

# Evolutionary origins of human apoptosis and genome stability gene networks

## Supporting Online Material

Mauro A. A. Castro<sup>1,3\*</sup>, Rodrigo J. S. Dalmolin<sup>1\*</sup>, José C. F. Moreira<sup>1</sup>, José C. M. Mombach<sup>4</sup> & Rita M. C. de Almeida<sup>2</sup>

<sup>1</sup>Departamento de Bioquímica, Universidade Federal do Rio Grande do Sul, Rua Ramiro Barcelos 2600-anexo, Porto Alegre 90035-003, Brazil. <sup>2</sup>Instituto de Física, Universidade Federal do Rio Grande do Sul, Avenida Bento Gonçalves 9500, Porto Alegre 91501-970, Caixa Postal 15051, Brazil. <sup>3</sup>Universidade Luterana do Brasil, Avenida Itacolomi 3600, Gravataí 94170-240, Brazil. <sup>4</sup>Universidade Federal de Santa Maria, Santa Maria 97105-900, Brazil.

\*These authors contributed equally to this work.

Correspondence to: Mauro A. A. Castro<sup>1</sup> (mauro@ufrgs.br) or Rita M. C. de Almeida<sup>2</sup> (rita@if.ufrgs.br).

# Contents

<b>1. SUPPLEMENTARY RESULTS AND DISCUSSION</b> .....	<b>3</b>
1.1. Exemplification of the evolutionary analysis.....	3
1.1.1. Parsimony analysis and evolutionary scenarios .....	3
1.1.2. From STNs to GNNs .....	4
1.1.3. Plasticity analysis .....	6
1.2. Consistency of evolutionary and functional data.....	6
1.2.1. Network statistics.....	6
1.2.2. Orthologous groups statistics .....	7
1.2.2.1 Comparing evolutionary scenarios: KOG x Inparanoid orthologous groups.....	7
1.2.3. Mouse and yeast statistics .....	9
1.2.4. Cancer statistics.....	10
1.3. The deep root of eukaryotes.....	11
1.4. The deep root of metazoans .....	12
<b>2. REFERENCES</b> .....	<b>105</b>

## LIST OF SUPPLEMENTARY FIGURES

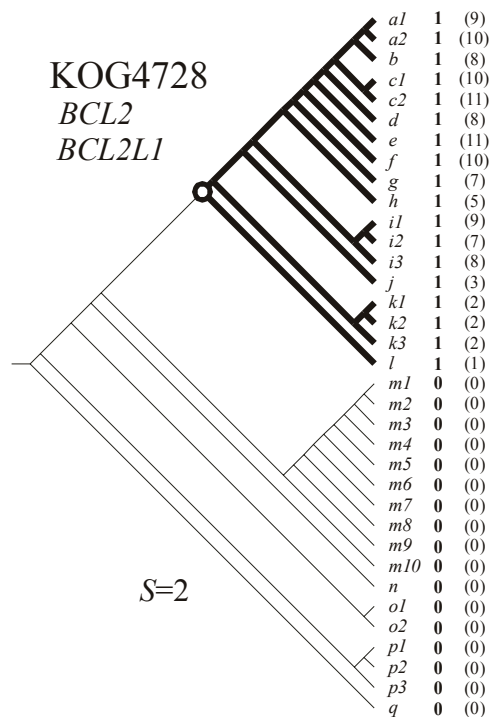
<b>Supplementary Figure S1.</b> Inferring ancestral state of human apoptosis and genome stability genes. ....	<b>3</b>
<b>Supplementary Figure S2.</b> Orthology information transferred for graphs.....	<b>5</b>
<b>Supplementary Figure S3.</b> Inferring ancestral state of human apoptosis and genome stability genes according to Inparanoid Database.....	<b>13</b>
<b>Supplementary Figure S4.</b> Inconsistency scores $S$ estimated from KOG and Inparanoid orthology analysis..	<b>14</b>
<b>Supplementary Figure S5.</b> Comparing the evolutionary roots inferred from KOG and Inparanoid orthology analysis.....	<b>15</b>
<b>Supplementary Figure S6.</b> Inparanoid evolutionary scenarios: from species-tree nodes (STNs) to gene-network nodes (GNNs). ....	<b>16</b>
<b>Supplementary Figure S7.</b> Summary of cancer statistics.....	<b>17</b>
<b>Supplementary Figure S8.</b> KOG-to-COG correspondence.....	<b>18</b>
<b>Supplementary Figure S9.</b> Subsample of the data with <i>Nematostella vectensis</i> .....	<b>19</b>
<b>Supplementary Figure S10.</b> Alignment of amino acid sequences of 40S ribosomal proteins.....	<b>20</b>
<b>Supplementary Figure S11.</b> Alignment of amino acid sequences of initiation factor 5A proteins.....	<b>21</b>
<b>Supplementary Figure S12.</b> Alignment of amino acid sequences of 5-3 exonucleases.....	<b>22</b>
<b>Supplementary Figure S13.</b> Divergence matrix according to amino acid alignments.....	<b>23</b>
<b>Supplementary Figures S14 to S49.</b> Parsimony analysis of KOG's orthologous groups.....	<b>24</b>
<b>Supplementary Figures S50 to S94.</b> Parsimony analysis of InParanoid's orthologous groups.....	<b>60</b>

# 1. Supplementary results and discussion

## 1.1. Exemplification of the evolutionary analysis

### 1.1.1. Parsimony analysis and evolutionary scenarios

To illustrate the parsimony analysis, consider the evolutionary scenario presented in **Supplementary Figure S1** for KOG4728. This KOG comprises 123 genes distributed among 18 species and the pattern of presence or absence of the species in the analyzed set of species is indicated by, respectively, zero or one. Observe that at least one ortholog is present in all extant metazoan species from *Homo sapiens* (*a1*) to *Caenorhabditis elegans* (*l*).



**Supplementary Figure S1.** Inferring ancestral state of human apoptosis and genome stability genes. Parsimony analysis of KOG 4728. The inconsistency value  $S$  is indicated. Branch codes: *a1* (*Homo sapiens*); *a2* (*Pan troglodytes*); *b* (*Macaca mulatta*); *c1* (*Rattus norvegicus*); *c2* (*Mus musculus*); *d* (*Canis familiaris*); *e1* (*Bos Taurus*); *f* (*Monodelphis domestica*); *g* (*Gallus gallus*); *h* (*Xenopus tropicalis*); *i1* (*Takifugu rubripes*); *i2* (*Tetraodon nigroviridis*); *i3* (*Danio rerio*); *j* (*Ciona intestinalis*); *k1* (*Drosophila melanogaster*); *k2* (*Anopheles gambiae*); *k3* (*Apis mellifera*); *l* (*Caenorhabditis elegans*); *m1* (*Kluyveromyces lactis*); *m2* (*Saccharomyces cerevisiae*); *m3* (*Candida glabrata*); *m4* (*Eremothecium gossypii*); *m5* (*Debaryomyces hansenii*); *m6* (*Yarrowia lipolytica*); *m7* (*Aspergillus fumigatus*); *m8* (*Schizosaccharomyces pombe*); *m9* (*Filobasidiella neoformans*); *m10* (*Encephalitozoon cuniculi*); *n* (*Dictyostelium discoideum*); *o1* (*Arabidopsis thaliana*); *o2* (*Cyanidioschyzon merolae*); *p1* (*Plasmodium falciparum*); *p2* (*Cryptosporidium hominis*); *p3* (*Thalassiosira pseudonana* CCMP1335); *q* (*Giardia lamblia* ATCC 50803).

