00010$ $\omega_2 = 0.00076$ $\kappa = 1.91676$
 $AIC = 2p + 2 * 2 + 5915.050978 = 2p + 5919.050978$

13-RM:

$\ln L = -2943.893767$ $\omega_0 = 0.00623$ $\omega_1 = 0.00010$ $\omega_2 = 0.00075$ $\kappa = 1.95339$
 $\omega_{3-12} \sim (0.02852 \ 999.00000 \ 999.00000 \ 0.06094 \ 0.00601 \ 0.02369 \ 999.00000 \ 0.00156 \ 999.00000 \ 0.02487)$
 $AIC = 2p + 2 * 12 + 5887.787534 = 2p + 5911.787534$

Compare:

LRT:

ThreeRM and 13-RM of OBSM method I are all significant better than TwoRM (calde meodel) in original study.

$Df = 1$ $2\Delta l = 8.11109$ $p\text{-value} = 4.400e-003$
 $Df = 11$ $2\Delta l = 35.374534$ $p\text{-value} = 2.149e-004$

FRM vs 13-RM

$Df = 13$ $2\Delta l = 6.117282$ $p\text{-value} = 0.9417$

AIC: the 13-RM is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we try six models according the original paper, and only branch model b and clade model a is shown.

c: more optimal models significant better than ORM is not shown.

CASE 32: Ketoacyl synthase domains Clade IV

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

Previous study:

ORM:

$$\ln L = -1572.043286 \quad \omega_0 = 0.02194 \quad \kappa = 2.52282 \\ AIC = 2p + 3144.086572$$

TwoRM:

$$\ln L = -1570.583052 \quad \omega_0 = 0.01887 \quad \omega_1 = 0.07576 \quad \kappa = 2.47106 \\ AIC = 2p + 2 + 3141.166104$$

FRM:

$$\ln L = -1566.875766 \\ AIC = 2p + 2 \cdot 10 + 3133.751532 = 2p + 3153.751532$$

OBSM Method I:

TwoRM:

$$\ln L = -1570.351702 \quad \omega_0 = 0.02856 \quad \omega_1 = 0.00010 \quad \kappa = 2.36845 \\ AIC = 2p + 2 + 3140.703404 = 2p + 3142.703404$$

OBSM Method II:

same with method I

OBSM Method III:

TwoRM (k=0.5):

$$\ln L = -1567.760247 \quad \omega_0 = 0.0347 \quad \omega_1 = 0.0011 \quad \kappa = 2.44028 \\ AIC = 2p + 2 + 3135.520494 = 2p + 3137.520494$$

Compare:

LRT:

TwoRM of OBSM method III is significant better than TwoRM in original study.

$$Df = 1 \quad 2\Delta l = 5.64561 \quad p\text{-value} = 0.0175$$

FRM vs TwoRM

$$Df = 9 \quad 2\Delta l = 1.768962 \quad p\text{-value} = 0.9946$$

AIC: the TwoRM of Method III is the best model

a: the sequences may have a little change (may be shorter than original paper)

CASE 33: Ketoacyl synthase domains Clade V

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -2356.171114 \quad \omega_0 = 0.01622 \quad \kappa = 4.75052 \\ AIC = 2p + 4712.342228$$

TwoRM: (branch a)

$$\ln L = -2355.599899 \quad \omega_0 = 0.01526 \quad \omega_1 = 0.14417 \quad \kappa = 4.81586 \\ AIC = 2p + 2 + 4711.199798 = 2p + 4713.199798$$

TwoRM: (clade a)

$$\ln L = -2353.627003 \quad \omega_0 = 0.00831 \quad \omega_1 = 0.02135 \quad \kappa = 4.81184 \\ AIC = 2p + 2 + 4707.254006 = 2p + 4709.254006$$

FRM:

$$\ln L = -2336.130641 \\ AIC = 2p + 2 * 26 + 4672.261282 = 2p + 4724.261282$$

OBSM Method I:

TwoRM:

$$\ln L = -2348.967332 \quad \omega_0 = 0.02106 \quad \omega_1 = 0.00010 \quad \kappa = 4.70463 \\ AIC = 2p + 2 + 4697.934664 = 2p + 4699.934664$$

FourRM:

$$\ln L = -2345.674448 \quad \omega_0 = 0.01878 \quad \omega_1 = 0.00010 \quad \omega_2 = 74.91367 \quad \omega_3 = 174.69408 \quad \kappa = 4.82045 \\ AIC = 2p + 2 * 3 + 4691.348896 = 2p + 4697.348896$$

SixRM:

$$\ln L = -2342.410437 \quad \omega_0 = 0.01991 \quad \omega_1 = 0.00010 \quad \omega_2 = 67.12320 \quad \kappa = 4.93669 \\ \omega_{3-5} = (212.64653 \ 0.05116 \ 0.00587) \\ AIC = 2p + 2 * 5 + 4684.820874 = 2p + 4694.820874$$

Compare:

LRT:

TwoRM, FourRM and SixRM of OBSM method I are all significant better than TwoRM (calde meodel) in original study.

$$Df=2 \quad 2\Delta l=15.90511 \quad p\text{-value}=3.518e-004 \\ Df=4 \quad 2\Delta l=22.433132 \quad p\text{-value}=1.643e-004$$

FRM vs SixRM

$$Df=21 \quad 2\Delta l=12.559592 \quad p\text{-value}=0.9233$$

AIC: the SixRM is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we try four models according the original paper, and only branch model a and clade model a is shown.

CASE 34: Ketoacyl synthase domains Clade VI

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -1776.430586 \quad \omega_0 = 0.04141 \quad \kappa = 4.22245$$

$$AIC = 2p + 3552.861172$$

TwoRM: (clade b)

$$\ln L = -1768.018146 \quad \omega_0 = 0.12261 \quad \omega_1 = 0.02464 \quad \kappa = 4.70537$$

$$AIC = 2p + 2 + 3536.036292 = 2p + 3538.036292$$

FRM:

$$\ln L = -1738.570051$$

$$AIC = 2p + 2 * 15 + 3477.140102 = 2p + 3507.140102$$

^cOBSM Method I:

TwoRM:

$$\ln L = -1763.087057 \quad \omega_0 = 0.08059 \quad \omega_1 = 0.00408 \quad \kappa = 5.01164$$

$$AIC = 2p + 2 + 3526.174114 = 2p + 3528.174114$$

FourRM:

$$\ln L = -1755.050278 \quad \omega_0 = 0.06845 \quad \omega_1 = 0.00010 \quad \omega_2 = 0.02008 \quad \omega_3 = 999.00000 \quad \kappa = 5.08283$$

$$AIC = 2p + 2 * 3 + 3510.100556 = 2p + 3516.100556$$

10-RM:

$$\ln L = -1738.752202 \quad \omega_0 = 0.04541 \quad \omega_1 = 0.00010 \quad \omega_2 = 0.01973 \quad \kappa = 5.27260$$

$$\omega_{3-9} = (999.00000 \ 1.32052 \ 0.54685 \ 999.00000 \ 0.17221 \ 999.00000 \ 0.00010)$$

$$AIC = 2p + 2 * 9 + 3477.504404 = 2p + 3495.504404$$

Compare:

LRT:

TwoRM, FourRM and 10-RM of OBSM method I are all significant better than TwoRM (calde meodel) in original study.

$$Df=2 \quad 2\Delta l=25.935736 \quad p\text{-value}=2.334e-006$$

$$Df=8 \quad 2\Delta l=58.531888 \quad p\text{-value}=9.039e-010$$

FRM vs 10-RM

$$Df=6 \quad 2\Delta l=0.364302 \quad p\text{-value}=0.9991$$

AIC: the 10-RM is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we try four models according the original paper, and only clade model b is shown.

c: more optimal models significant better than ORM is not shown.

CASE 35: Ketoacyl synthase domains Clade VII

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$\ln L = -6919.600817$ $\omega_0 = 0.04925$ $\kappa = 2.76949$
 $AIC = 2p + 13839.201634$

TwoRM: (clade a)

$\ln L = -6917.240527$ $\omega_0 = 0.05346$ $\omega_1 = 0.03344$ $\kappa = 2.75765$
 $AIC = 2p + 2 + 13834.481054 = 2p + 13836.481054$

TwoRM: (branch a)

$\ln L = -6917.612101$ $\omega_0 = 0.04847$ $\omega_1 = 999.00000$ $\kappa = 2.78387$
 $AIC = 2p + 2 + 13835.224202 = 2p + 13837.224202$

FRM:

$\ln L = -6862.802488$
 $AIC = 2p + 2 * 39 + 13725.604976 = 2p + 13803.604976$

^cOBSM Method I:

TwoRM:

$\ln L = -6906.202348$ $\omega_0 = 0.04641$ $\omega_1 = 0.80833$ $\kappa = 2.80510$
 $AIC = 2p + 2 + 13812.404696 = 2p + 13814.404696$

ThreeRM:

$\ln L = -6902.812662$ $\omega_0 = 0.04821$ $\omega_1 = 0.80835$ $\omega_2 = 0.01364$ $\kappa = 2.80989$
 $AIC = 2p + 2 * 2 + 13805.625324 = 2p + 13809.625324$

13-RM:

$\ln L = -6879.796276$ $\omega_0 = 0.04593$ $\omega_1 = 0.81655$ $\omega_2 = 0.01367$ $\kappa = 2.98419$
 $\omega_{3-12} \sim (0.28053 \ 999.00000 \ 0.00117 \ 0.01360 \ 0.00010 \ 0.01618 \ 0.00010 \ 999.00000 \ 0.18339 \ 45.39811)$
 $AIC = 2p + 2 * 12 + 13759.592552 = 2p + 13783.592552$

Compare:

LRT:

TwoRM, ThreeRM and 13-RM of OBSM method I are all significant better than TwoRM (calde meodel) in original study.

$Df = 1$ $2\Delta l = 28.85573$ $p\text{-value} = 7.797e-008$
 $Df = 11$ $2\Delta l = 74.888502$ $p\text{-value} = 1.424e-011$

FRM vs 13-RM

$Df = 27$ $2\Delta l = 33.987576$ $p\text{-value} = 0.1664$

AIC: the 13-RM is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we try ten models according the original paper, and two of them are shown.

c: more optimal models significant better than ORM is not shown.

CASE 36: Ketoacyl synthase domains Clade VIII

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

Previous study:

ORM:

$$\ln L = -4250.179073 \quad \omega_0 = 0.00844 \quad \kappa = 0.80641 \\ AIC = 2p + 8500.358146$$

TwoRM:

$$\ln L = -4245.816088 \quad \omega_0 = 0.00805 \quad \omega_1 = 999.00000 \quad \kappa = 0.80984 \\ AIC = 2p + 2 + 8491.632176 = 2p + 8493.632176$$

FRM:

$$\ln L = -4215.743015 \\ AIC = 2p + 2 * 43 + 8431.48603 = 2p + 8517.48603$$

^bOBSM Method I:

TwoRM:

$$\ln L = -4245.816088 \quad \omega_0 = 0.00805 \quad \omega_1 = 999.00000 \quad \kappa = 0.80984 \\ AIC = 2p + 2 + 8491.632176 = 2p + 8493.632176$$

ThreeRM:

$$\ln L = -4242.093240 \quad \omega_0 = 0.00744 \quad \omega_1 = 936.71839 \quad \omega_2 = 0.14035 \quad \kappa = 0.80719 \\ AIC = 2p + 2 * 2 + 8484.18648 = 2p + 8488.18648$$

16-RM:

$$\ln L = -4222.098324 \quad \omega_0 = 0.00626 \quad \omega_1 = 0.03564 \quad \omega_2 = 0.17681 \quad \kappa = 0.86957 \\ \omega_{3-15} \sim (0.20547 \ 48.59960 \ 2.33064 \ 30.88344 \ 0.00062 \ 0.00084 \ 0.00012 \ 0.00010 \ 0.13668 \ 46.83536 \\ 0.00010 \ 0.00010 \ 0.00393) \\ AIC = 2p + 2 * 15 + 8444.196648 = 2p + 8474.196648$$

Compare:

LRT:

ThreeRM and 16-RM of OBSM method I are all significant better than TwoRM in original study.

$$Df=1 \quad 2\Delta l=7.445696 \quad p\text{-value}=6.359e-003 \\ Df=14 \quad 2\Delta l=47.435528 \quad p\text{-value}=1.629e-005$$

FRM vs 16-RM

$$Df=28 \quad 2\Delta l=12.710618 \quad p\text{-value}=0.994$$

AIC: the 16-RM is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: more optimal models significant better than ORM is not shown.

CASE 37: Ketoacyl synthase domains Clade IX

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -3801.247660 \quad \omega_0 = 0.03241 \quad \kappa = 2.98811$$

$$AIC = 2p + 7602.49532$$

TwoRM: (branch b)

$$\ln L = -3799.831950 \quad \omega_0 = 0.03418 \quad \omega_1 = 0.00980 \quad \kappa = 2.99947$$

$$AIC = 2p + 2 + 7599.6639 = 2p + 7601.6639$$

TwoRM: (clade c)

$$\ln L = -3790.660485 \quad \omega_0 = 0.01990 \quad \omega_1 = 0.06138 \quad \kappa = 3.04410$$

$$AIC = 2p + 2 + 7581.32097 = 2p + 7583.32097$$

FRM:

$$\ln L = -3770.506192$$

$$AIC = 2p + 2 \times 26 + 7541.012384 = 2p + 7593.012384$$

^cOBSM Method I:

TwoRM:

$$\ln L = -3793.903278 \quad \omega_0 = 0.03598 \quad \omega_1 = 0.00343 \quad \kappa = 2.98488$$

$$AIC = 2p + 2 + 7587.806556 = 2p + 7589.806556$$

ThreeRM:

$$\ln L = -3787.908586 \quad \omega_0 = 0.03335 \quad \omega_1 = 0.00344 \quad \omega_2 = 999.00000 \quad \kappa = 2.99748$$

$$AIC = 2p + 2 \times 2 + 7575.817172 = 2p + 7579.817172$$

10-RM:

$$\ln L = -3774.202946 \quad \omega_0 = 0.04314 \quad \omega_1 = 0.00371 \quad \omega_2 = 999.00000 \quad \kappa = 3.08358$$

$$\omega_{3-9} \sim (0.08615 \ 0.11755 \ 0.00437 \ 0.01750 \ 0.00911 \ 0.01522 \ 0.00977)$$

$$AIC = 2p + 2 \times 9 + 7548.405892 = 2p + 7566.405892$$

Compare:

LRT:

ThreeRM and 10-RM of OBSM method I are all significant better than TwoRM (clade model) in original study.

$$Df=1 \quad 2\Delta l=5.503798 \quad p\text{-value}=1.898e-002$$

$$Df=8 \quad 2\Delta l=32.915078 \quad p\text{-value}=6.380e-005$$

FRM vs 10-RM

$$Df=172 \Delta l=7.393508 \quad p\text{-value}=0.9778$$

AIC: the 10-RM is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we try six models according the original paper, and two of them are shown.

c: more optimal models significant better than ORM is not shown.

CASE 38: Ketoacyl synthase domains Clade X

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM: $\ln L = -5368.317970$ $\omega_0 = 0.02285$ $\kappa = 1.88448$

$AIC = 2p + 10736.63594$

TwoRM: (clade b)

$\ln L = -5350.866731$ $\omega_0 = 0.00635$ $\omega_1 = 0.09546$ $\kappa = 1.90825$

$AIC = 2p + 2 + 10701.733462 = 2p + 10703.733462$

FRM:

$\ln L = -5341.104326$

$AIC = 2p + 2 \times 21 + 10682.208652 = 2p + 10724.208652$

^cOBSM Method I:

TwoRM:

$\ln L = -5361.540123$ $\omega_0 = 0.01938$ $\omega_1 = 0.81478$ $\kappa = 1.89504$

$AIC = 2p + 2 + 10723.080246 = 2p + 10725.080246$

ThreeRM:

$\ln L = -5357.617433$ $\omega_0 = 0.03022$ $\omega_1 = 0.78823$ $\omega_2 = 0.00366$ $\kappa = 1.87018$

$AIC = 2p + 2 \times 2 + 10715.234866 = 2p + 10719.234866$

SevenRM:

$\ln L = -5348.280085$ $\omega_0 = 0.02861$ $\omega_1 = 0.64904$ $\omega_2 = 0.00589$ $\kappa = 1.90195$

$\omega_{3-6} = (0.01076 \ 999.00000 \ 0.19920 \ 0.00213)$

$AIC = 2p + 2 \times 6 + 10696.56017 = 2p + 10708.56017$

^cOBSM Method II:

TwoRM:

It was same with Method I.

ThreeRM:

$\ln L = -5356.377877$ $\omega_0 = 0.01532$ $\omega_1 = 0.81557$ $\omega_2 = 999.00000$ $\kappa = 1.90255$

$AIC = 2p + 2 \times 2 + 10712.755754 = 2p + 10716.755754$

SixRM:

$\ln L = -5346.670985$ $\omega_0 = 0.00621$ $\omega_1 = 0.64300$ $\omega_2 = 999.00000$ $\kappa = 1.90844$

$\omega_{3-5} = (0.20178 \ 0.06858 \ 0.05229)$

$AIC = 2p + 2 \times 5 + 10693.34197 = 2p + 10703.34197$

^dOBSM Method III: ($k=0.5$)

TwoRM:

$\ln L = -5349.174539$ $\omega_0 = 0.00422$ $\omega_1 = 0.10201$ $\kappa = 1.91186$

$AIC = 2p + 2 + 10698.349078 = 2p + 10700.349078$

Compare:

LRT:

In this case, hypothesis TwoRM (clade model) is obviously significant better than some optimal models explored by OBSM Method I or Method II, but is not better than final models, the final optimal model of Method III is better than this hypothesis model, but not significant if set the $df=1$.

$Df=1$ $2\Delta l=3.384384$ $p\text{-value}=6.582e-002$

FRM vs TwoRM

$Df=202$ $\Delta l=16.140426$ $p\text{-value}=0.7078$

AIC: the TwoRM of Method III is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we tried four models according the original paper, and only the best one is shown.

c: more optimal models significant better than ORM is not shown.

d: we also tried set $k=0.2$ and $k=1$, the results also didn't get better, but interestingly, we observed a sudden log likelihood increase in NineRM ($k=1$) and in EightRM ($k=0.2$) when the previous optimal models have very slim log likelihood difference.

CASE 39: Ketoacyl synthase domains Clade XI

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -2436.063975 \quad \omega_0 = 0.00849 \quad \kappa = 1.81073$$

$$AIC = 2p + 4872.12795$$

TwoRM: (clade model)

$$\ln L = -2434.195583 \quad \omega_0 = 0.01402 \quad \omega_1 = 0.00530 \quad \kappa = 1.83238$$

$$AIC = 2p + 2 + 4868.391166 = 2p + 4870.391166$$

FRM:

$$\ln L = -2427.362661$$

$$AIC = 2p + 2 \times 17 + 4854.725322 = 2p + 4888.725322$$

OBSM Method I:

No model is significant better than ORM.

^cOBSM Method II:

FourRM:

$$\ln L = -2430.452590 \quad \omega_0 = 0.00406 \quad \omega_1 = 681.43273 \quad \omega_2 = 926.92276 \quad \omega_3 = 0.06654 \quad \kappa = 1.83103$$

$$AIC = 2p + 2 \times 3 + 4860.90518 = 2p + 4866.90518$$

OBSM Method III: ($\kappa=0.5$)

TwoRM:

$$\ln L = -2430.552363 \quad \omega_0 = 0.00391 \quad \omega_1 = 33.20042 \quad \kappa = 1.83445$$

$$AIC = 2p + 2 + 4861.104726 = 2p + 4863.104726$$

Compare:

LRT:

The hypothesis clade model is not significant better than ORM but is very close, the p value is $5.323e-002$.

FourRM of Method II, TwoRM of Method III of OBSM are all significant better than ORM.

$$\text{FourRM of Method II vs ORM} \quad Df=3 \quad 2\Delta l=11.22277 \quad p\text{-value}=1.058e-002$$

$$\text{TwoRM of Method III vs ORM} \quad Df=1 \quad 2\Delta l=11.023224 \quad p\text{-value}=8.998e-004$$

And FourRM of Method II, TwoRM of Method III of OBSM are all significant better than TwoRM in original study.

$$\text{FourRM of Method II vs TwoRM} \quad Df=2 \quad 2\Delta l=7.485986 \quad p\text{-value}=2.368e-002$$

$$\text{TwoRM of Method III vs TwoRM} \quad Df=1 \quad 2\Delta l=7.28644 \quad p\text{-value}=6.948e-003$$

FRM vs TwoRM

$$Df=16 \quad 2\Delta l=6.379404 \quad p\text{-value}=0.9834$$

AIC: the TwoRM of the Method III is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we try three models according the original paper, and the best model is shown.

c: more optimal models significant better than ORM is not shown.

CASE 40: Ketoacyl synthase domains Clade XII

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -2334.288011 \quad \omega_0 = 0.00388 \quad \kappa = 3.16771$$
$$AIC = 2p + 4668.576022$$

TwoRM: (branch a)

$$\ln L = -2332.734490 \quad \omega_0 = 0.00780 \quad \omega_1 = 0.00027 \quad \kappa = 3.19185$$
$$AIC = 2p + 2 + 4665.46898 = 2p + 4667.46898$$

FRM:

$$\ln L = -2331.698694$$
$$AIC = 2p + 2 \cdot 11 + 4663.397388 = 2p + 4685.397388$$

^cOBSM Method:

TwoRM:

$$\ln L = -2332.734490 \quad \omega_0 = 0.00780 \quad \omega_1 = 0.00027 \quad \kappa = 3.19185$$
$$AIC = 2p + 2 + 4665.46898 = 2p + 4667.46898$$

Compare:

LRT:

In this case, we try all three OBSM methods and no model is found out that significant better than ORM.
The maximum log likelihood in TwoRM suggest by three methods are all congruent with hypothesis model.

$$Df=1 \quad 2\Delta l=3.107042 \quad p\text{-value}=7.795e-002$$

FRM vs FRM

$$Df=102 \Delta l=2.071592 \quad p\text{-value}=0.9957$$

AIC: the TwoRM is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we try four models according the original paper, and the best model is shown.

c: we try all three method, and no model is significant better than ORM.

CASE 41: Hepcidin Gene in mammal

Sequences: according the id in paper

^aAlignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -1310.463212 \quad \omega_0 = 0.34291 \quad \kappa = 2.32131$$

$$AIC = 2p + 2620.926424$$

FRM:

$$\ln L = -1306.357211$$

$$AIC = 2p + 2 \cdot 24 + 2612.714422 = 2p + 2660.714422$$

^cOBSM Method I:

TwoRM:

$$\ln L = -1309.645825 \quad \omega_0 = 0.36984 \quad \omega_1 = 0.18439 \quad \kappa = 2.32462$$

$$AIC = 2p + 2 + 2619.29165 = 2p + 2621.29165$$

^cOBSM Method II:

TwoRM:

$$\ln L = -1308.499025 \quad \omega_0 = 0.3883 \quad \omega_1 = 0.1342 \quad \kappa = 2.44803$$

$$AIC = 2p + 2 + 2616.99805 = 2p + 2618.99805$$

OBSM Method III: ($k=0.5$)

TwoRM:

$$\ln L = -1306.082389 \quad \omega_0 = 0.49775 \quad \omega_1 = 0.14775 \quad \kappa = 2.35330$$

$$AIC = 2p + 2 + 2612.164778 = 2p + 2614.164778$$

Compare:

LRT:

$$\text{TwoRM of Method II vs ORM} \quad Df=1 \quad 2\Delta l=3.928374 \quad p\text{-value}=4.748e-002$$

$$\text{TwoRM of Method III vs TwoRM} \quad Df=1 \quad 2\Delta l=8.761646 \quad p\text{-value}=3.076e-003$$

FRM vs TwoRM

$$Df=23 \quad 2\Delta l=-0.549644 \quad p\text{-value}=1$$

AIC: the TwoRM of Method III is the best model

a: the alignment may have a little different

b: no hypothesis model, and free ratio model is not significant better than ORM.

c: in this case, there're two pair of sequences are same with each other respectively, and the log likelihood of some models are so instable that the optimal two model of Method I and Method II is not the same.

CASE 42: Hepcidin Gene in Pleuronectiformes and Perciformes

Sequences: according the id in paper

^aAlignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -1692.468704 \quad \omega_0 = 0.56045 \quad \kappa = 2.55100$$

$$AIC = 2p + 3384.937408$$

FRM:

$$\ln L = -1659.752558$$

$$AIC = 2p + 2 \cdot 60 + 3319.505116 = 2p + 3439.505116$$

^cOBSM Method I:

TwoRM:

$$\ln L = -1688.457755 \quad \omega_0 = 0.59137 \quad \omega_1 = 0.00010 \quad \kappa = 2.55587$$

$$AIC = 2p + 2 + 3376.91551 = 2p + 3378.91551$$

ThreeRM:

$$\ln L = -1685.432393 \quad \omega_0 = 0.56484 \quad \omega_1 = 0.00010 \quad \omega_2 = 999.00000 \quad \kappa = 2.55776$$

$$AIC = 2p + 2 \cdot 2 + 3370.864786 = 2p + 3374.864786$$

SevenRM(final optimal model):

$$\ln L = -1675.308796 \quad \omega_0 = 0.65674 \quad \omega_1 = 0.00010 \quad \omega_2 = 999.00000 \quad \kappa = 2.56295$$

$$\omega_{2-6} \sim (0.13097 \ 0.13689 \ 0.00010 \ 2.60183)$$

$$AIC = 2p + 2 \cdot 6 + 3350.617592 = 2p + 3362.617592$$

Compare:

LRT:

TwoRM, ThreeRM and SevenRM of OBSM Method I are all significant better than ORM.

$$Df = 1 \quad 2\Delta l = 8.021898 \quad p\text{-value} = 4.622e-003$$

$$Df = 2 \quad 2\Delta l = 14.072622 \quad p\text{-value} = 8.794e-004$$

$$Df = 6 \quad 2\Delta l = 34.319816 \quad p\text{-value} = 5.835e-006$$

FRM vs SevenRM

$$Df = 54 \quad 2\Delta l = 31.112476 \quad p\text{-value} = 0.9947$$

AIC: the SevenRM is the best model

a: the alignment may have a little different

b: no hypothesis model, and free ratio model is not significant better than ORM.

c: more optimal models significant better than ORM is not shown.

CASE 43: SRPX2_gene
Sequences: according the id in paper
Alignment: MEGA4
Phylogeny: is congruent with Fig 4

Previous study:

ORM:

$$\ln L = -2172.869055 \quad \omega_0 = 0.05323 \quad \kappa = 27.93800 \\ AIC = 2p + 4345.73811$$

TwoRM: (Human lineage)

$$\ln L = -2170.444813 \quad \omega_0 = 0.04253 \quad \omega_1 = 999.00000 \quad \kappa = 27.96774 \\ AIC = 2p + 2 + 4340.889626 = 2p + 4342.889626$$

FRM:

$$\ln L = -2166.238379 \\ AIC = 2p + 2 \cdot 10 + 4332.476758 = 2p + 4352.476758$$

OBSM Method I:

Same with hypothesis model.

OBSM Method II:

FourRM:

$$\ln L = -2167.312563 \quad \omega_0 = 0.02178 \quad \omega_1 = 999.00000 \quad \omega_2 = 0.90182 \quad \omega_3 = 0.49328 \quad \kappa = 28.13717 \\ AIC = 2p + 2 \cdot 3 + 4334.625126 = 2p + 4340.625126$$

OBSM Method III: ($\kappa=0.5$)

TwoRM:

$$\ln L = -2167.835109 \quad \omega_0 = 0.02177 \quad \omega_1 = 0.96753 \text{ (red arrow)} \quad \kappa = 28.14329 \\ AIC = 2p + 2 + 4335.670218 = 2p + 4337.670218$$

Compare:

LRT:

FourRM of Method II, TwoRM of Method III of OBSM are all significant better than hypothesis model.

FourRM of Method II vs TwoRM	Df=2	$2\Delta l = 6.2645$	p-value=4.362e-002
------------------------------	------	----------------------	--------------------

TwoRM of Method III vs TwoRM	Df=1	$2\Delta l = 5.219408$	p-value=2.234e-002
------------------------------	------	------------------------	--------------------

FRM vs TwoRM

Df=9	$2\Delta l = 3.19346$	p-value= 0.9561
------	-----------------------	-----------------

AIC: the TwoRM of Method III is the best model

CASE 44: Cholesterol Metabolism Gene

Sequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 3

Previous study:

ORM:

$$\ln L = -1377.506105 \quad \omega_0 = 0.38583 \quad \kappa = 8.01562$$

$$AIC = 2p + 2755.01221$$

^aTwoRM:

$$\ln L = -1374.183004 \quad \omega_0 = 0.35223 \quad \omega_g = 999.00000 \quad \kappa = 7.99858$$

$$AIC = 2p + 2 + 2748.366008 = 2p + 2750.366008$$

FRM:

$$\ln L = -1360.065636$$

$$AIC = 2p + 2 \cdot 24 + 2720.131272 = 2p + 2768.131272$$

^bOBSM Method I:

TwoRM:

Same with hypothesis model

ThreeRM:

$$\ln L = -1371.448491 \quad \omega_0 = 0.32659 \quad \omega_1 = 999.00000 \quad \omega_2 = 999.00000 \quad \kappa = 7.99851$$

$$AIC = 2p + 2 \cdot 2 + 2742.896982 = 2p + 2746.896982$$

SevenRM:

$$\ln L = -1364.613059 \quad \omega_0 = 0.49116 \quad \omega_1 = 999.00000 \quad \omega_2 = 999.00000 \quad \kappa = 7.99626$$

$$\omega_{3-6} = (0.09249 \ 0.06185 \ 0.00010 \ 0.15082)$$

$$AIC = 2p + 2 \cdot 6 + 2729.226118 = 2p + 2741.226118$$

Compare:

LRT:

ThreeRM and SevenRM of OBSM Method I are all significant better than hypothesis.

$$Df=1 \quad 2\Delta l=5.469026 \quad p\text{-value}=1.936e-002$$

$$Df=5 \quad 2\Delta l=19.13989 \quad p\text{-value}=1.810e-003$$

FRM vs SevenRM

$$Df=18 \quad 2\Delta l=9.094846 \quad p\text{-value}=0.9575$$

AIC: the SevenRM is the best model

a: we try several models according the original paper, and the best model is shown.

b: more optimal models significant better than hypothesis models are not shown.

CASE 45: MOXD2

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 5A

^bPrevious study:

ORM:

$$\ln L = -1437.821252 \quad \omega_0 = 0.22956 \quad \kappa = 7.46279$$

$$AIC = 2p + 2875.642504$$

TwoRM:

$$\ln L = -1432.478425 \quad \omega_0 = 0.17723 \quad \omega_1 = 2.84197 \quad \kappa = 7.49018$$

$$AIC = 2p + 2 + 2864.95685 = 2p + 2866.95685$$

FRM:

$$\ln L = -1430.473242$$

$$AIC = 2p + 2 \cdot 9 + 2860.946484 = 2p + 2878.946484$$

OBSM Method I, II and III:

TwoRM:

Same with hypothesis model

Compare:

LRT:

All the three methods and the hypothesis model suggest the same model.

FRM vs TwoRM

$$Df = 8 \quad 2\Delta l = 4.010366 \quad p\text{-value} = 0.8562$$

AIC: the best model is the TwoRM

a: the id given by original paper is too confused to get identical sequences and the data we got have some changes (have some region missing and shorter than original paper)

b: we try several models according the original paper, and the best model is shown.

CASE 46: S100A15A

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -316.981093 \quad \omega_0 = 0.10216 \quad \kappa = 4.55632$$

$$AIC = 2p + 633.962186$$

TwoRM:

$$\ln L = -315.411448 \quad \omega_0 = 0.07520 \quad \omega_G = 0.53909 \quad \kappa = 4.60778$$

$$AIC = 2p + 2 + 630.822896 = 2p + 632.822896$$

TwoRM:

$$\ln L = -315.562527 \quad \omega_0 = 0.07124 \quad \omega_H = 0.34965 \quad \kappa = 4.66719$$

$$AIC = 2p + 2 + 631.125054 = 2p + 633.125054$$

FRM:

$$\ln L = -310.909248$$

$$AIC = 2p + 2 \cdot 9 + 621.818496 = 2p + 639.818496$$

OBSM Method I:

TwoRM:

$$\ln L = -314.968454 \quad \omega_0 = 0.13024 \quad \omega_1 = 0.00010 \quad \kappa = 4.45714$$

$$AIC = 2p + 2 + 629.936908 = 2p + 631.936908$$

OBSM Method II:

Same with Method I

OBSM Method III: ($\kappa=0.5$)

TwoRM:

$$\ln L = -311.941308 \quad \omega_0 = 0.2752 \quad \omega_1 = 0.0001 \quad \kappa = 4.58316$$

$$AIC = 2p + 2 + 623.882616 = 2p + 625.882616$$

Compare:

LRT:

TwoRM of OBSM Method III is significant better than hypothesis model.

$$Df=1 \quad 2\Delta l=11.023224 \quad p\text{-value}=8.428e-003$$

FRM vs TwoRM

$$Df=8 \quad 2\Delta l=2.06412 \quad p\text{-value}=0.97898$$

AIC: the TwoRM of Method III is the best model

a: the sequences may have a little change (may be shorter than original paper)

b: we try three models according the original paper, and the best model is shown.

c: more optimal models significant better than ORM is not shown.

CASE 47: Penaeidin antimicrobial peptides

Sequences: according the id in paper

^aAlignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$\ln L = -1493.075007$ $\omega_0 = 0.79029$ $\kappa = 1.64528$

$AIC = 2p + 2986.150014$

FRM:

$\ln L = -1461.569406$

$AIC = 2p + 2 \cdot 69 + 2923.138812 = 2p + 3061.138812$

^cOBSM Method I:

TwoRM:

$\ln L = -1485.990663$ $\omega_0 = 0.92749$ $\omega_1 = 0.00010$ $\kappa = 1.66723$

$AIC = 2p + 2 + 2971.981326 = 2p + 2973.981326$

FourRM:

$\ln L = -1482.637071$ $\omega_0 = 0.86416$ $\omega_1 = 0.00010$ $\omega_2 = 999.00000$ $\kappa = 1.66892$

$AIC = 2p + 2 \cdot 3 + 2965.274142 = 2p + 2971.274142$

EightRM:

$\ln L = -1473.714700$ $\omega_0 = 1.05542$ $\omega_1 = 0.00010$ $\omega_2 = 999.00000$ $\kappa = 1.65460$

$\omega_{3-7} = (0.25444 \ 0.00010 \ 999.00000 \ 0.00010 \ 0.21093)$

$AIC = 2p + 2 \cdot 7 + 2947.4294 = 2p + 2961.4294$

Compare:

LRT:

TwoRM FourRM and Eight of OBSM Method I are all significant better than ORM.

Df=1 $2\Delta l = 14.168688$ p-value=1.671e-004

Df=3 $2\Delta l = 20.875872$ p-value=1.117e-004

Df=7 $2\Delta l = 38.720614$ p-value=2.210e-006

FRM vs EightRM

Df=62 $2\Delta l = 24.290588$ p-value=0.999995

AIC: the EightRM is the best model

a: the alignment may have a little different

b: no model is significant better than ORM in original study

c: more optimal models significant better than ORM is not shown.

CASE 48: NYD SP12

Sequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 1

Previous study:

ORM:

$$\ln L = -3222.614087 \quad \omega_0 = 0.62004 \quad \kappa = 4.50420 \\ AIC = 2p + 6445.228174$$

TwoRM:

$$\ln L = -3214.198480 \quad \omega_0 = 0.4589 \quad \omega_1 = 4.8091 \quad \kappa = 4.51495 \\ AIC = 2p + 2 + 6428.39696 = 2p + 6430.39696$$

FRM:

$$\ln L = -3209.257860 \\ AIC = 2p + 2 \cdot 11 + 6418.51572 = 2p + 6440.51572$$

OBSM Method I:

TwoRM:

$$\ln L = -3217.680711 \quad \omega_0 = 0.54921 \quad \omega_1 = 999.00000 \quad \kappa = 4.50692 \\ AIC = 2p + 2 + 6435.361422 = 2p + 6437.361422$$

ThreeRM:

$$\ln L = -3215.065138 \quad \omega_0 = 0.49872 \quad \omega_1 = 999.00000 \quad \omega_2 = 3.56377 \quad \kappa = 4.51073 \\ AIC = 2p + 2 \cdot 2 + 6430.130276 = 2p + 6434.130276$$

OBSM Method II:

Same with Method I

OBSM Method III: ($\kappa=0.5$)

TwoRM:

$$\ln L = -3213.688700 \quad \omega_0 = 0.4524 \quad \omega_1 = 4.9861 \quad \kappa = 4.51422 \\ AIC = 2p + 2 + 6427.3774 = 2p + 6429.3774$$

FourRM:

$$\ln L = -3209.920138 \quad \omega_0 = 0.2135 \quad \omega_1 = 4.8632 \quad \omega_2 = 0.8290 \quad \omega_3 = 0.0001 \quad \kappa = 4.51602 \\ AIC = 2p + 2 \cdot 3 + 6419.840276 = 2p + 6425.840276$$

Compare:

LRT:

FourRM of OBSM Method III are all significant better than hypothesis model

$$Df=2 \quad 2\Delta l=8.556684 \quad p\text{-value}=1.387e-002$$

FRM vs FourRM

$$Df=8 \quad 2\Delta l=1.324556 \quad p\text{-value}=0.9952$$

AIC: the FourRM of Method III is the best model

a: we try three models according the original paper, and the best model is shown.