To explain this common genetic trace, the parsimonious evolutionary scenario assumes that the ortholog was present in the last common ancestor (LCA) of metazoans, and was genetically transmitted to the descendants up to the species we know today. Given that two human apoptotic genes are listed in KOG4728 (*BCL2* and *BCL2L1*), the evolutionary root of these two human genes is inferred in the origin of metazoans as well. The evolutionary scenarios for the remaining 141 KOGs are provided in [Supplementary Figures S14-S49](#), and the parsimony analysis required to reconcile orthologous groups with the species tree topology is resolved according to the gain/penalty approach (Mirkin et al. 2003), where the most parsimonious scenario of presence/absence of all the genes at all ancestral nodes of the tree was predicted by

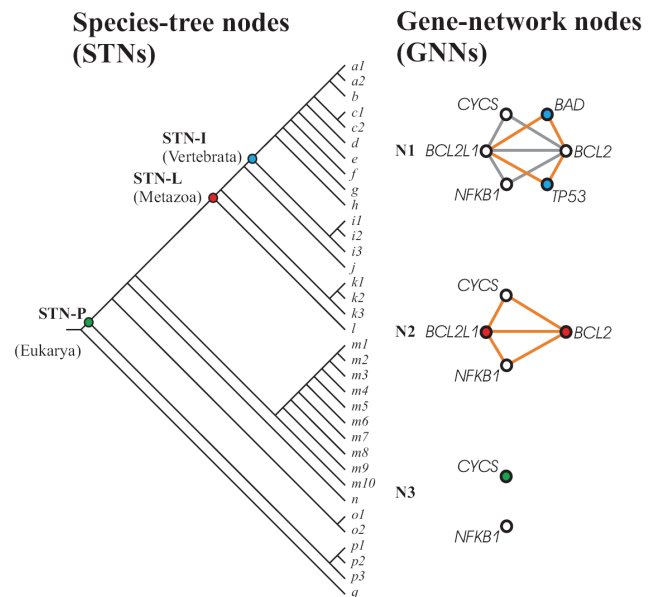
$$S = \lambda + g\gamma, \quad (1)$$

where  $\lambda$  is the number of gene losses,  $\gamma$  is the number of gene gains and  $g$  is the gain penalty. Then, a parsimonious scenario must minimize the total score according to the lowest evolutionary cost. The minimal score is referred to as the inconsistency value  $S$  for the given orthologous group and is associated with the number/types of events required to reconcile the evolutionary scenario for the given orthologous group with the species tree. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (*i.e.*  $g=2$ ), and one cost unit for gene loss (Snel et al. 2002; Kunin and Ouzounis 2003). Therefore, in the case of KOG4728 illustrated in [Supplementary Figure S1](#),  $S=2$ .

### 1.1.2. From STNs to GNNs

In order to exemplify the analysis presented in Figure 3 in the paper, here we applied the same approach but in a small set of six apoptotic genes ([Supplementary Figure S2](#)). Note that, as defined in the paper, any *species-tree node* (STN) represents a LCA, while any *gene-network node* (GNN) represents an ortholog in the human apoptosis and genome stability gene network (*i.e.* one STN can comprise more than one GNN roots). The orthology information from STNs is then overlaid on the network graphs, which are arranged in [Supplementary Figure S2](#) to optimally display the intersection within the species tree.

**Supplementary Figure S2.** Orthology information transferred for graphs. The orthology information from STNs is overlaid on GNNs, which are arranged for best observing the intersection within the species tree. It illustrates the projection of orthology information inferred for six orthologs. Roots of an ortholog: color nodes; presence of an ortholog: white nodes. Branch codes as in **Supplementary Figure S1**.



Therefore, the orthology information from STN-I is transferred for graph N1 and shows that *BAD* and *TP53* genes are rooted in the LCA of vertebrates (blue GNNs), while the other four genes are rooted prior to this species tree level (white GNNs). Also, in graphs N2 and N3 we show the orthology information inferred in STN-L and -P, respectively. Observe that *BAD* and *TP53* genes were removed from graph N2 because they are not placed in the LCA of metazoans. By the same reason, only *NFKB1* and *CYCS* genes are placed in graph N3, which have no links because these two ancient components are not functionally associated in the network-based model of human apoptosis and genome stability genes (Castro et al. 2007). Therefore, every STN intersection can produce an orthology projection onto the network graph. We remark that all orthology analysis should be restricted to the network components, given that the links represent association among human proteins. The strategy is then to follow the evolutionary steps that lead to the human gene/protein association network, *i.e.*, the roots of the human genes. In the paper we present the results for the set of 180 human genes.

### 1.1.3. Plasticity analysis

To exemplify the plasticity analysis, consider the KOG4728 illustrated in [Supplementary Figure S1](#). As stated before, this KOG comprises 123 genes distributed in 18 species. However, the orthologous genes are not equally distributed, as long as an orthologous group may have more genes associated with some species than others. In theory, if a given orthologous group has only one ortholog in every organism of the species tree, then its diversity will be maximum ( $H_\alpha=1.0$ ). In contrast, if all genes are present in only one organism, then its diversity will be minimum ( $H_\alpha=0.0$ ); note that  $H_\alpha$  zero will never be achieved, as long as at least three species are required to build an orthologous group (Tatusov et al. 2003). Thus, the former indicates a broad distribution in the species tree comparing to the latter, a possible measure of the evolutionary conservation. As a complementary quantity, we also estimate the abundance  $D_\alpha$ , which is simply obtained by the ratio between the number of orthologous genes and the number of organisms. If an orthologous group has only one ortholog per organism, *i.e.*,  $D_\alpha=1.0$ , then it indicates poor representation in any particular species. In contrast, KOGs with  $D_\alpha \gg 1.0$  indicates many variants per organisms, a possible measure of the evolutionary plasticity. These statistics are also extended to Inparanoid database in order to assess a different data source.