CASE 49: Myxovirus resistance gene

Sequences: according the id in paper

^aAlignment: MEGA4

Phylogeny: is congruent with Fig 1

^bPrevious study:

ORM:

$$\ln L = -11761.353476 \quad \omega_0 = 0.30812 \quad \kappa = 2.40593$$

$$AIC = 2p + 23522.706952$$

TwoRM: (Branch duck)

$$\ln L = -11743.731336 \quad \omega_0 = 0.27584 \quad \omega_1 = 30.81443 \quad \kappa = 2.42547$$

$$AIC = 2p + 2 + 23487.462672 = 2p + 23489.462672$$

TwoRM: (Branch chicken)

$$\ln L = -11734.736366 \quad \omega_0 = 0.26402 \quad \omega_1 = 1.24721 \quad \kappa = 2.40555$$

$$AIC = 2p + 2 + 23469.472732 = 2p + 23471.472732$$

FRM:

$$\ln L = -11697.659375$$

$$AIC = 2p + 2 \times 21 + 23395.31875 = 2p + 23437.31875$$

^cOBSM Method I:

TwoRM:

Same with hypothesis model (chicken)

ThreeRM:

$$\ln L = -11730.837896 \quad \omega_0 = 0.25735 \quad \omega_1 = 0.74941 \quad \omega_2 = 0.70825 \quad \kappa = 2.41932$$

$$AIC = 2p + 2 \times 2 + 23461.675792 = 2p + 23465.675792$$

TenRM:

$$\ln L = -11703.168029 \quad \omega_0 = 0.28169 \quad \omega_1 = 0.74325 \quad \omega_2 = 0.73485 \quad \kappa = 2.45458$$

$$\omega_{3-7} = (0.00010 \ 0.27599 \ 0.14481 \ 0.15634 \ 0.50811 \ 0.53044 \ 0.06674)$$

$$AIC = 2p + 2 \times 9 + 23406.336058 = 2p + 23424.336058$$

Compare:

LRT:

ThreeRM and EightRM of OBSM Method I are all significant better than hypothesis model (chicken).

$$Df=1 \quad 2\Delta l=7.79694 \quad p\text{-value}=5.233e-003$$

$$Df=8 \quad 2\Delta l=63.136674 \quad p\text{-value}=1.126e-010$$

FRM vs TenRM

$$Df=122 \Delta l=11.017308 \quad p\text{-value}=0.5274$$

AIC: the TenRM is the best model

a: the alignment may have a little different

b: we try several models according the original paper, and the best model is shown.

c: more optimal models significant better than hypothesis models are not shown.

CASE 50: Last Case

^aSequences: according the id in paper

Alignment: MEGA4

Phylogeny: is congruent with Fig 2

^bPrevious study:

ORM:

$$\ln L = -9972.141967 \quad \omega_0 = 0.11693 \quad \kappa = 1.44964$$

$$AIC = 2p + 19944.283934$$

TwoRM: (Branch GRCD1)

$$\ln L = -9964.958754 \quad \omega_0 = 0.11208 \quad \omega_1 = 0.35186 \quad \kappa = 1.44837$$

$$AIC = 2p + 2 + 19929.917508 = 2p + 19931.917508$$

TwoRM: (Branch VvMADS4)

$$\ln L = -9963.962501 \quad \omega_0 = 0.12148 \quad \omega_1 = 0.02850 \quad \kappa = 1.44404$$

$$AIC = 2p + 2 + 19927.925002 = 2p + 19929.925002$$

FRM:

$$\ln L = -9905.929279$$

$$AIC = 2p + 2 * 53 + 19811.858558 = 2p + 19917.858558$$

^cOBSM Method I:

TwoRM:

Same with hypothesis model (Branch VvMADS4)

ThreeRM:

$$\ln L = -9957.697019 \quad \omega_0 = 0.11654 \quad \omega_1 = 0.02935 \quad (\text{VvMADS4}) \quad \omega_2 = 0.33541 \quad (\text{GRCD1}) \quad \kappa = 1.44278$$

$$AIC = 2p + 2 * 2 + 19915.394038 = 2p + 19919.394038$$

16-RM:

$$\ln L = -9920.246557 \quad \omega_0 = 0.11331 \quad \omega_1 = 0.03090 \quad \omega_2 = 0.34508 \quad \kappa = 1.42649$$

$$\omega_{3-15} \sim (0.01912 \ 2.74301 \ 0.28361 \ 0.26519 \ 0.05442 \ 0.58769 \ 0.02559 \ 0.01744 \ 0.03899 \ 0.06861 \\ 0.17050 \ 0.36419 \ 0.19309)$$

$$AIC = 2p + 2 * 15 + 19840.493114 = 2p + 19870.493114$$

Compare:

LRT:

ThreeRM of 15-RM of OBSM Method I are all significant better than hypothesis model

$$Df = 3 \quad 2\Delta l = 12.530964 \quad p\text{-value} = 5.769e-003$$

$$Df = 14 \quad 2\Delta l = 87.431888 \quad p\text{-value} = 4.321e-013$$

FRM vs 15-RM

$$Df = 38 \quad 2\Delta l = 28.634556 \quad p\text{-value} = 0.8642$$

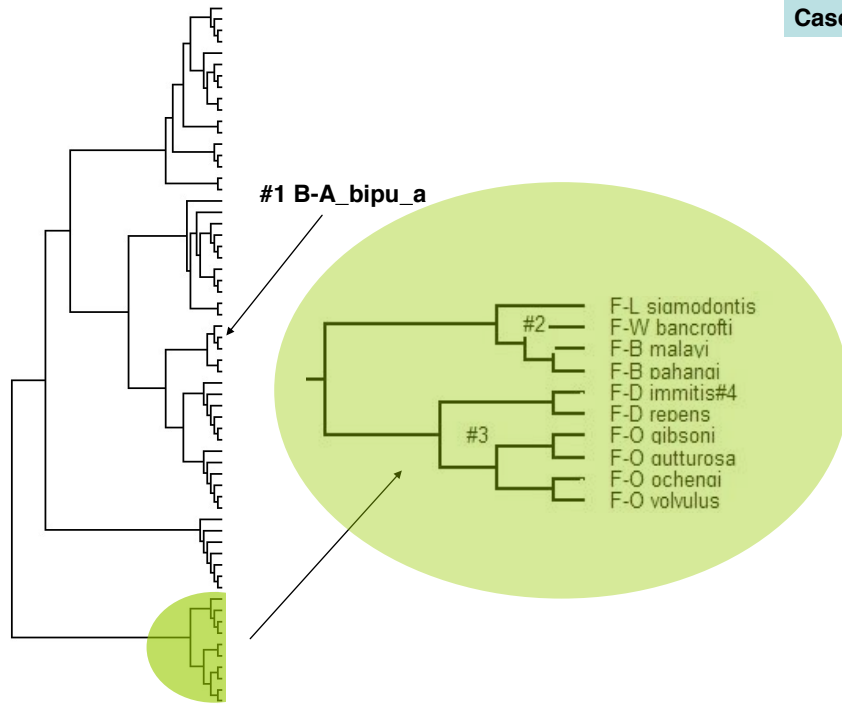
AIC: the 16-RM is the best model

a: the sequences may have a little change (some id given by original paper have changed)

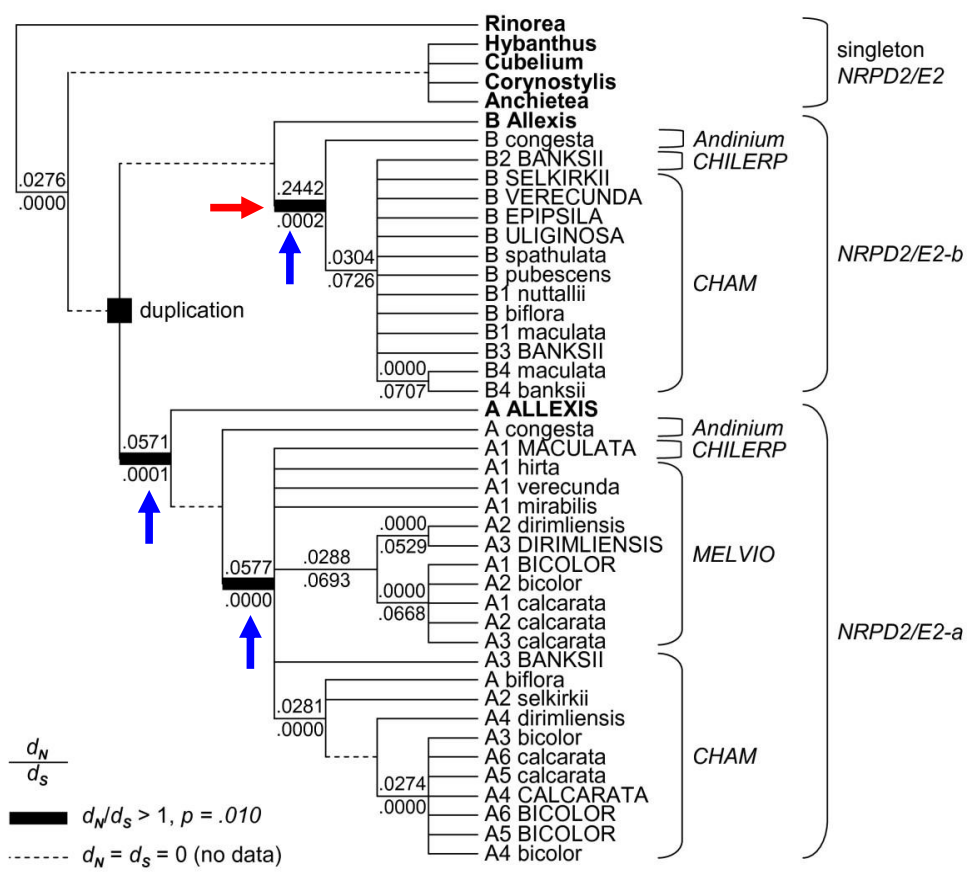
b: we try several models according the original paper, and the best two model are shown.

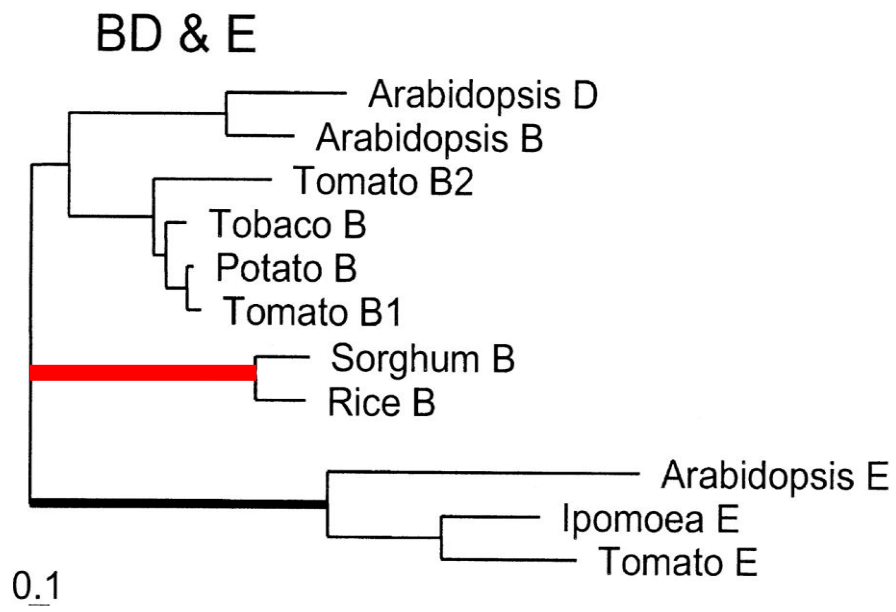
c: more optimal models significant better than ORM is not shown.

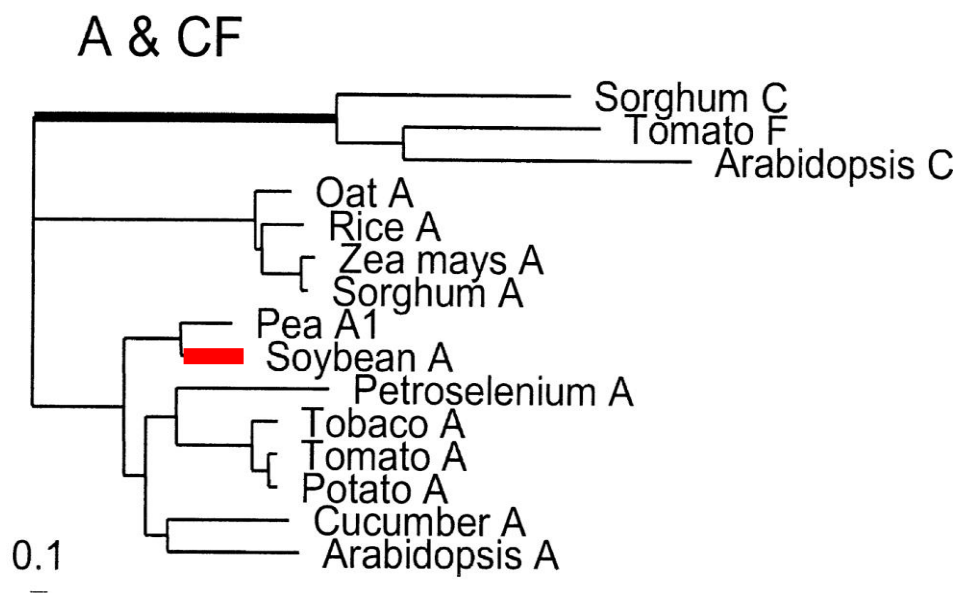
Case 1



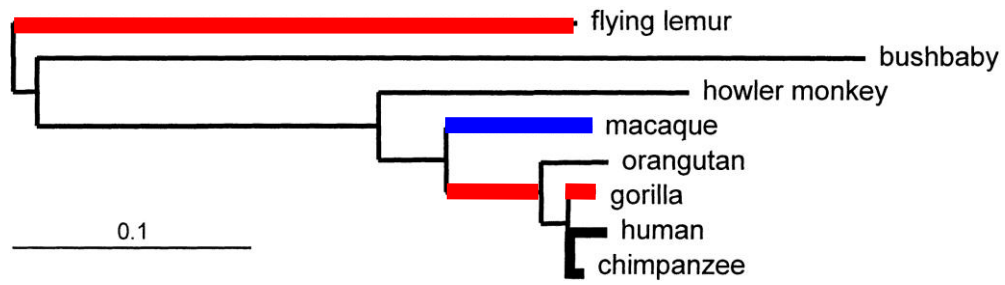
Case 2



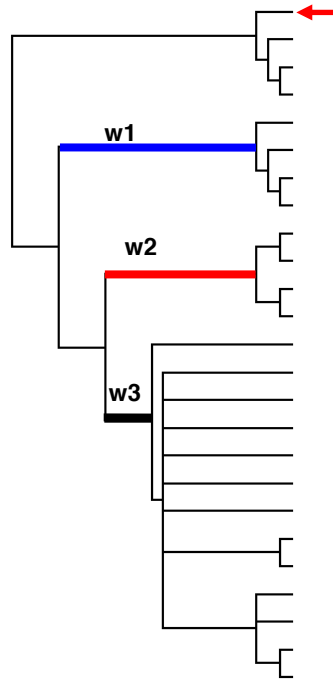




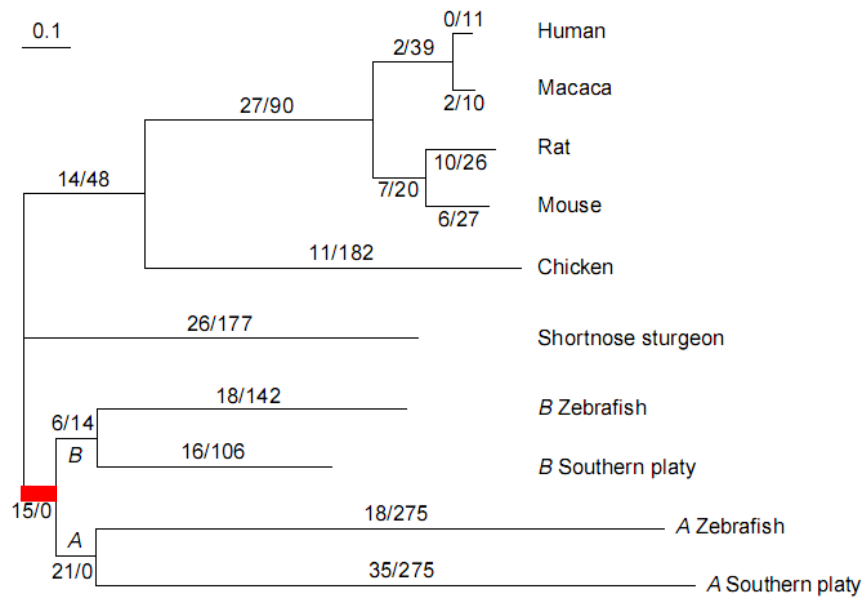
Case 5



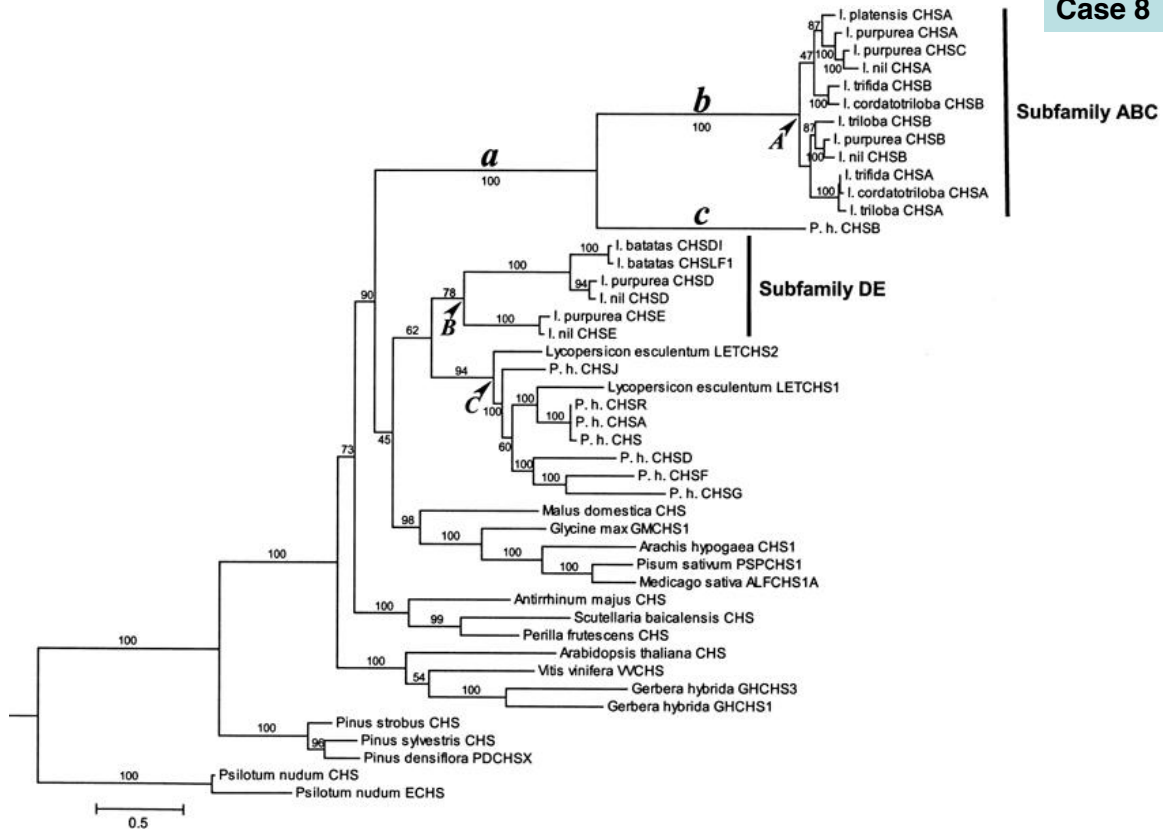
Case 6



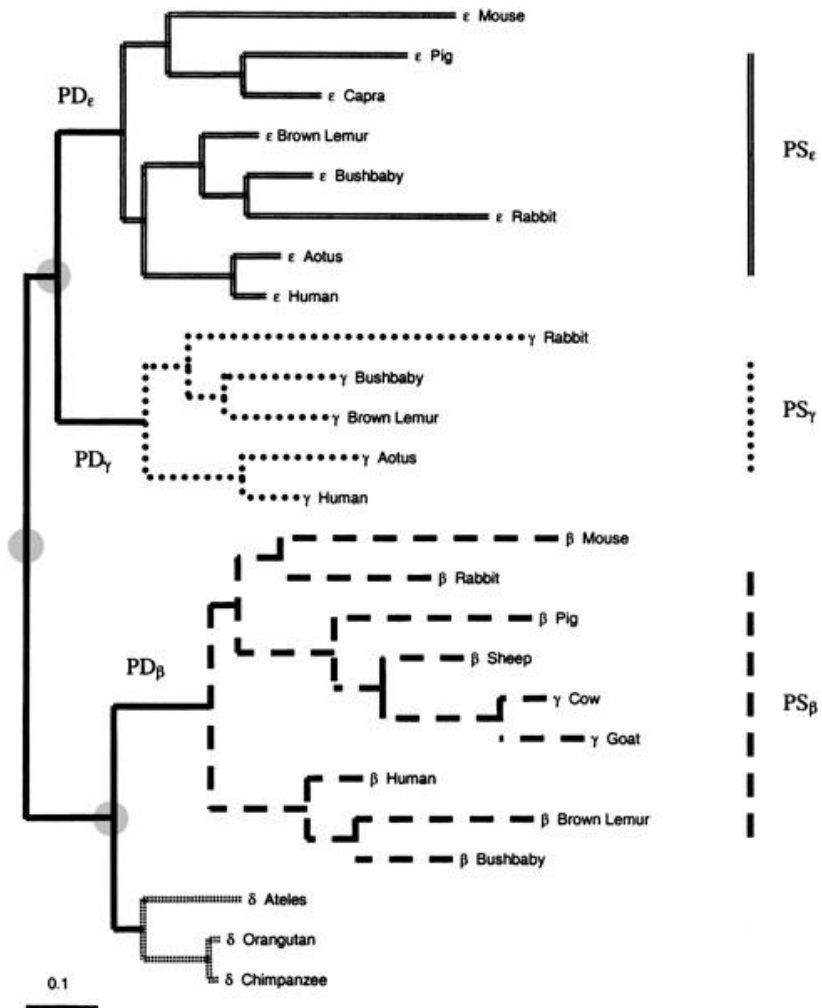
Case 7



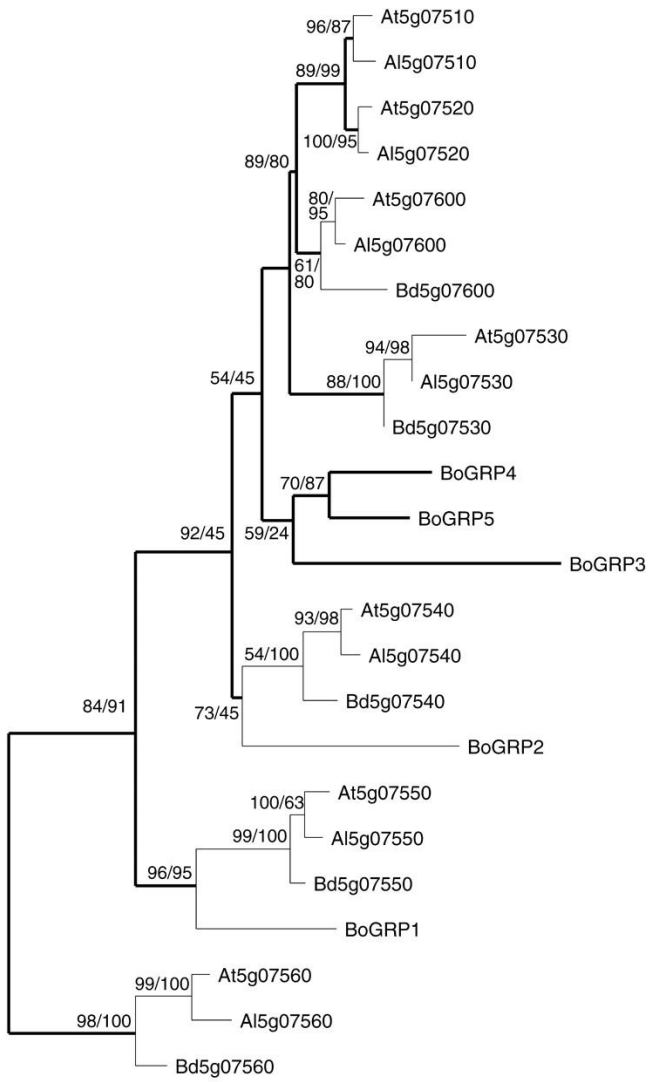
Case 8



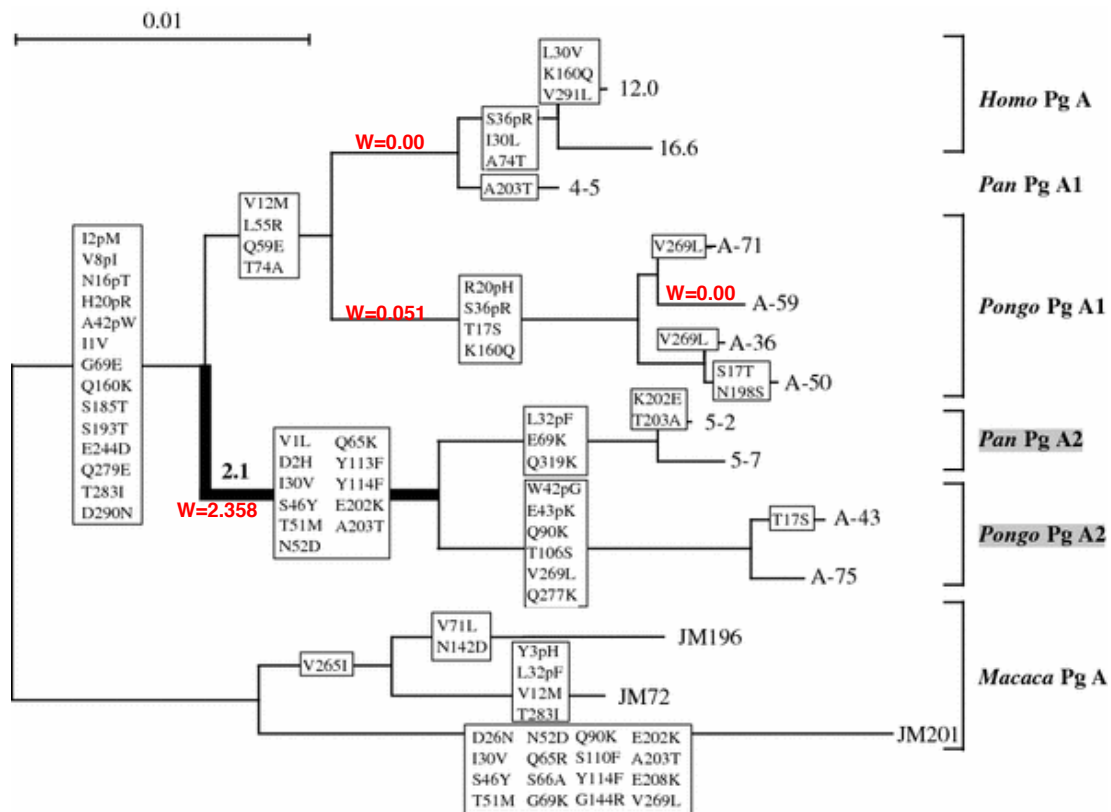
Case 9



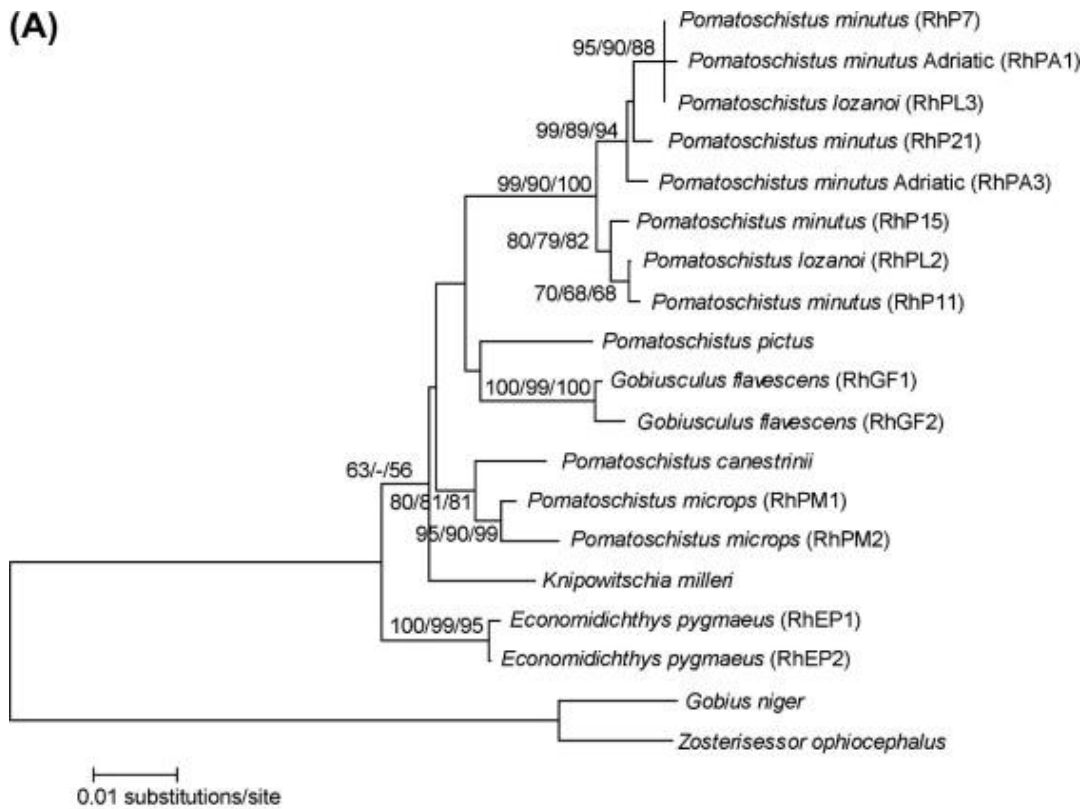
Case 10



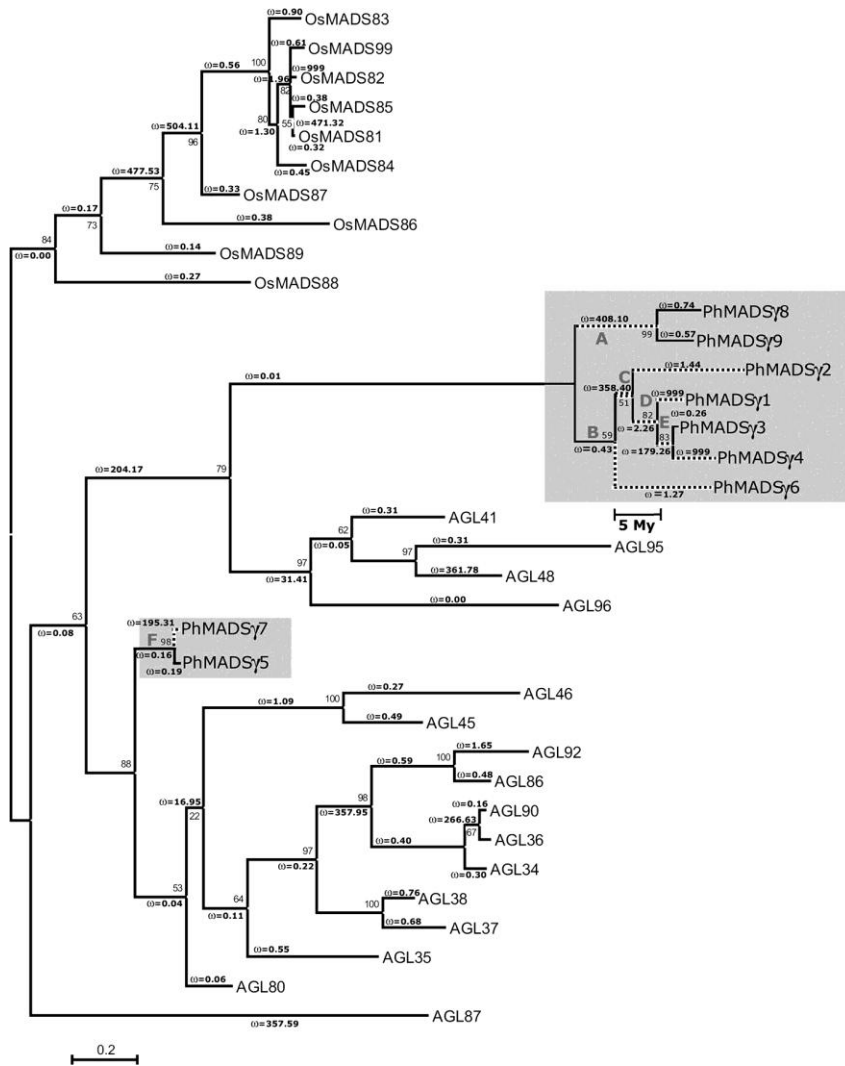
0.2

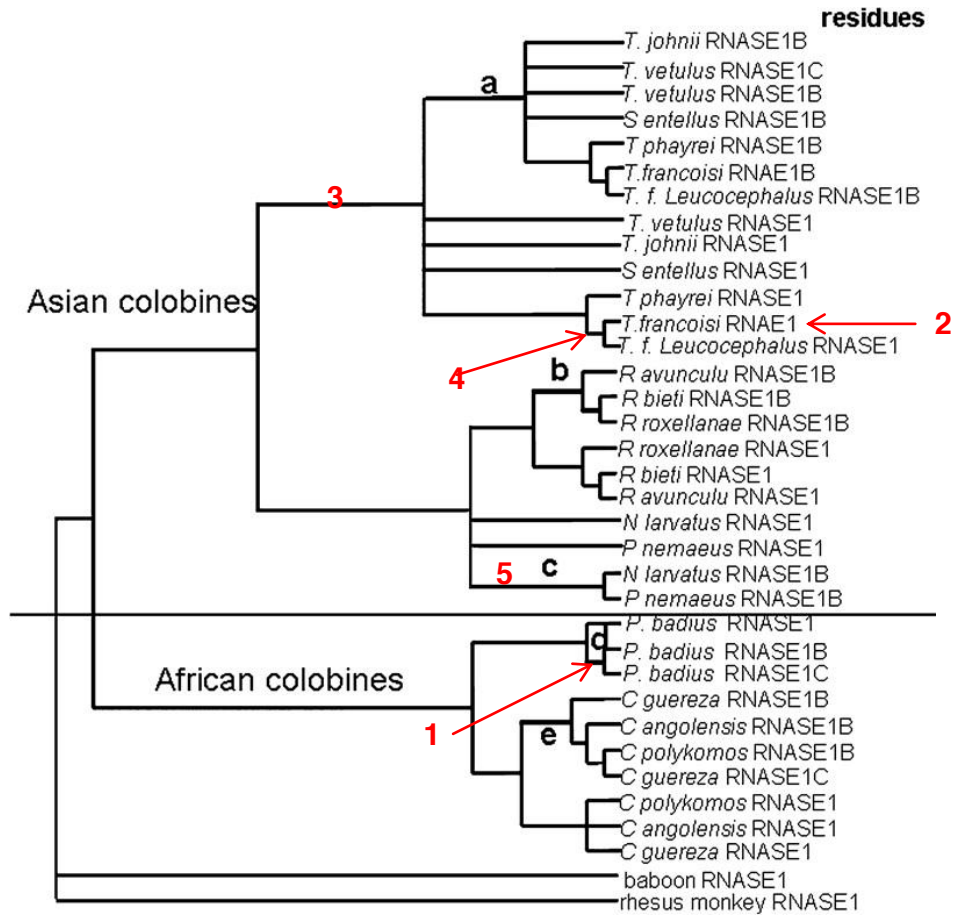


(A)