## 1.2. Consistency of evolutionary and functional data

### 1.2.1. Network statistics

It is important to remark that some apoptotic genes are missing from KEGG map, which is used to construct the apoptosis gene network. However, given our network approach we can not describe genes individually. To consider other genes we must consider firstly the protein interaction network. This problem is critical for reliable protein network reconstruction and KEGG represent curated reconstructions of protein-protein interactions. This feature made KEGG a reference source that benchmarks numerous system biology databases (*e.g.* Cancer Genome Projected at NIH/CGAP and STRING database). Also, such option for KEGG

potentially enhances the primary applicability of our orthology map, that is, the transferability of functional information from several organisms to human. A complete description of the construction of the genome maintenance gene network is described in our previous work (Castro et al. 2007).

### *1.2.2. Orthologous groups statistics*

Many of the current ortholog databases that uses genome-wide analysis will likely contain false-positives due to the limitations of the reciprocal-best-hits (RBH) approach (*e.g.* predict paralog as ortholog). However, in absence of golden standards, actually it is not possible to discriminate the best way of deriving orthologous groups, and whether there is in fact one way that works best for all applications. Although KOGs provide a sensible approach that does not rely on arbitrary score cut-offs, due to potential bias we extensively confronted our evolutionary scenarios obtained from KOGs with other independent sources (*i.e.* SGD, MGD, CGP, XP databases, as discusses in the paper, and Inparanoid database). Next, we present the evolutionary scenario for apoptosis and genome stability orthologs obtained from Inparanoid database, which provides a fully automatic method for finding orthologs (Remm et al. 2001).

#### *1.2.2.1 Comparing evolutionary scenarios: KOG x Inparanoid orthologous groups*

To test the robustness of the evolutionary scenarios predicted by our analysis we investigated the evolutionary roots of the same set of genes but using a different orthology detection approach. Here we consider the Inparanoid database. In contrast to KOG algorithm, Inparanoid is designed to find orthologs and in-paralogs between two species and to separate in-paralogs from out-paralogs (Remm et al. 2001).

In [Supplementary Figure S3A](#) we present the topology of the species tree for the available organisms at Inparanoid database. The evolutionary roots inferred for human apoptosis and genome stability genes are indicated in [Supplementary Figure S3B,C](#). Comparing to the species tree derived from KOGs, Inparanoid derived species tree shows the same human ancestral nodes, except for the basal position. Also, the overall results of the pooled orthologs are very

similar comparing to the results presented in the paper (see Figure 2). To assess quantitatively the contrasts between KOG and Inparanoid evolutionary scenarios we also compared the inconsistency scores  $S$  (**Supplementary Figure S4**). For the entire gene set,  $S=4.39(\pm 2.07$ ; KOG analysis), and  $S=5.27(\pm 2.48$ ; Inparanoid analysis). However, this statistics estimates the inconsistency of the evolutionary scenarios considering the entire species tree. As long as the resulting species trees are slightly different between databases, then we should consider a complementary quantity in order to estimate the inconsistency of the roots inferred for each gene. The inconsistency of each gene root can be defined as

$$\Delta STN_i = | KOG_{STN} - Inparanoid_{STN} | \quad (2)$$

where  $KOG_{STN}$  represents the species-tree node where is rooted a given gene  $i$ , according to KOG database, and  $Inparanoid_{STN}$  represents the root of the same gene  $i$  according to Inparanoid database. The  $\Delta STN$  average value for the entire gene set  $n$  gives the evolutionary inconsistency score  $R$  of the evolutionary scenarios.  $R$  is given by

$$R = \frac{\sum_i^n \Delta STN_i}{n} \quad (3)$$

In **Supplementary Figure S5** we present  $R$  for apoptosis and genome stability genes. This figure estimates the divergence between KOG and Inparanoid derived scenarios, that is,  $R=1.709$  for apoptosis (approximately two STNs up- and down-ward to the rooting point in the species tree) and  $R= 0.807$  for genome stability (approximately one STNs up- and down-ward).

To address qualitatively the differences among the databases, we reconstructed the network graphs of the evolutionary scenarios using the evolutionary roots inferred from Inparanoid database (**Supplementary Figure S6**). Accordingly, the two major increments in apoptosis network are aligned between KOGs and Inparanoid databases: the emergence of the apoptosis intrinsic pathway at the metazoan origin and the emergence of the extrinsic pathway at the vertebrata origin. Although not the same components are present in both scenarios, the overall results in the network topology are equivalent. As addressed in the paper, here the evolutionary scenario shows that the genome maintenance mechanisms could be divided into two major segments: i) the evolution of genome stability gene network placed in the basal position of the species tree and ii) the evolution of apoptosis gene network with components rooted throughout eukaryotic evolution. Likewise, the network core of both systems is rooted before the divergence



of metazoans, while GNNs placed in the periphery of the networks represent more recent evolutionary innovations. Further details for all Inparanoid orthologous groups, the evolutionary scenarios and the corresponding  $S$  value are presented in [Supplementary Figures S50 to S94](#) and also provided in spreadsheet format ([Supplementary Table S4](#)).

### 1.2.3. Mouse and yeast statistics

Any evolutionary scenario providing the putative history that has given rise to the species placed in the phylogenetic tree must reflect, to some extent, the known functional differences and similarities observed in the living organism. Incongruence between evolutionary and functional data indicates a biased construction and should be revised. In the scenario described in the paper, the network regions of high and low evolutionary plasticity indicate that the network was not equally tolerant to genetic changes (*e.g.* mutations), as long as apoptosis have produced more variants of its components than genome stability. One possible explanation for the difference in plasticity between apoptosis and genome stability genes is to assume that the evolutionary plasticity of the network is proportional to the genetic changes that a living organism may support without disrupting its viability. Thus, to investigate this assumption, we assessed the lethality data of the unicellular eukaryotic model *Saccharomyces cerevisiae* available in the *Saccharomyces* Genome Database (SGD) (Hirschman et al. 2006), as well as the lethality data of *Mus musculus* orthologs in Mouse Genome Database (MGD) (Eppig et al. 2007).

Both databases provide gene knock-out data used to map genetic essentiality in the genomes. However, care should be taken when considering SGD together with MGD database. The available lethality statistics for this unicellular eukaryote comes from systematic experiments that cover almost all ORFs in the yeast genome, representing a comprehensive cellular lethality map. In contrast to yeast, mouse genes are not equally studied as long as *Mus musculus* data did not arise from systematic experiments. This explain why the network area that concentrates several yeast essential orthologs corresponds mainly to those in mouse that lacks knock-out data (see, in the paper, **Figure 5**: white GNNs). Also, the phenotypic statistics in MGD database consider lethal any allele that causes death anytime after fertilization and before the postnatal day 2; thus, knock-out alleles may indicate “developmental lethality” or “essentiality” to embryonic stem cells. Furthermore, evidences for mouse lethality data are obtained according to the

frequency expected by Mendelian genetics (*i.e.* zygosity and allelic distribution observed in the offspring): any significant deviation from the expected frequency for the knock-out allele indicates lethality, but it does not mean that the lethal genotype does not arise in the offspring.

#### 1.2.4. Cancer statistics

In the same line discussed above, care should be taken in order to consider cancer statistics together with yeast and mouse due to differences among data sources. For instance, *CAN* gene statistics comes mainly from epidemiological data and shows exclusively genes in which mutations that are causally implicated in oncogenesis have been described at least in two independent reports, showing mutations in primary patient material (Futreal et al. 2004). According to CGP census, the underlying rationale for interpreting a mutated gene as causal in cancer development is that the number and pattern of mutations in the gene are likely to have been selected because they confer a growth advantage on the cell population from which the cancer has developed (Futreal et al. 2004). Also, in contrast to mouse and yeast knock-out alleles, *CAN* gene may have a range of mutations, from a single nucleotide substitution to a complete transcript disruption (*i.e.* null alleles).

In order to circumvent such data limitations and improve the analysis we further investigated the human statistics assessing the genotypic profile of several *CAN* gene loci. We attempt to obtain the proportion of null and non-null alleles in human following the strategy used in mouse to infer lethality according to the expected frequency in a Mendelian distribution. We focus the analysis in the set of *CAN* genes placed in  $\rho$  module, given that they are collectively represented in the same locus-specific mutation database – XP mutation database (<http://www.xpmutations.org>). These *CAN* genes are also associated with the same DNA repair function (nucleotide-excision repair) and are related to three rare autosomal recessive human clinical disorders (Xeroderma pigmentosum, Cockayne Syndrome and Trichothiodystrophy), which may turn reliable the obtaining of a representative human sample. Due to obvious reason, human data do not come from uniform samples but represent a historical collection (*e.g.* case reports documented all around the world). We retrieved 182 mutated genotypes available in that database, which is then pooled according to the zygosity and the presence of null and non-null alleles (see Table 1A in the paper). Sample number is also compared to a second database in

order to attest the representativeness of the database (see **Table 1B** in the paper). In **Supplementary Figure S7** we present the cancer statistics according to the presence of somatic and/or germline *CAN* genes in the network (*i.e.* apoptosis and genome stability). In order to inspect these results together with the main results present in the paper, this figure also shows the statistical contrasts for plasticity data in relation to abundance  $D\alpha$ , diversity  $H\alpha$ , gene function, and gene essentiality.

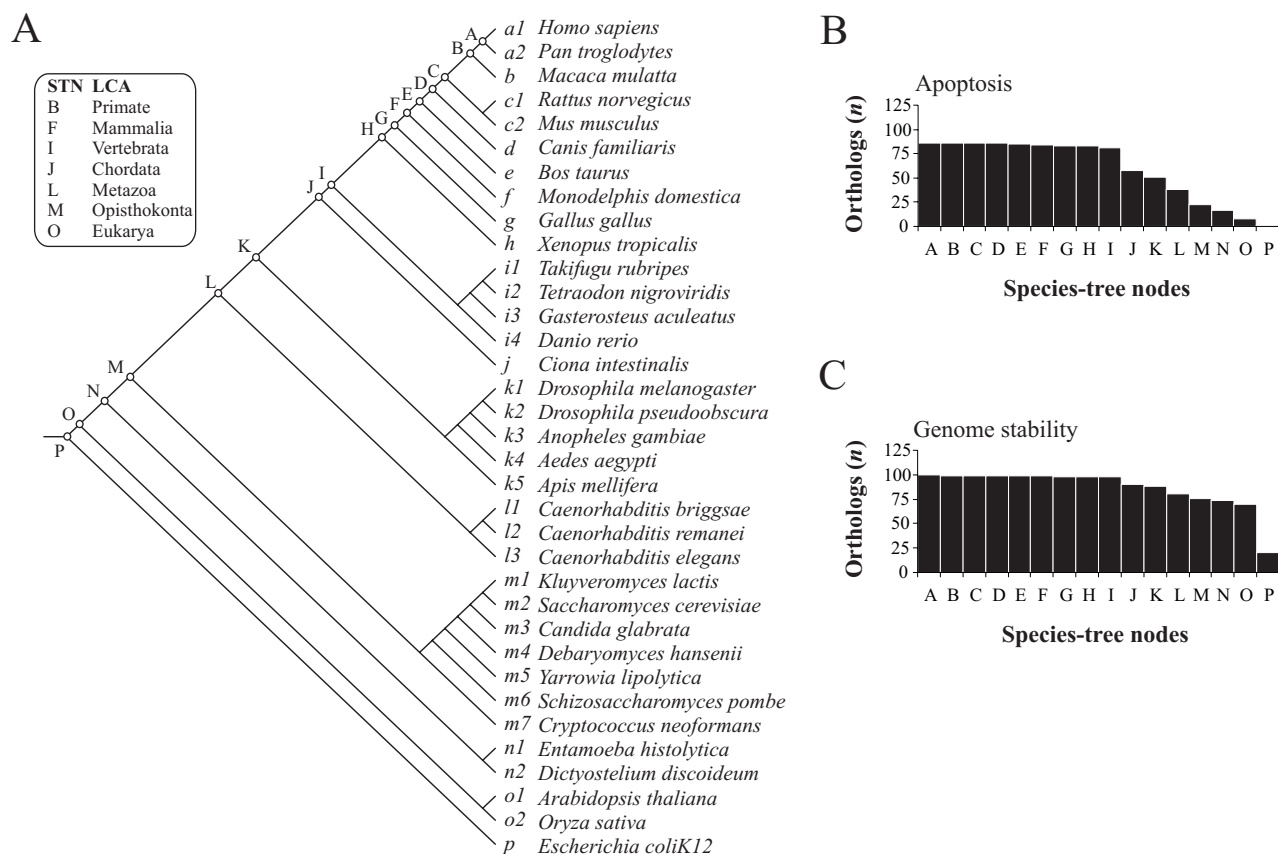
### 1.3. The deep root of eukaryotes

The deep root of eukaryotes is controversial, showing low resolution, and thus should be considered with caution (Doolittle 1999; Baldauf 2003). For instance, *Giardia lamblia* is an extant species that diverged in the consensus tree at STN-Q (eukarya). Due to this early divergence, *Giardia lamblia* carries important information on the evolutionary tree. However, *Giardia* is a parasite, prone to gene loss (Best et al. 2004). This means that its ancestors may have had more genes than their parasite descendants. In this case information is lost on the number of orthologs inferred in the basal position of the species tree. To bypass this limitation it has been suggested that *Giardia lamblia* and other related organisms should be jointly considered to constitute the representative lineage of the descendent species: the gene loss would mutually compensate if genes were lost differentially by different species (Makarova et al. 2005). Accordingly, in a combined scenario (*e.g.* STN-P and -Q), the contrast described in the paper between apoptosis and genome stability pathways could be extended to the origins of all eukaryotes in the species tree. Additional evidence can be inferred considering the likely origin of the ancestral eukaryotic KOGs by identifying their closest prokaryotic orthologous groups (COGs). The KOG-to-COG correspondence is presented in **Supplementary Figure S8** and **Supplementary Table S2**, and shows that 77.0 % of the genome stability orthologs have identifiable prokaryotic orthologous groups, against 39.5% for apoptotic orthologous genes. Furthermore, this scenario is consistent with the origins described for a small set of COGs that can be traced back to the universal ancestor and that recapitulate the three-domain phylogeny (Harris et al. 2003), since in our list only eukaryotic orthologous groups associated with genome