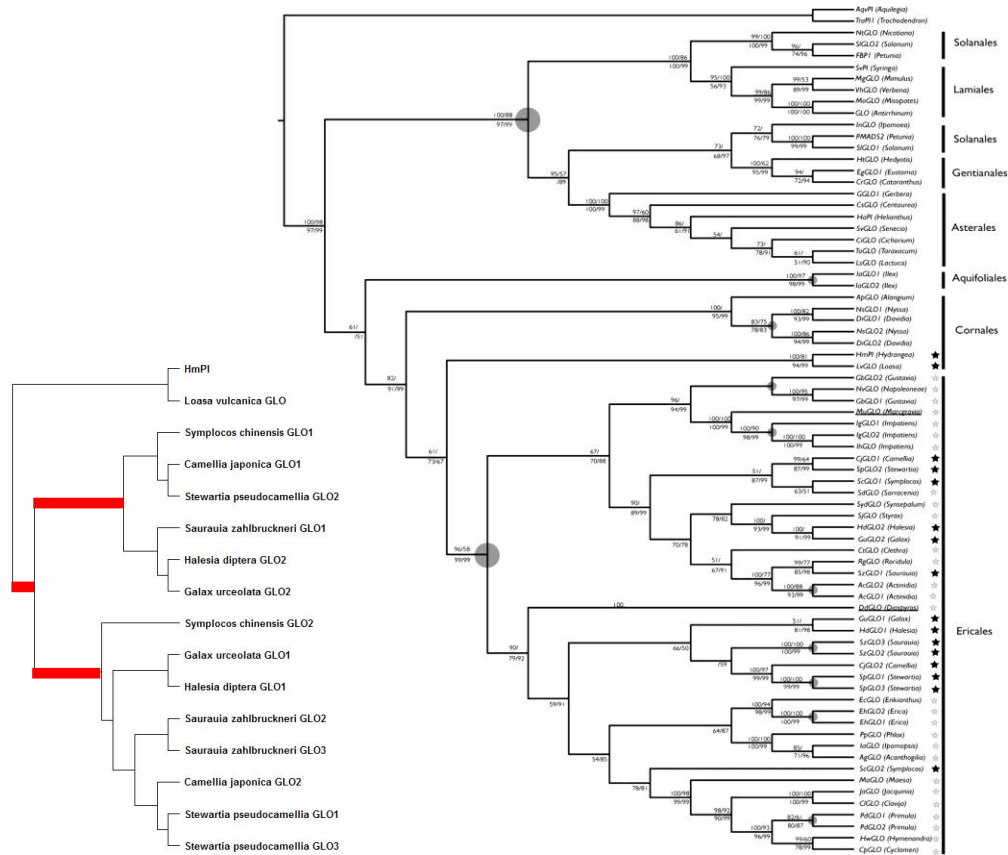


Case 13



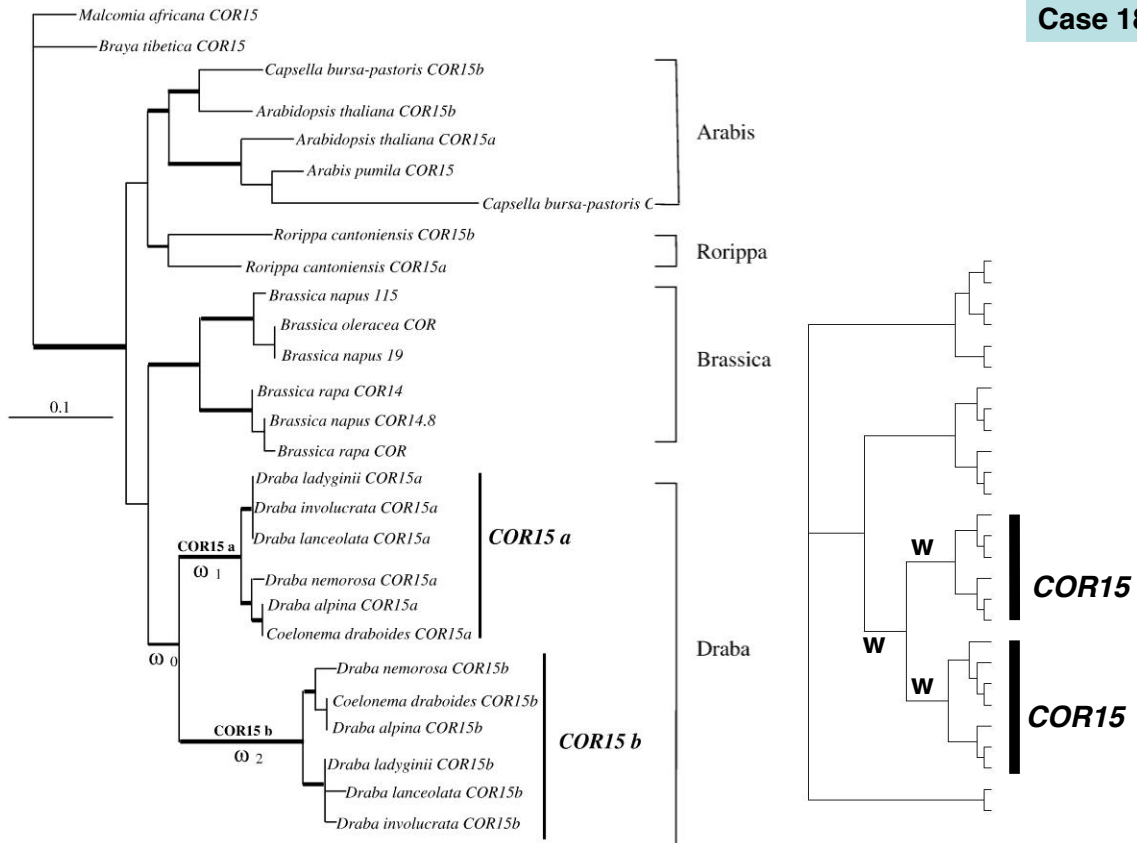


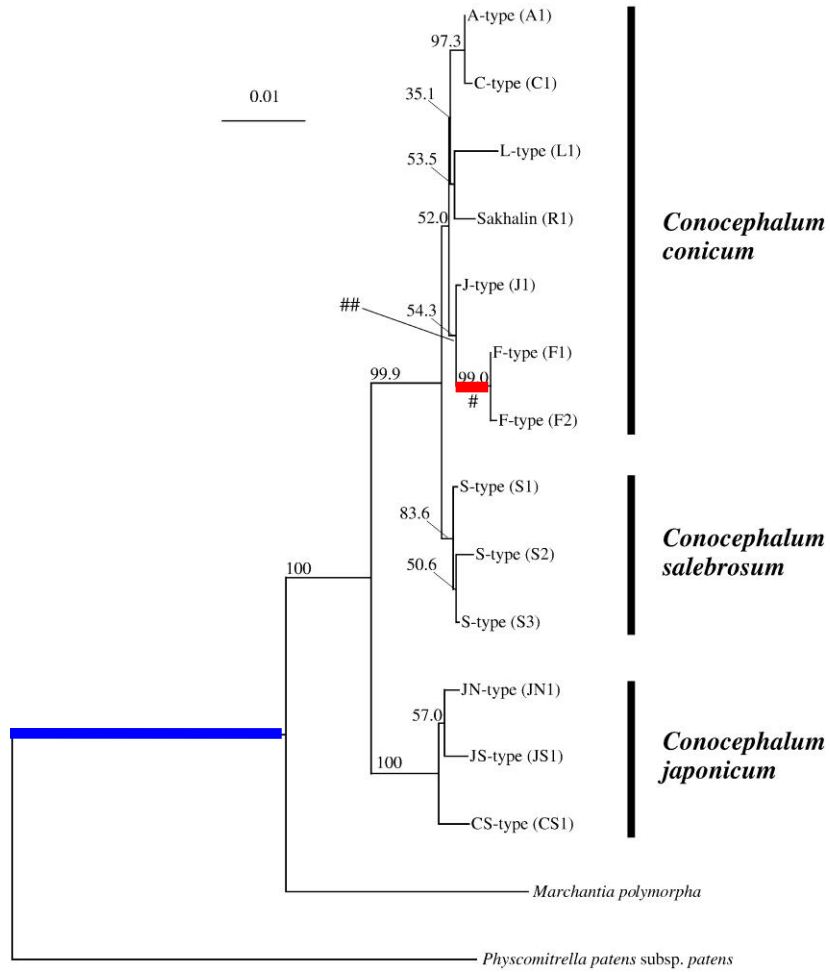
Case 15



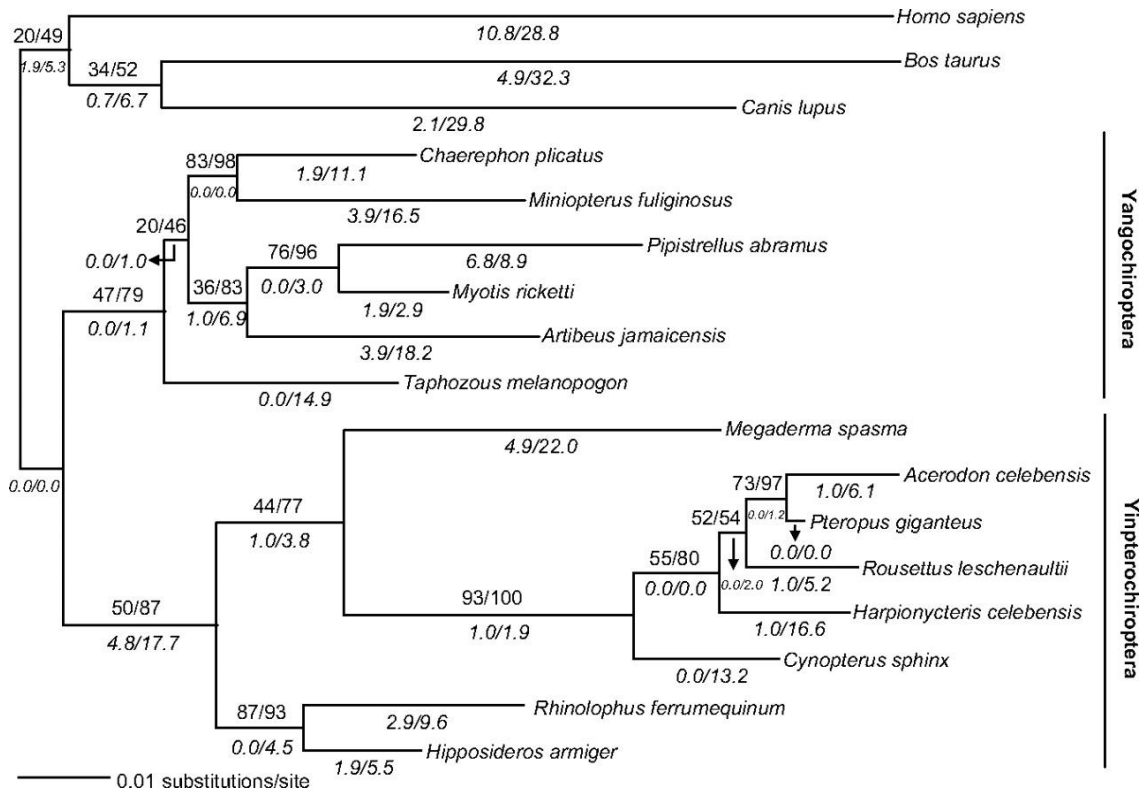




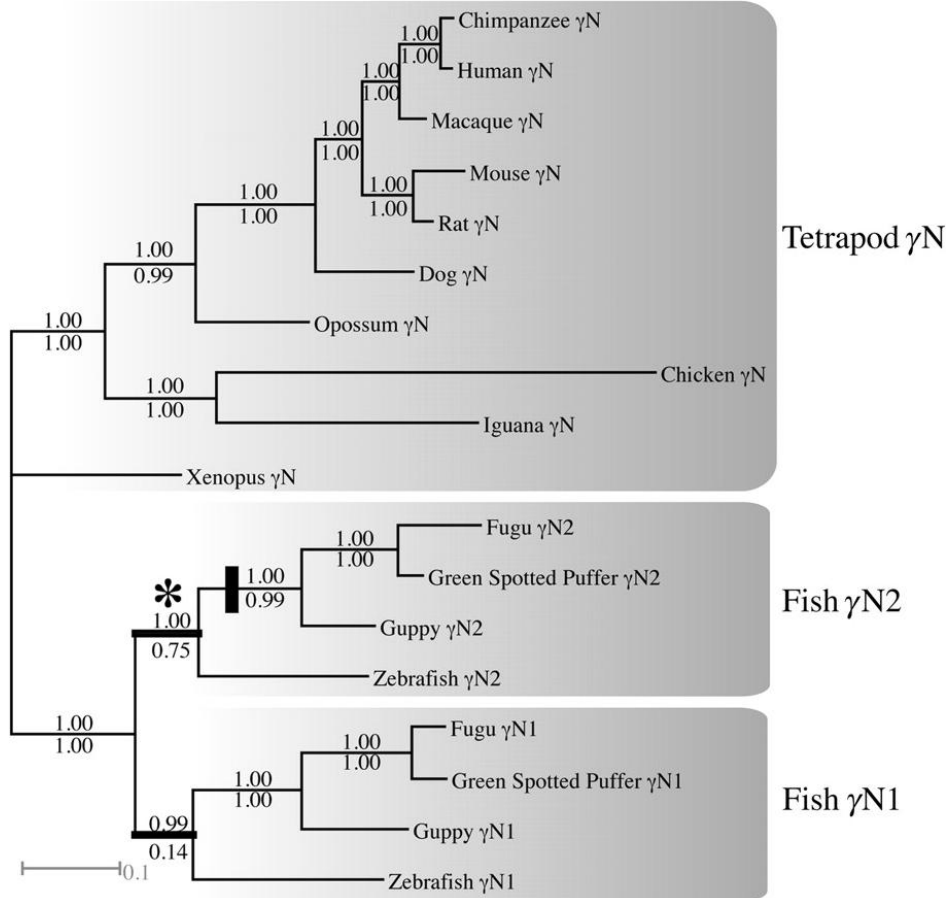


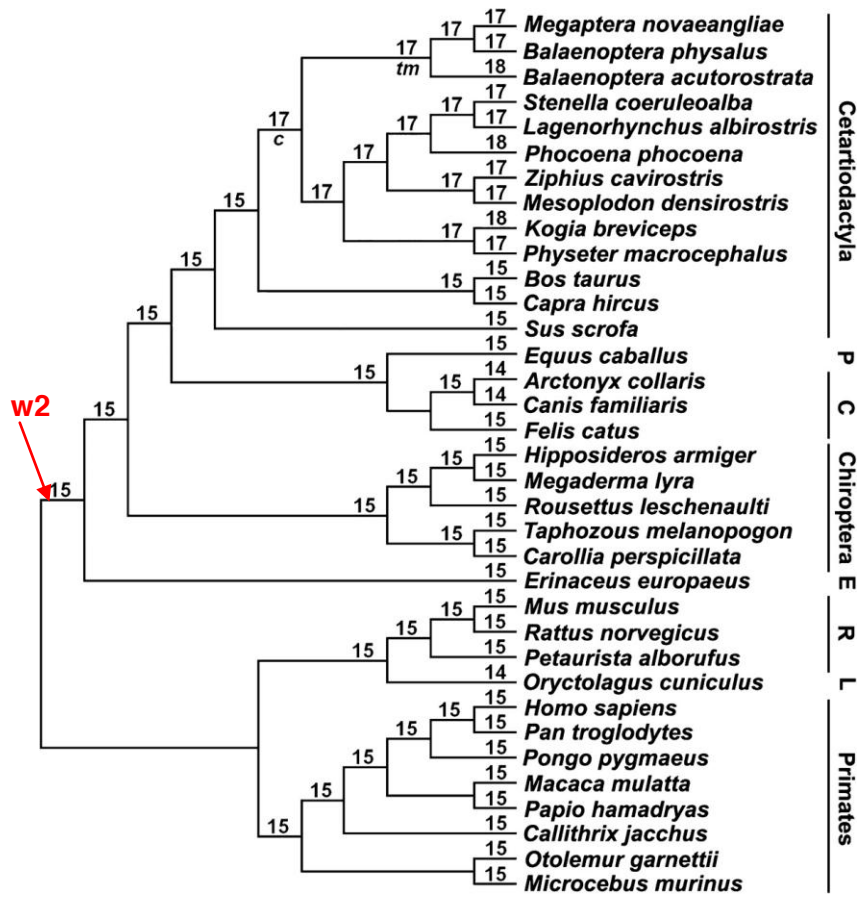


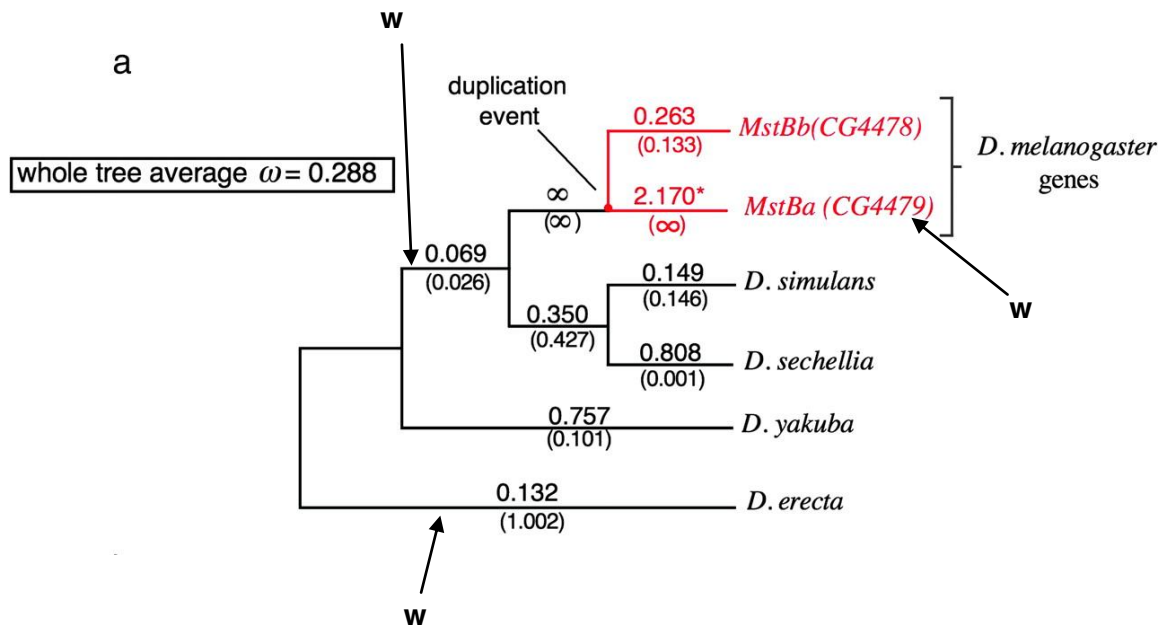
Case 20

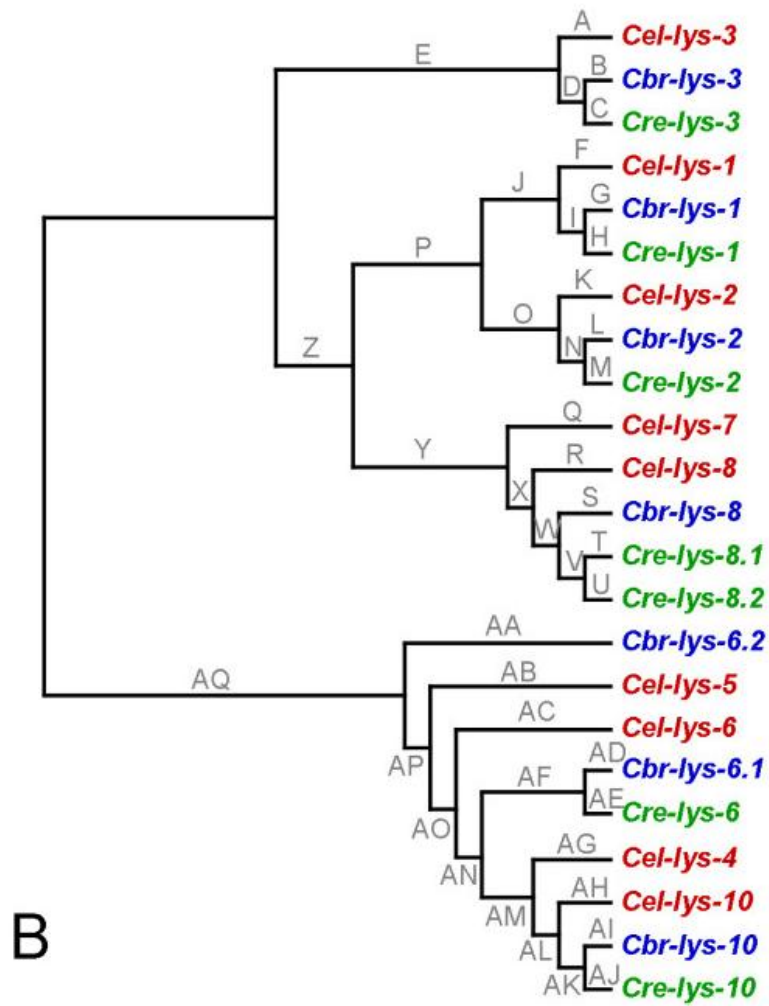


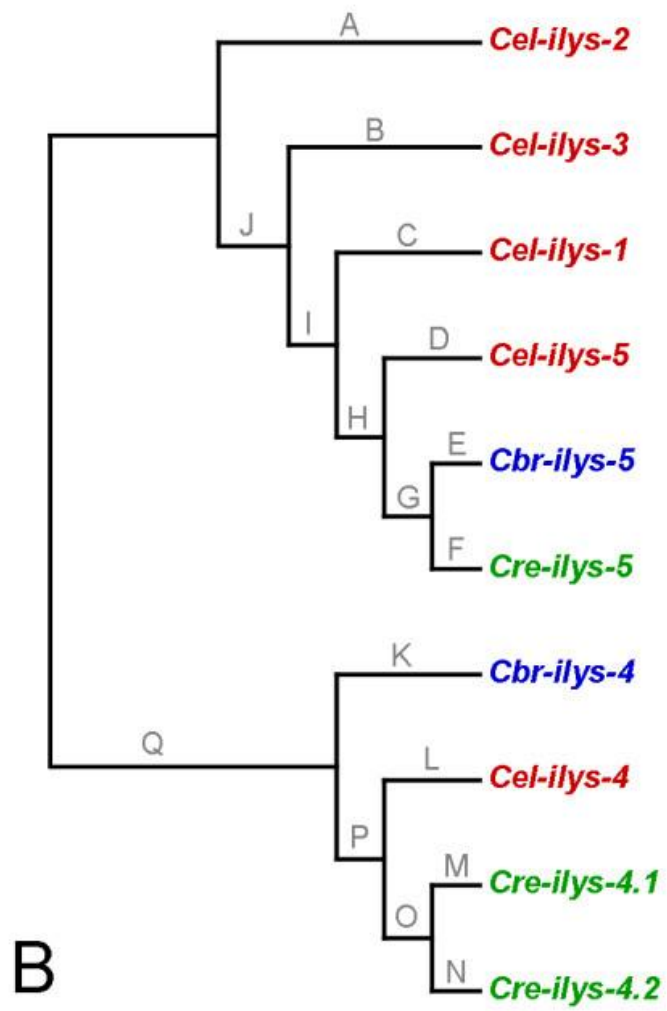