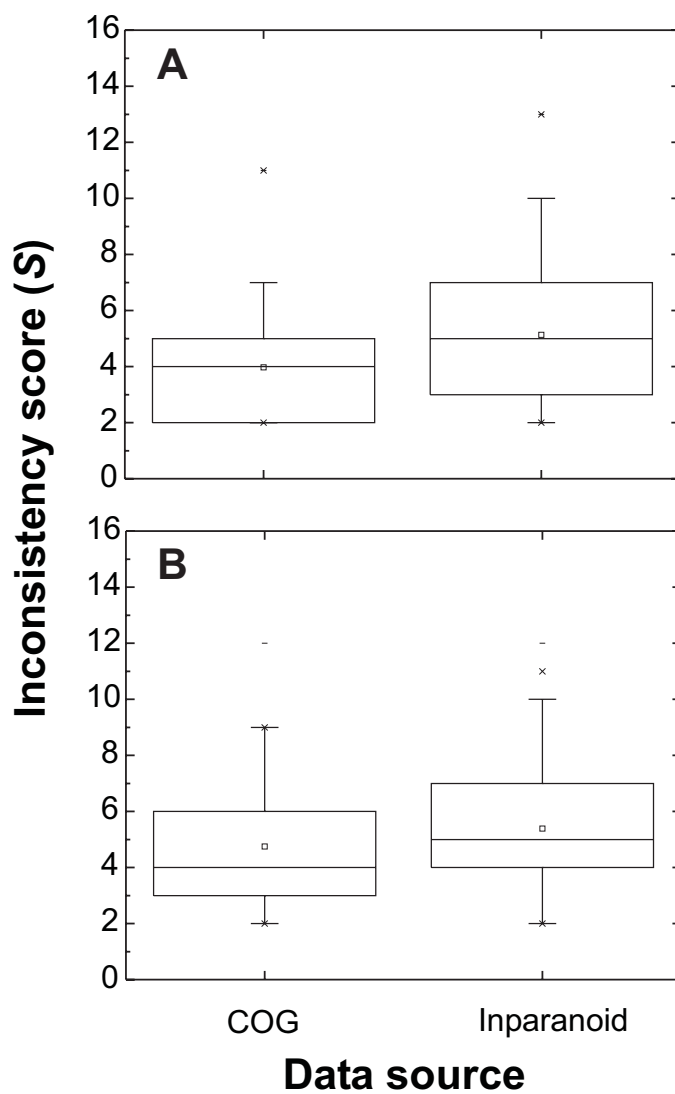
stability functions match these universally conserved COGs ([Supplementary Figure S8](#), red stripes).

## 1.4. The deep root of metazoans

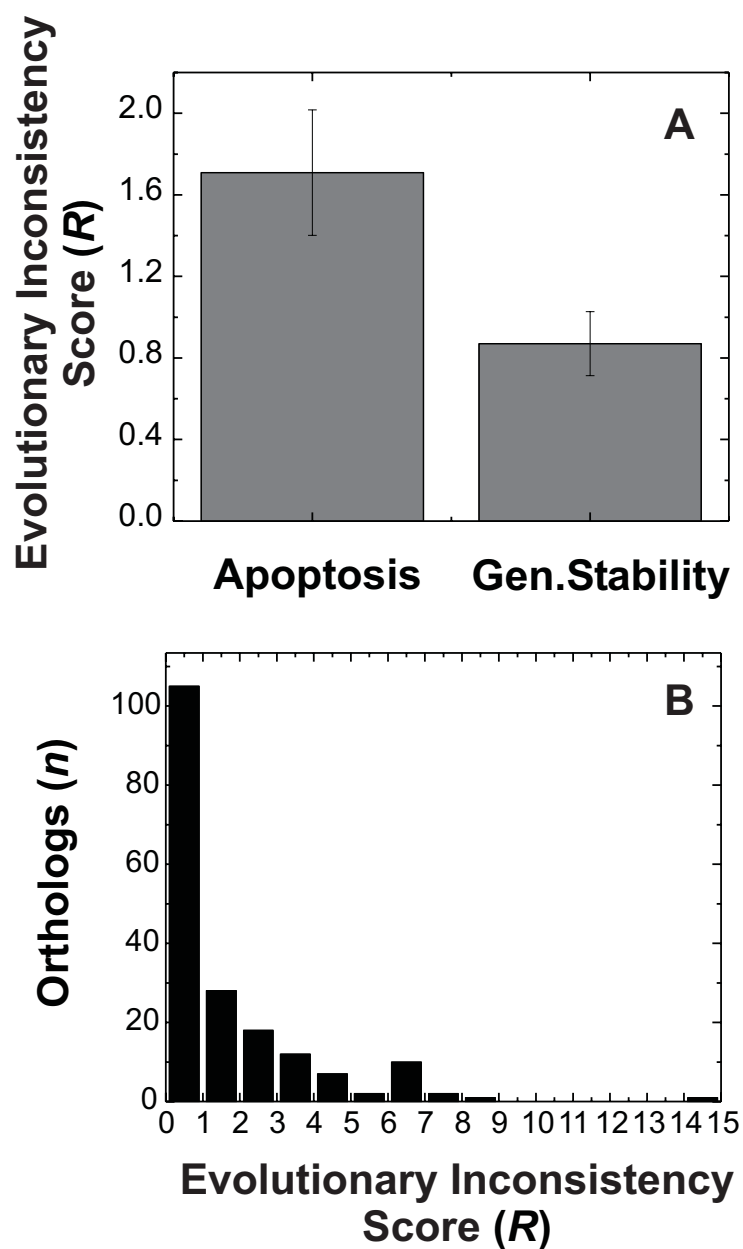
The deep root of metazoan is another controversial issue. However, the current work do not propose the phylogeny of the organisms; instead, here is considered an integration of a variety of phylogenies (Katinka et al. 2001; Pennisi 2003; Baldauf 2003; Delsuc et al. 2005; Ciccarelli et al. 2006; Letunic and Bork 2007). As mentioned in the *Material and Methods*, for each orthologous group associated with the human apoptosis and genome stability genes, our problem is how to find the earliest ortholog in the eukaryote phylogeny (*i.e.* given the species tree, where the human gene roots is placed). Therefore, we focused in the interpretation of our data given the already described species tree topology. Despite these limitations, a different phylogeny in which *C.elegans* isn't at the metazoan root would be welcome, since this species tree node showed the origin of many orthologs of the human apoptosis gene network, which could potentially affect the interpretation of our results. Hence, we carefully included *Nematostella vectensis* as a subsample, which thus changes the base of metazoa ([Supplementary Figure S9](#)). Comparing this figure with the original analysis (**Figure 2B**), the distribution of orthologs remains almost the same (the complete analysis is available at **Supplementary Table S5**). The presence or absence in *Nematostella* was unknown for 34 KOGs, in which case a change from *C. elegans* was assumed (and in each case this meant assuming the presence of the KOG absent in *C. elegans*). The result after this process is that only 1 KOG was found to be different between *C. elegans* and *Nematostella*, however, even assuming all 34 KOGs where *Nematostella* data were unavailable were also different from *C. elegans* did not change the overall pattern (see [Supplementary Figure S9](#)).