Case 21



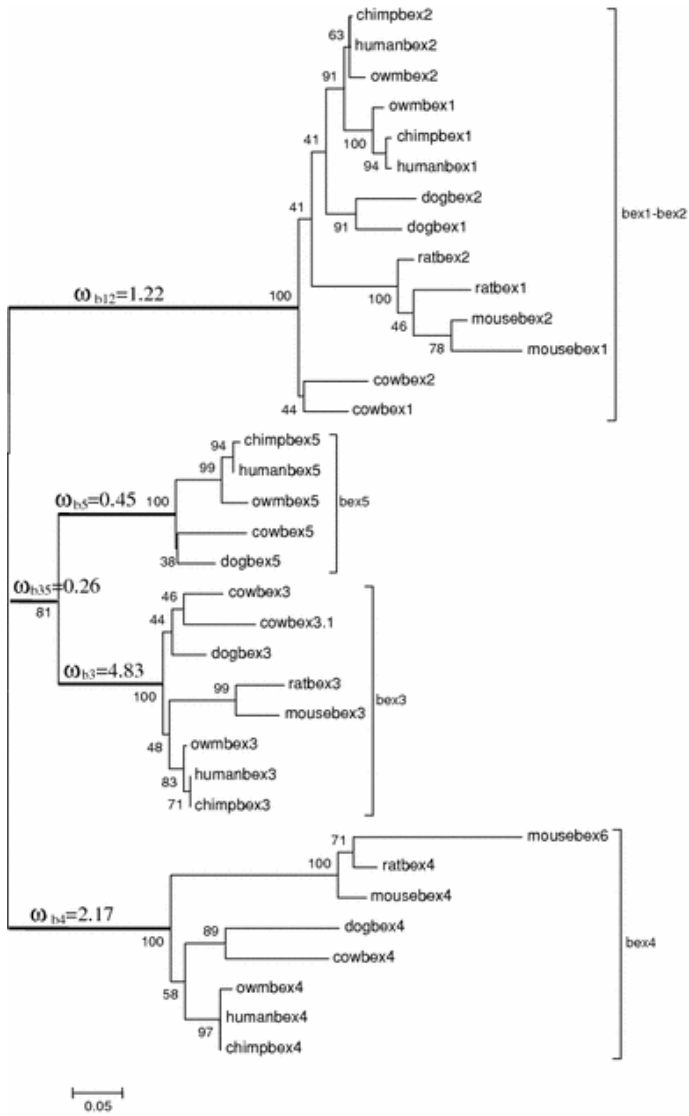




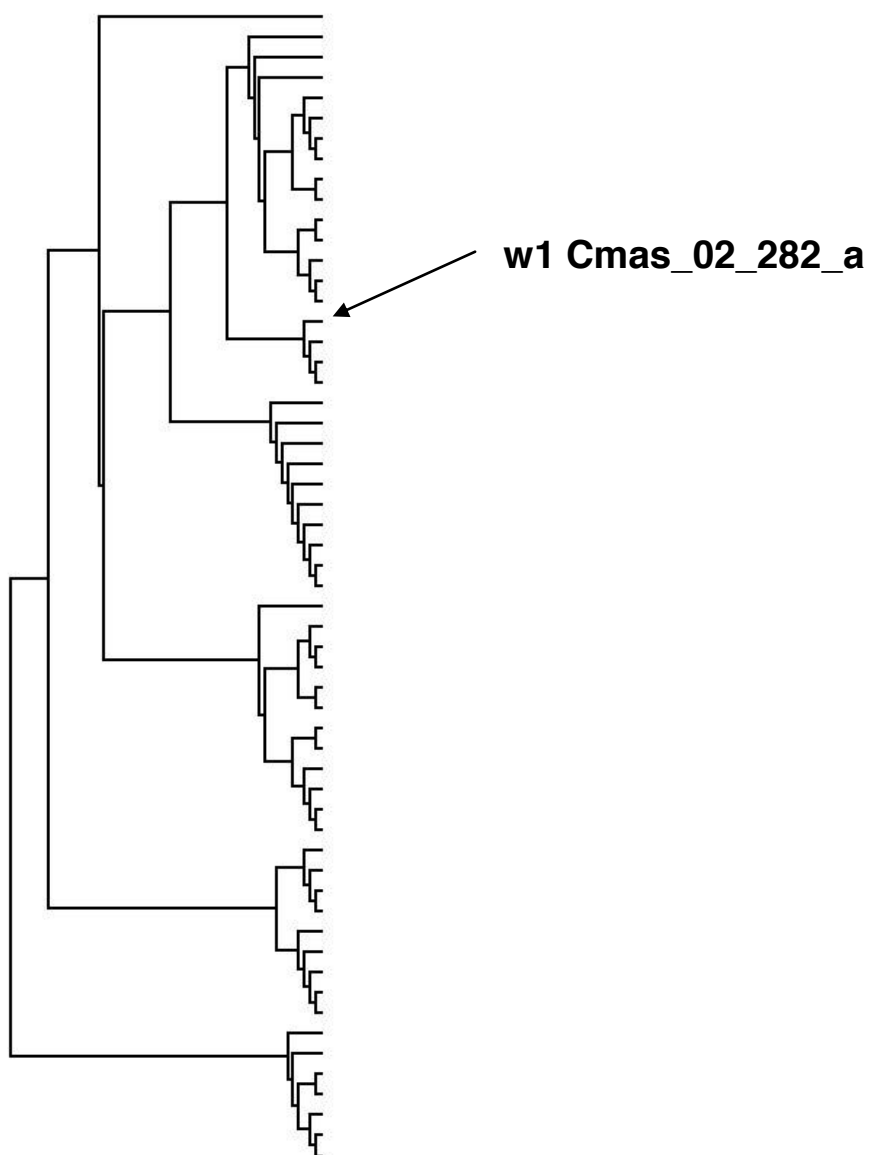




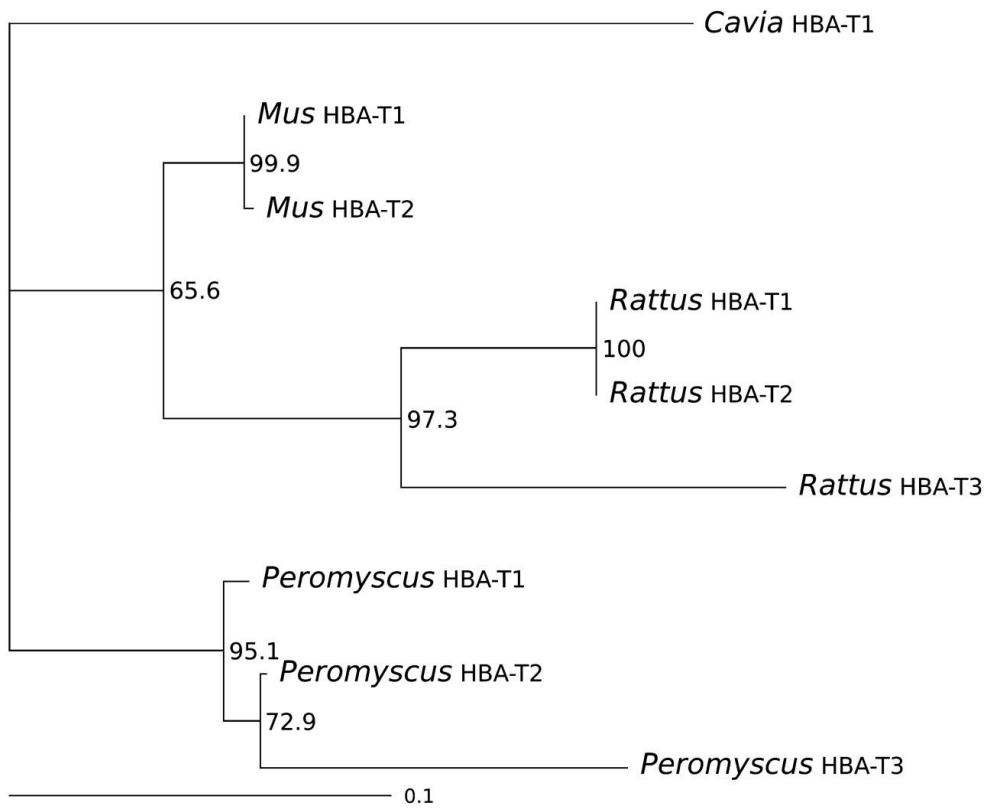
Case 26

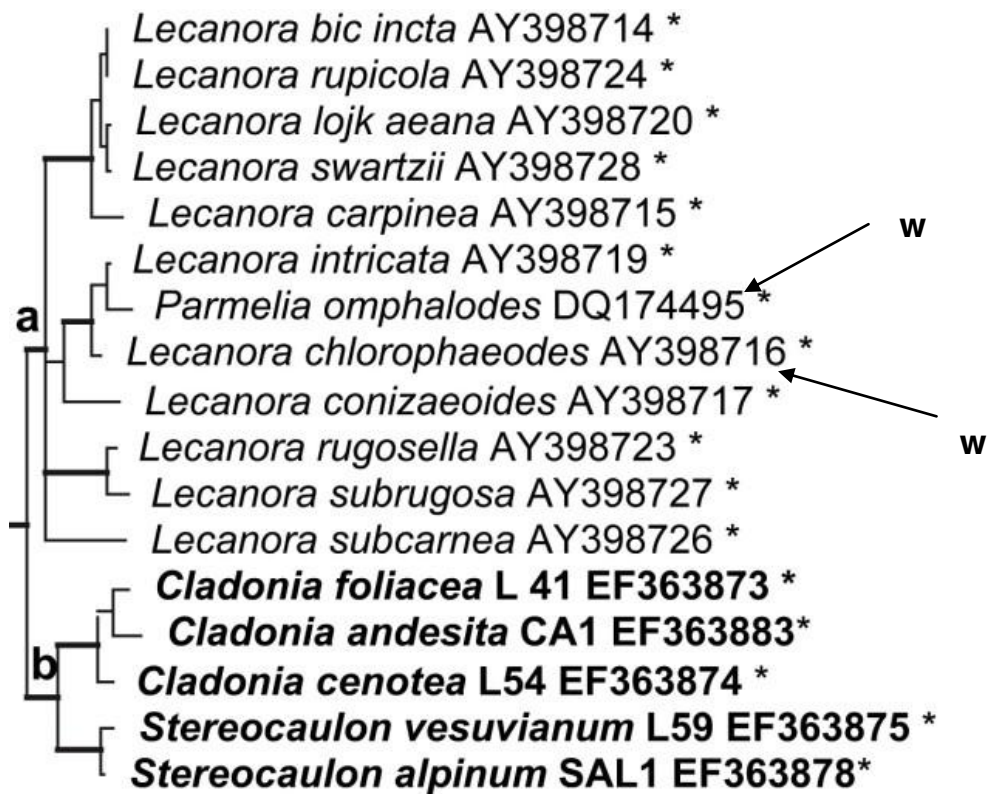


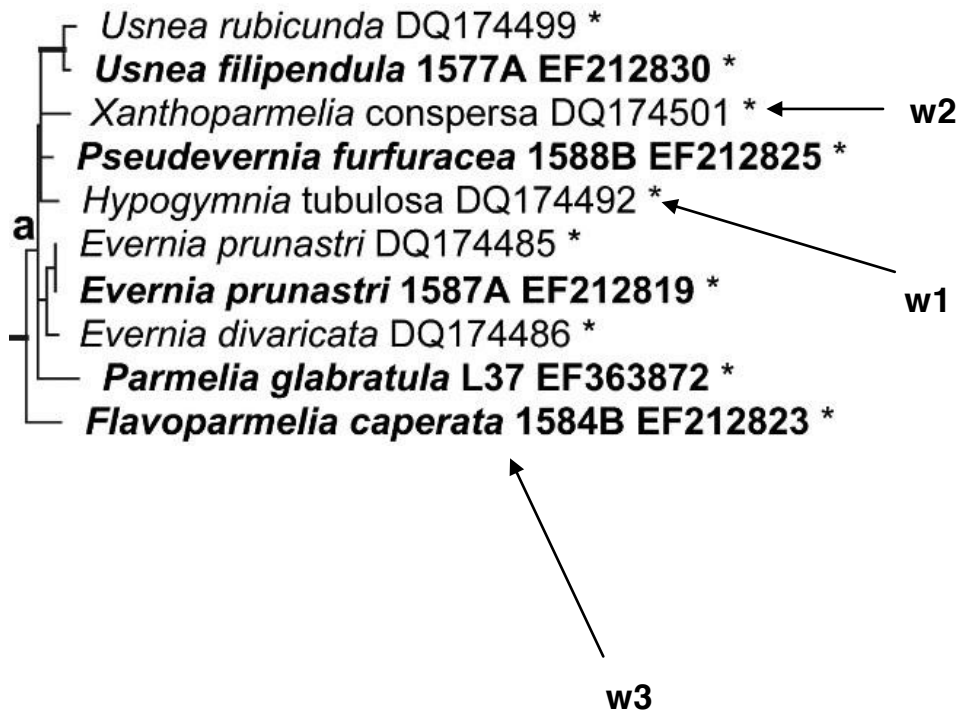
Case 27

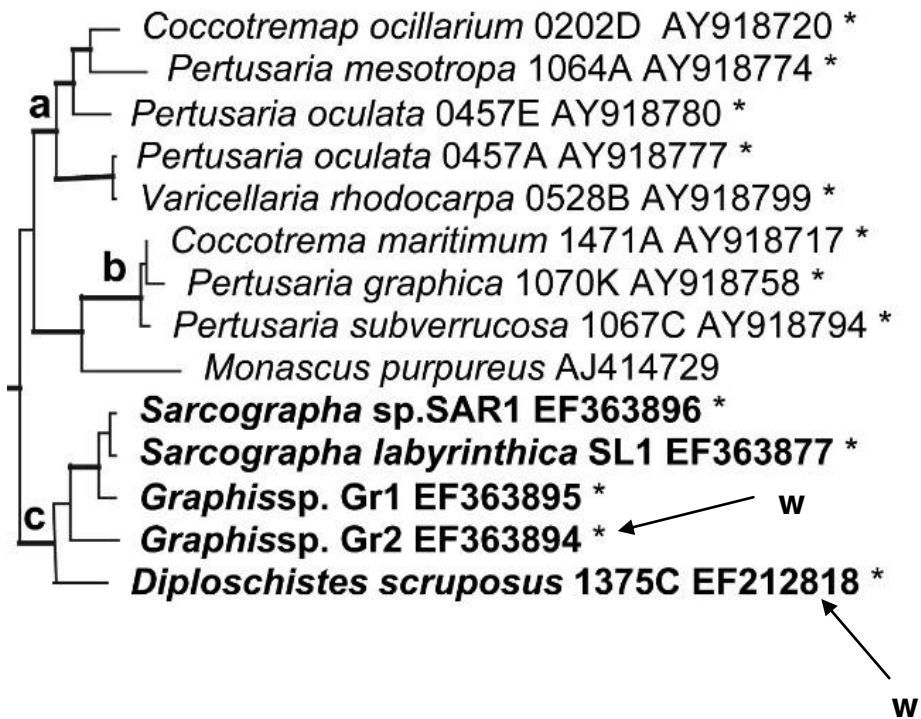


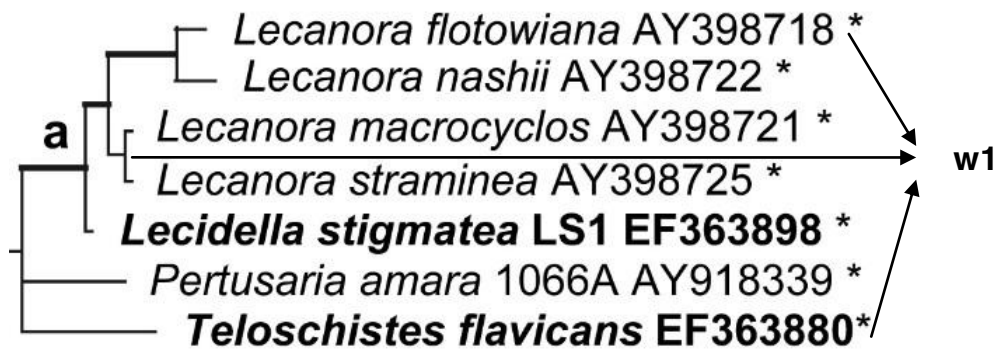
Case 28

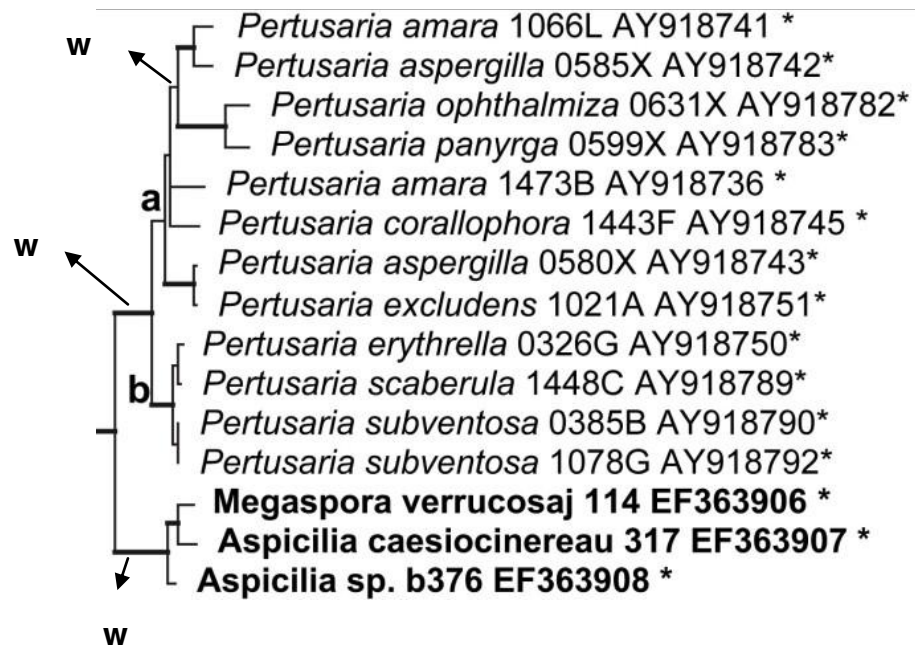