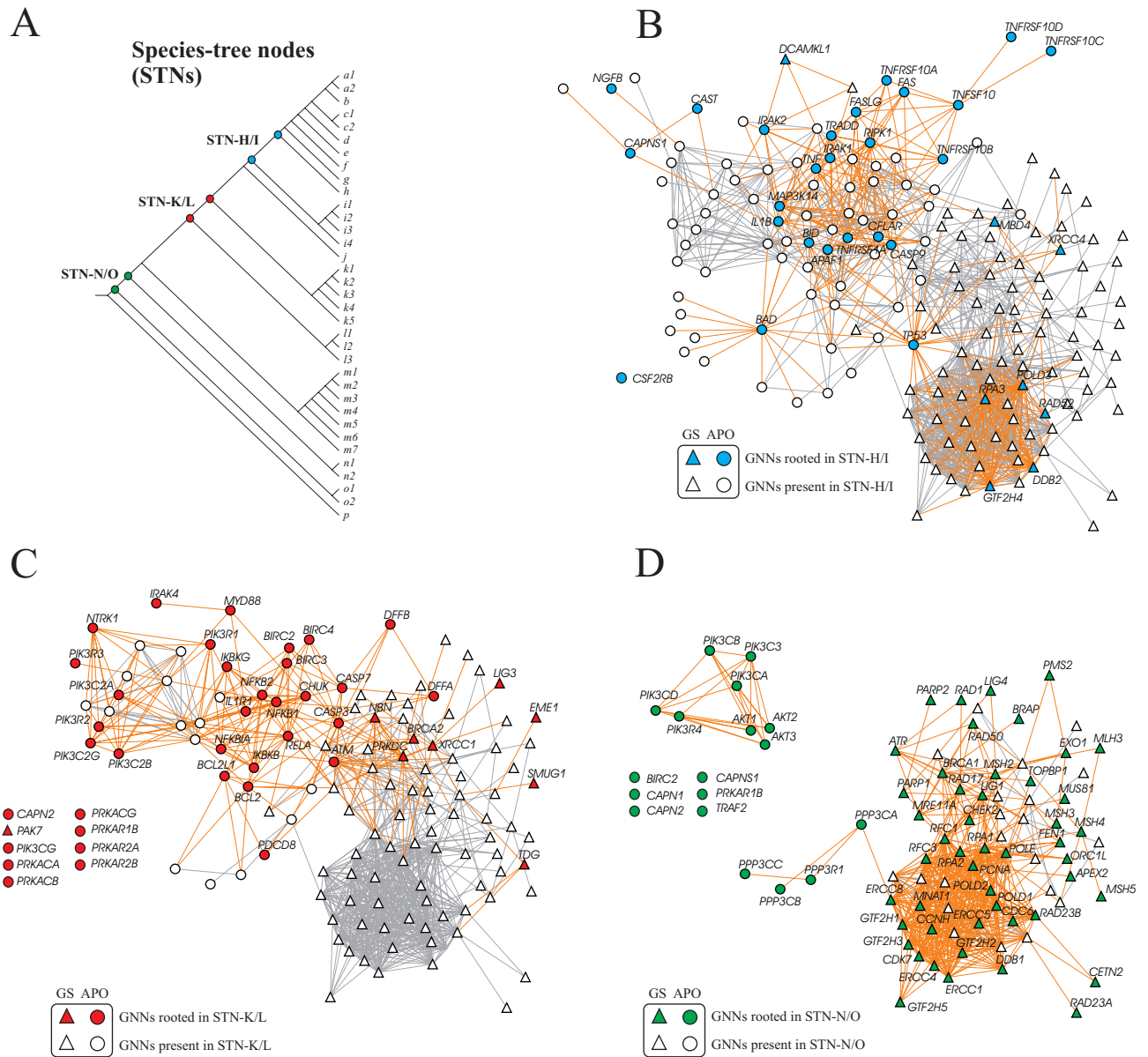
**Supplementary Figure S3.** Inferring ancestral state of human apoptosis and genome stability genes according to Inparanoid Database. **(A)** Eukaryote species tree topology used in the parsimony analysis. Species-Tree Nodes (STNs) and the corresponding Last Common Ancestor (LCA) are indicated. **(B)** Distribution of apoptosis orthologs according to the roots inferred in the species tree. **(C)** Distribution of genome stability orthologs, as in B.



**Supplementary Figure S4.** Inconsistency scores  $S$  estimated from KOG and Inparanoid orthology analysis. Apoptosis (**A**) and genome stability gene set (**B**). For the entire gene set,  $S=4.39(\pm 2.07$ ; KOG analysis), and  $S=5.27(\pm 2.48$ ; Inparanoid analysis).