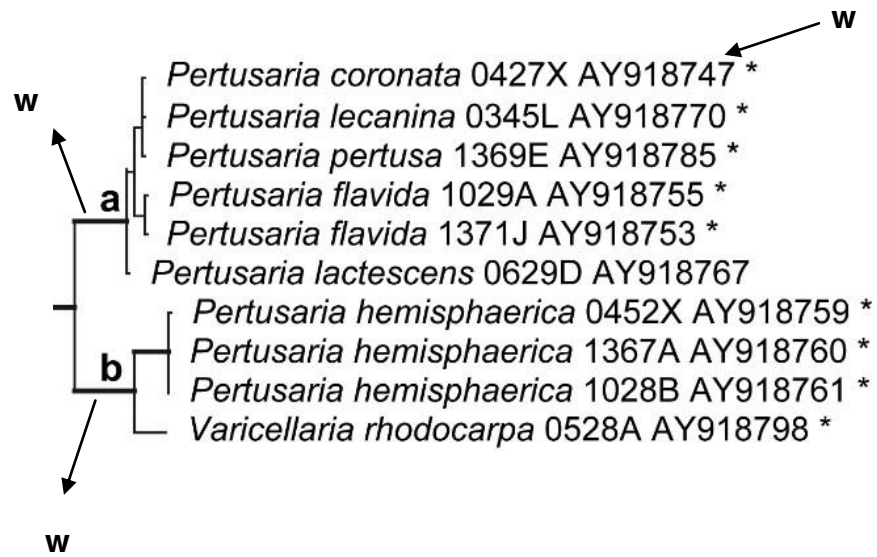


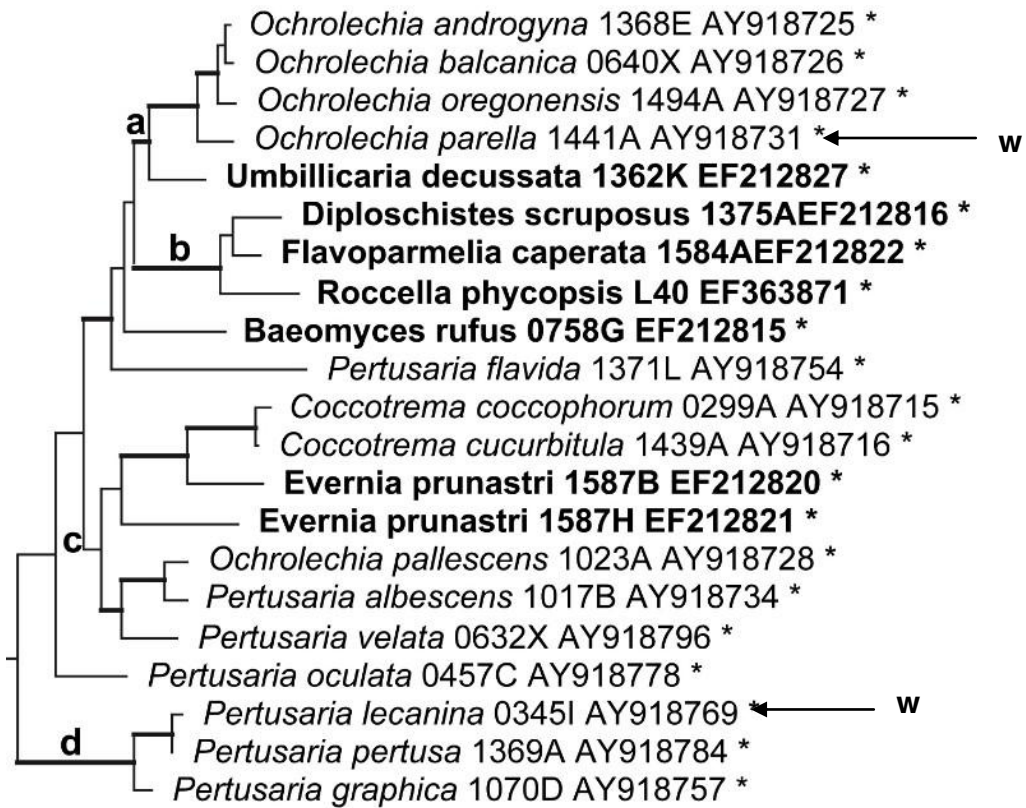


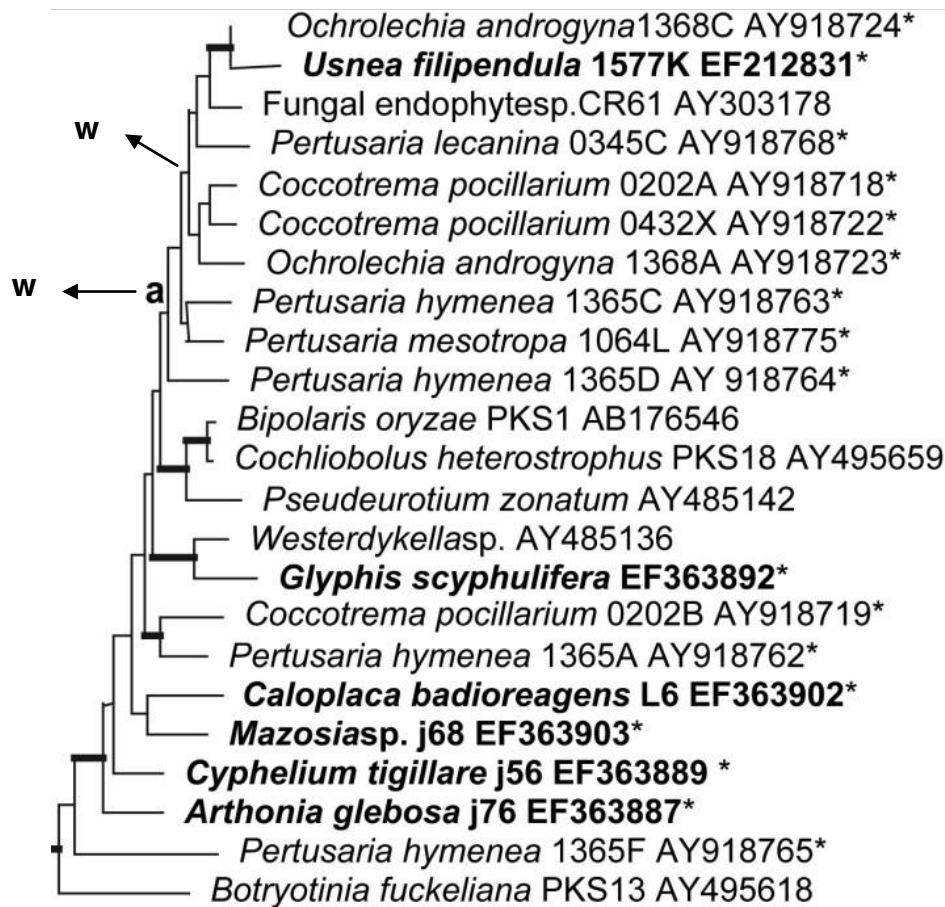


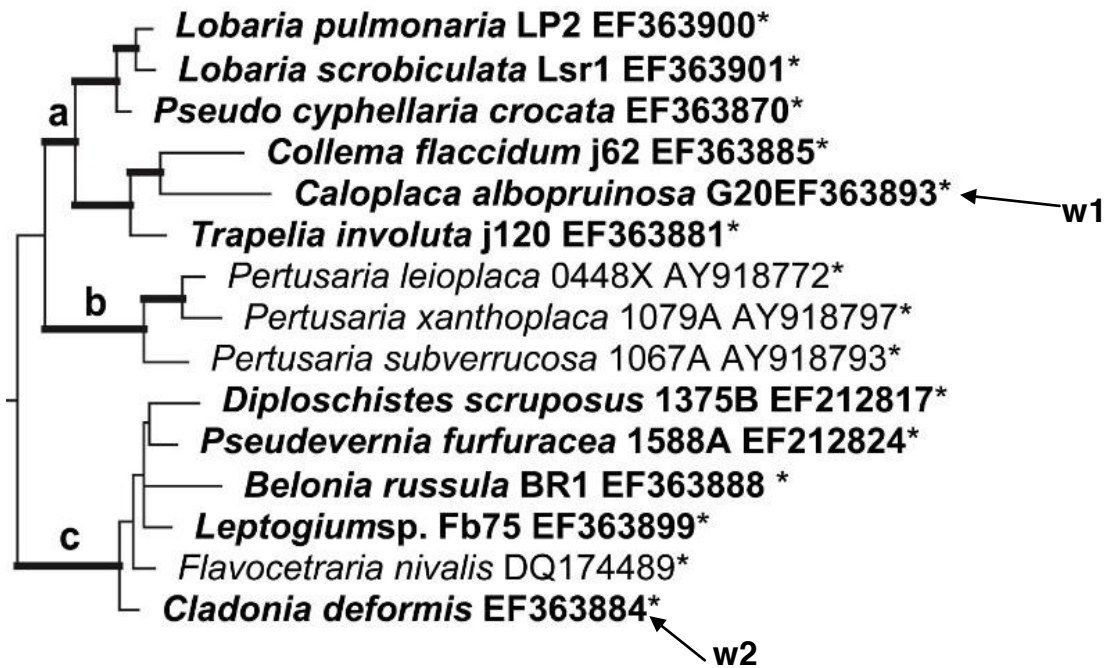


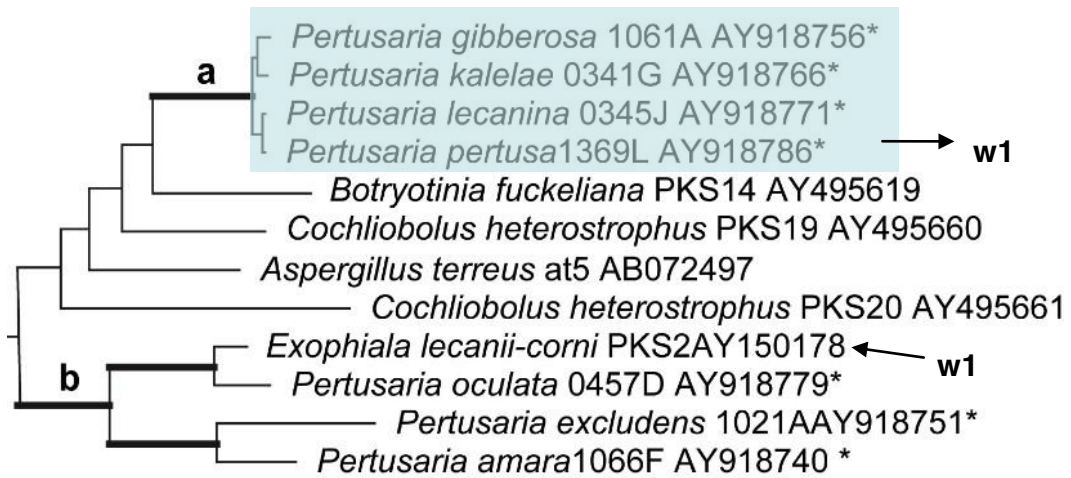


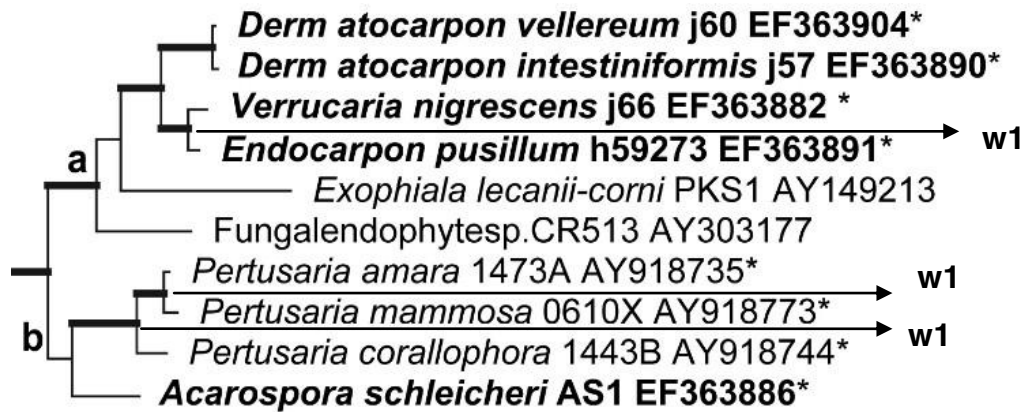


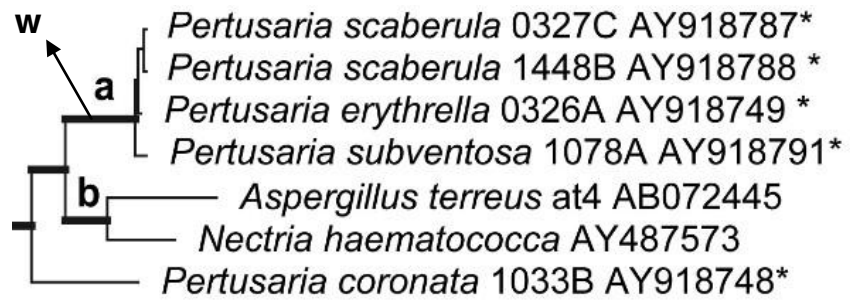




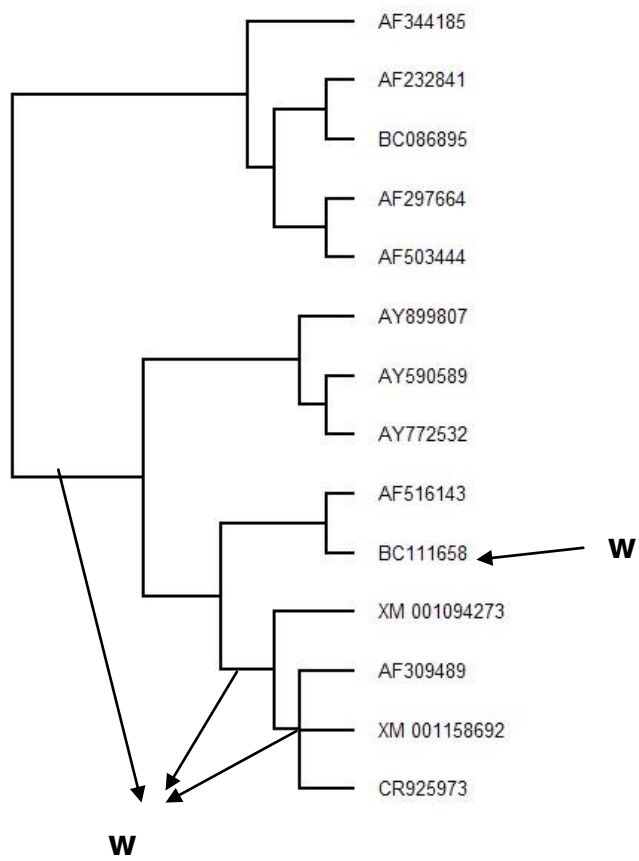




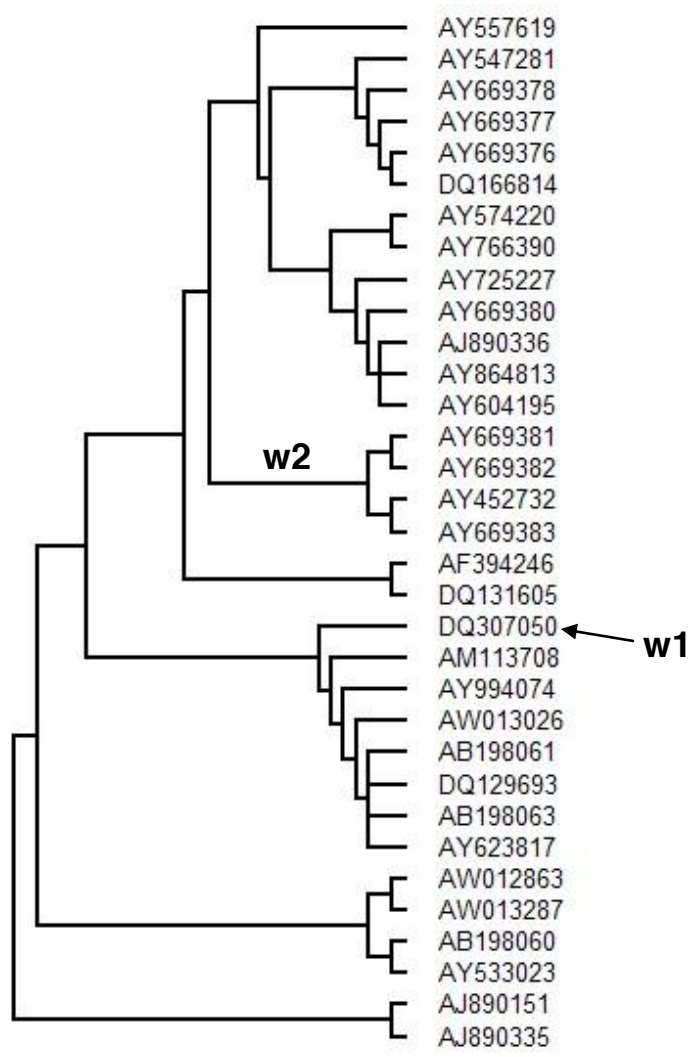




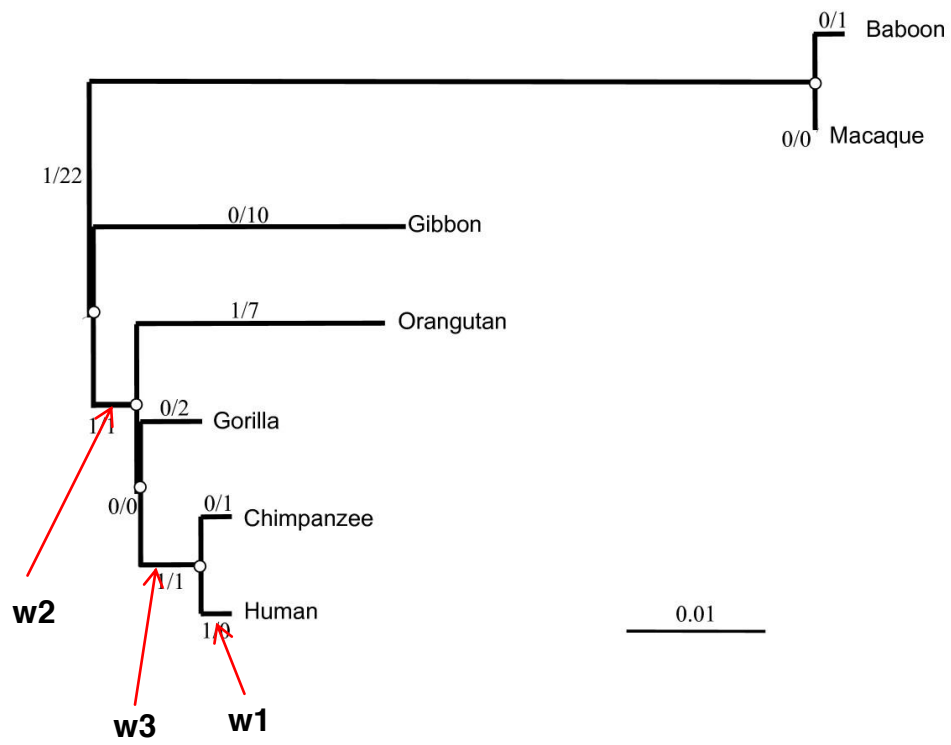
Case 41



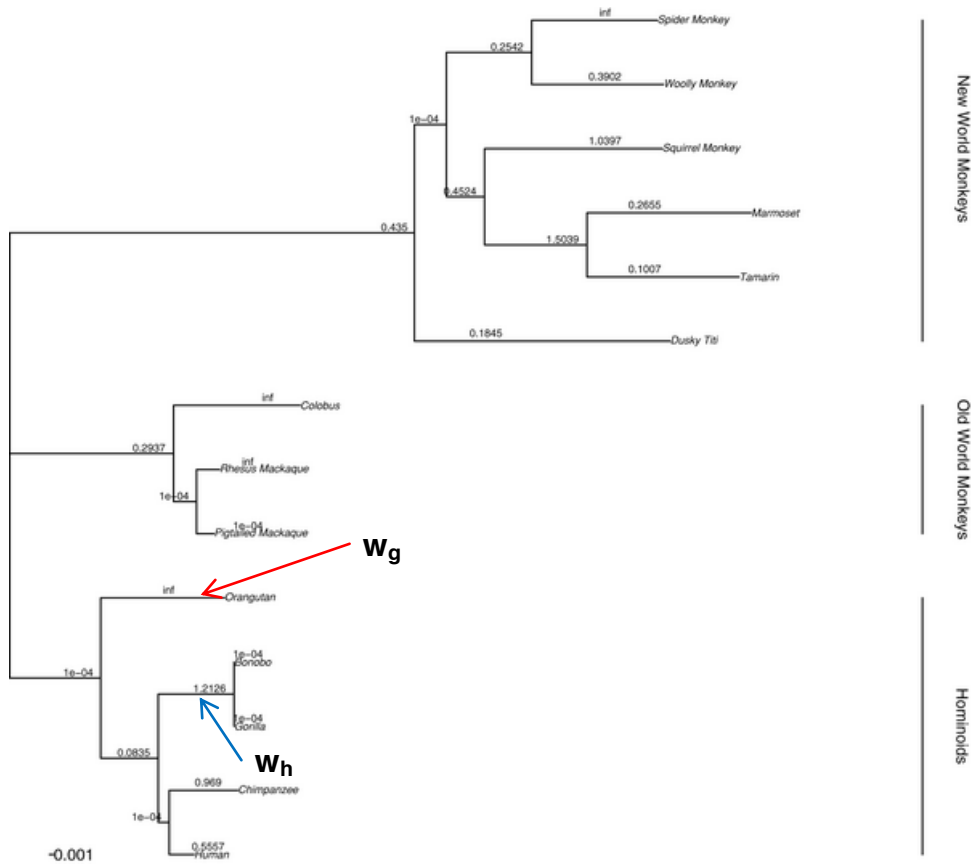
Case 42

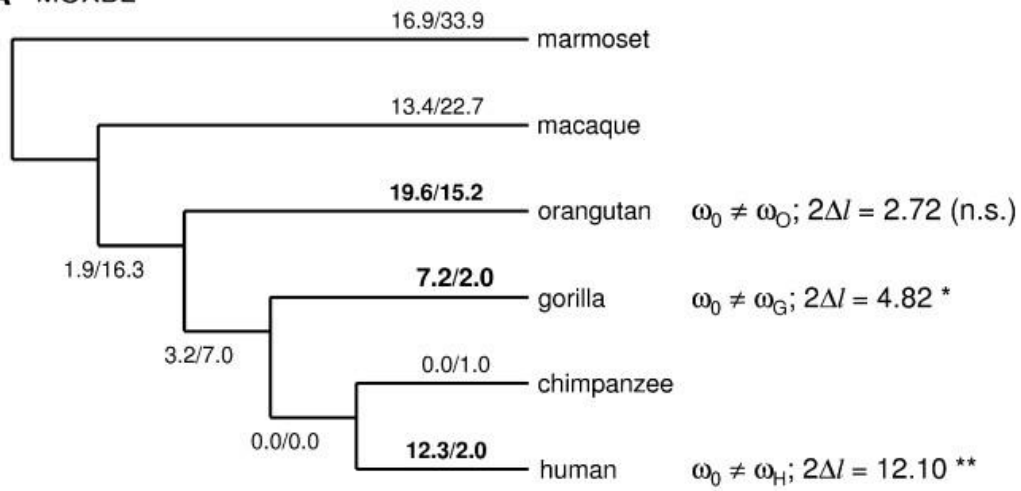


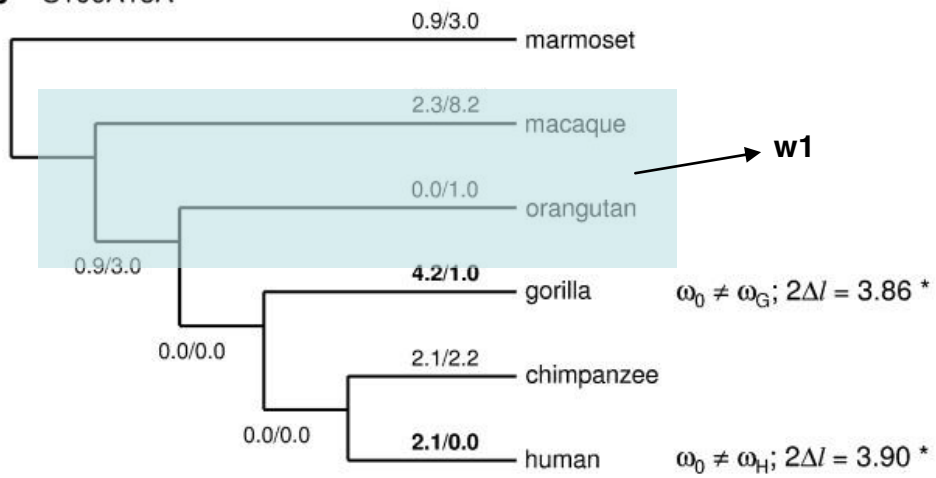
Case 43



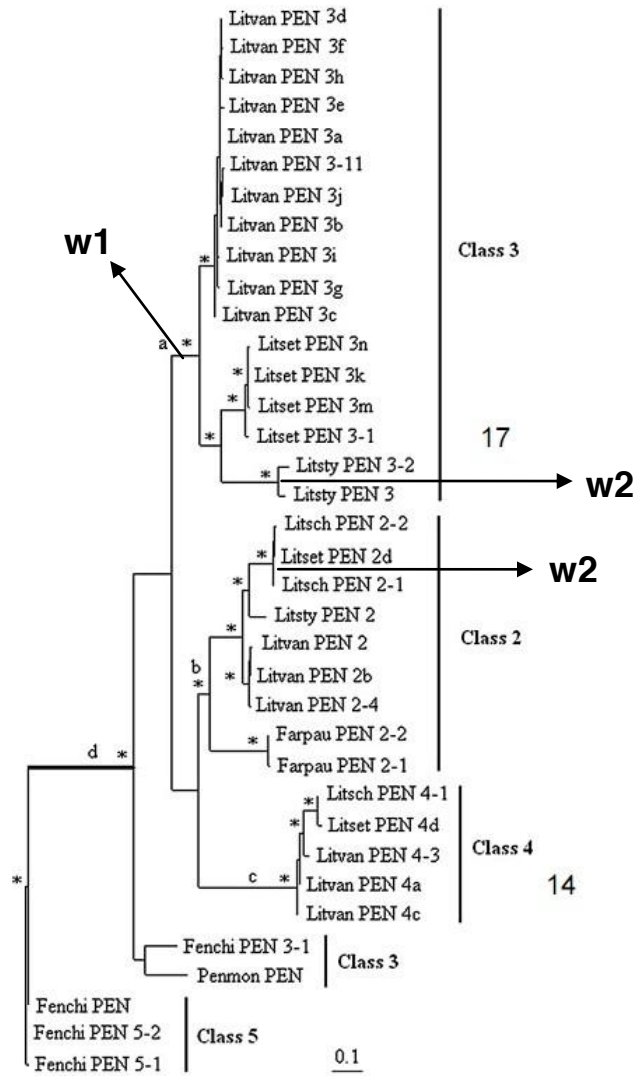
Case 44



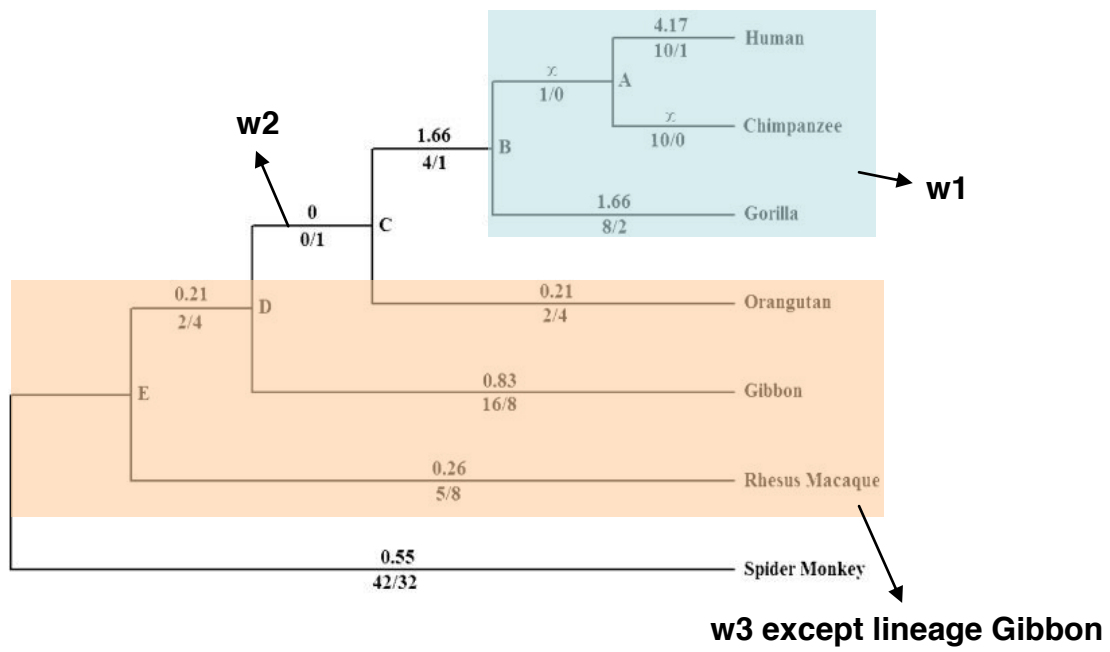
A MOXD2

B S100A15A

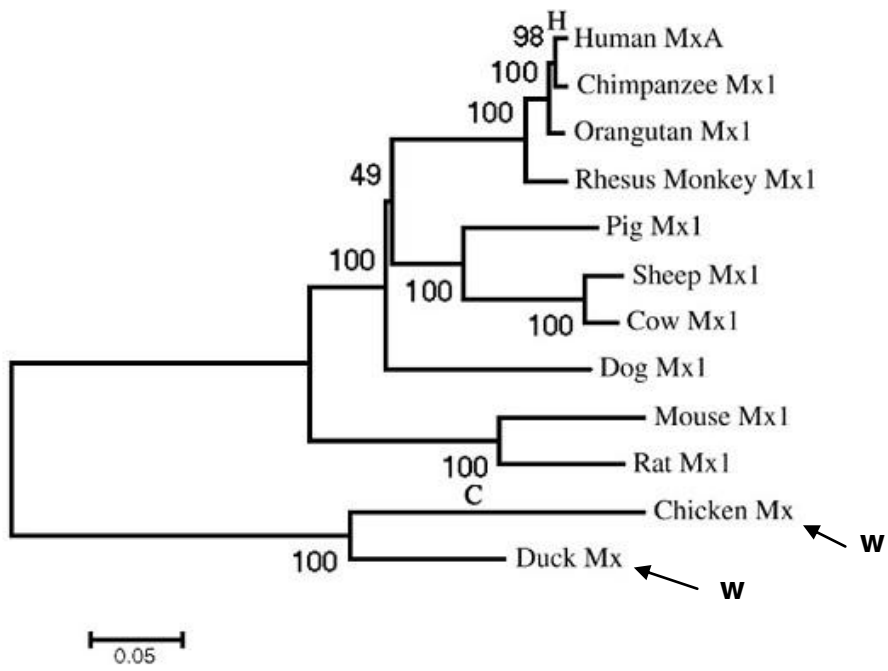
Case 47



Case 48



Case 49



Case 50