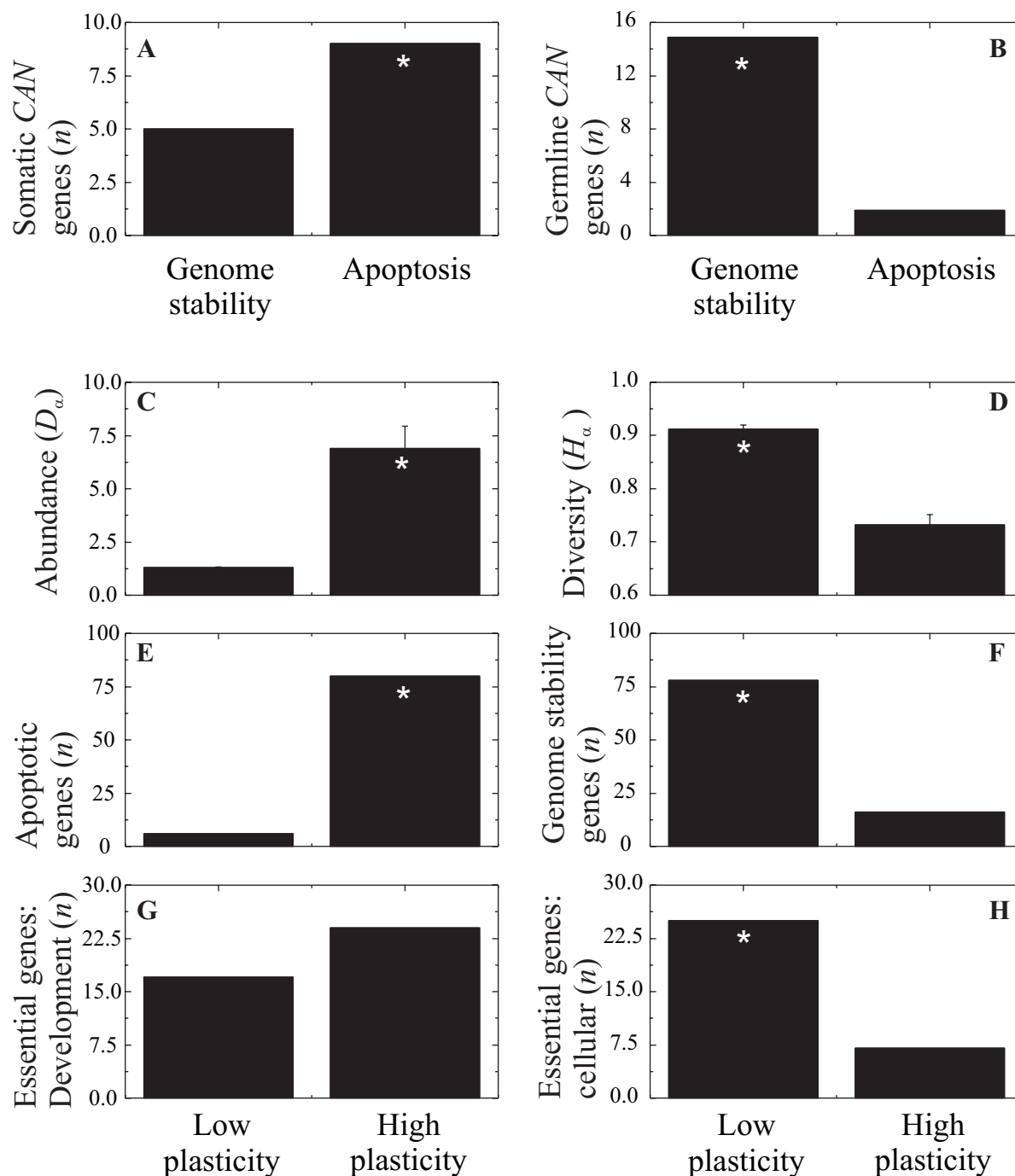


**Supplementary Figure S5.** Comparing the evolutionary roots inferred from COG and Inparanoid orthology analysis. (A) Apoptosis and genome stability gene sets: evolutionary inconsistency score  $R$  is presented as mean  $\pm$ SEM. (B) Frequency distribution of orthologs for the entire gene set according to  $R$ .



**Supplementary Figure S6.** Inparanoid evolutionary scenarios: from species-tree nodes (STNs) to gene-network nodes (GNNs). The orthology information from STNs (**A**) is overlaid on GNNs (**B**, **C** and **D**), which are arranged to optimally display the intersection within the species tree. Roots of an ortholog: color nodes; presence of an ortholog: white nodes. B: Projection of STN-H/I; C: Projection of STN-K/L; D: Projection of STN-O/N.





**Supplementary Figure S7.** Summary of cancer statistics. Contrast between gene functions according to the presence of somatic (A) and germline (B) *CAN* genes. In order to observe the main results present in the paper, we also show the contrasts between plasticity data according to abundance  $D_\alpha$  (C), diversity  $H_\alpha$  (D), gene function (E, F), and gene essentiality (G, H). Network plasticity groups follow the categories defined in Fig.4C: Low plasticity groups contain those genes described in class-a gene set, while high plasticity groups contain those genes described in class-a/c gene set. \* $P < 0.05$  (Mann-Whitney  $U$  Test).

# Apoptosis

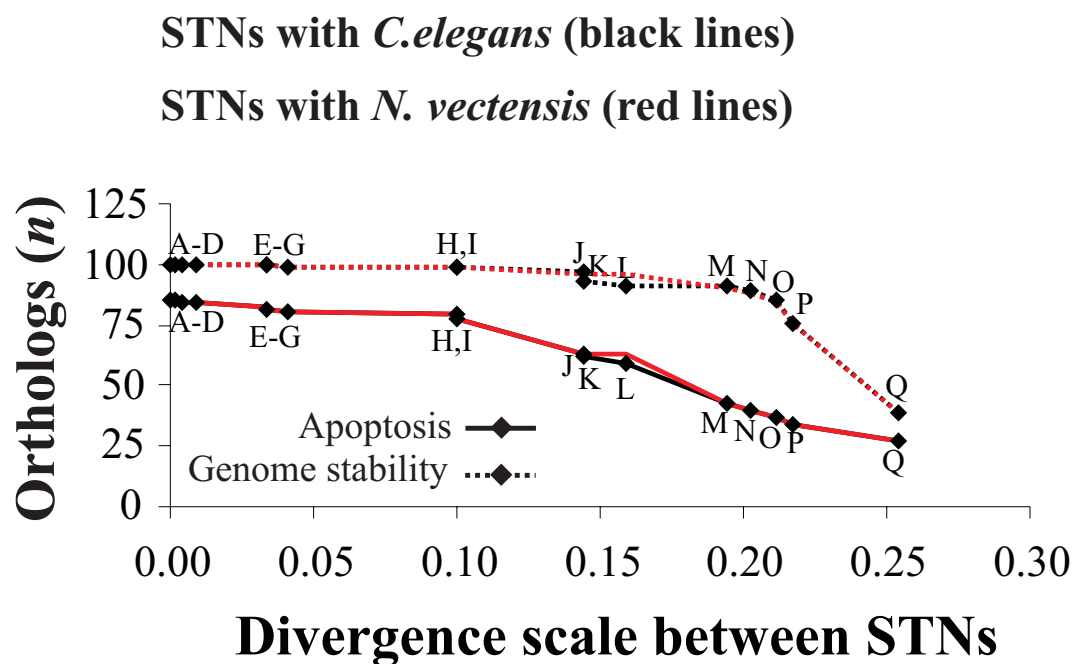
Gene Symbol	KOG-toCOG	Gene Symbol	KOG-toCOG	Gene Symbol	KOG-toCOG	Gene Symbol	KOG-toCOG
BAD	NOG10429	IL1A	NOG08276	AKT1	KOG0690 / COG0515	PRKAR2B	KOG1113 / COG0664
BAX	NOG16176	IL1B	NOG09287	AKT2	KOG0690 / COG0515	RIPK1	KOG0192 / COG0515
BCL2	KOG4728	IL1R1	NOG04905	AKT3	KOG0690 / COG0515	PIK3C2A	KOG0905 / COG5032
BCL2L1	KOG4728	IL1RAP	KOG4641	APAF1	KOG4155 / COG2319	PIK3C2B	KOG0905 / COG5032
BID	NOG10549	IL3	NOG21098	ATM	KOG0892 / COG5032	PIK3C3	KOG0906 / COG5032
BIRC2	KOG1101	IL3RA	NOG21096	CHUK	KOG4250 / COG0515	PIK3CB	KOG0904 / COG5032
BIRC3	KOG1101	MYD88	NOG07319	CYCS	KOG3453 / COG3474	PIK3CD	KOG0904 / COG5032
BIRC4	KOG1101	NGFB	NOG07040	IKBKB	KOG4250 / COG0515	PIK3CG	KOG0904 / COG5032
CAPN1	KOG0045	PIK3R1	KOG4637	IRAK1	KOG1187 / COG0515		
CAPN2	KOG0045	PIK3R2	KOG4637	IRAK2	KOG1187 / COG0515		
CAPNS1	KOG0037	PIK3R3	KOG4637	IRAK4	KOG1187 / COG0515		
CASP10	KOG3573	PIK3R4	KOG1240	MAP3K14	KOG0198 / COG0515		
CASP3	KOG3573	PIK3R5	NOG05250	NFKB1	KOG0504 / COG0666		
CASP6	KOG3573	PPP3R1	NOG04634	NFKB2	KOG0504 / COG0666		
CASP7	KOG3573	PRKAR1B	KOG1113	NFKBIA	KOG0504 / COG0666		
CASP8	KOG3573	RELA	NOG04893	NTRK1	KOG1026 / COG0515		
CASP9	KOG3573	TNF	NOG08240	PDCD8	KOG1346 / COG0446		
CAST	KOG1181	TNFRSF10A	NOG13097	PIK3C2G	KOG0905 / COG5032		
CFLAR	KOG3573	TNFRSF10B	NOG13097	PIK3CA	KOG0904 / COG5032		
CSF2RB	NOG04828	TNFRSF10C	-	PPP3CA	KOG0375 / COG0639		
DFFA	NOG04137	TNFRSF10D	NOG36564	PPP3CB	KOG0375 / COG0639		
DFFB	NOG05130	TNFRSF1A	NOG06963	PPP3CC	KOG0375 / COG0639		
FADD	NOG10546	TNFSF10	NOG08316	PRKACA	KOG0616 / COG0515		
FAS	NOG10520	TP53	NOG07483	PRKACB	KOG0616 / COG0515		
FASLG	NOG07555	TRADD	NOG04168	PRKACG	KOG0616 / COG0515		
IKBKG	KOG0161	TRAF2	KOG0297	PRKAR2A	KOG1113 / COG0664		

KOGs that have a counterpart in prokaryotic groups  
 KOGs without a counterpart in prokaryotic groups  
 KOGs that match universally conserved COGs

# Genome stability

Gene Symbol	KOG-toCOG	Gene Symbol	KOG-toCOG	Gene Symbol	KOG-toCOG	Gene Symbol	KOG-toCOG
BRAP	KOG0804	APEX1	KOG1294 / COG0708	LIG1	KOG0967 / COG1793	PRKDC	KOG0891 / COG5032
BRCA1	KOG4362	APEX2	KOG1294 / COG0708	LIG3	KOG4437 / COG1793	RAD17	KOG1970 / COG0470
BRCA2	KOG4751	ATR	KOG0890 / COG5032	LIG4	KOG0966 / COG1793	RAD23A	KOG0011 / COG5272
DDB1	KOG1897	CCNH	KOG2496 / COG5333	MILH1	KOG1979 / COG0323	RAD23B	KOG0011 / COG5272
EME1	NOG05923	CDC6	KOG2227 / COG1474	MLH3	KOG1977 / COG0323	RAD50	KOG0962 / COG0419
GTF2H1	KOG2074	CDK7	KOG0659 / COG0515	MNAT1	KOG3800 / COG5220	RAD51	*KOG1434 / COG0468
GTF2H5	KOG3451	CETN2	KOG0028 / COG5126	MPG	KOG4486 / COG2094	RAD52	KOG4141 / COG5055
MBD4	KOG4161	CHEK1	KOG0590 / COG0515	MRE11A	KOG2310 / COG0420	RAD54B	KOG0390 / COG0553
NBN	NOG06900	CHEK2	KOG0615 / COG0515	MSH2	KOG0219 / COG0249	RAD54L	KOG0390 / COG0553
PARP1	KOG1037	DCAMKL1	KOG3757 / COG0515	MSH3	KOG0218 / COG0249	RFC1	KOG1968 / COG0470
PARP2	KOG1037	DCLRE1C	KOG1361 / COG1236	MSH4	KOG0220 / COG0249	RFC3	KOG2035 / COG0470
POLD3	NOG05166	DDB2	KOG4328 / COG2319	MSH5	KOG0221 / COG0249	RFC4	KOG0989 / COG0470
POLD4	NOG12441	DMC1	*KOG1434 / COG0468	MSH6	KOG0217 / COG0249	RFC5	KOG0990 / COG0470
RAD1	KOG3194	ERCC1	KOG2841 / COG5241	MUS81	KOG2379 / COG1948	RPA1	KOG0851 / COG1599
RPA3	NOG10964	ERCC2	KOG1131 / COG1199	MUTYH	KOG2457 / COG1194	RPA2	KOG3108 / COG5235
SHFM1	KOG4764	ERCC3	KOG1123 / COG1061	NTHL1	KOG1921 / COG0177	TDG	KOG4120 / COG3663
SMUG1	NOG04402	ERCC4	KOG0442 / COG1948	OGG1	KOG2875 / COG0122	UNG	KOG2994 / COG0692
TOPBP1	KOG1929	ERCC5	*KOG2520 / COG0258	ORC1L	KOG1514 / COG1474	XPA	KOG4017 / COG5145
XAB2	KOG2047	ERCC6	KOG0387 / COG0553	PAK7	KOG0578 / COG0515	XPC	KOG2179 / COG5535
XRCC1	KOG3226	ERCC8	KOG4283 / COG2319	PCNA	*KOG1636 / COG0592	PIK3C2A	KOG0905 / COG5032
XRCC4	NOG07367	EXO1	*KOG2518 / COG0258	PMS2	KOG1978 / COG0323	PIK3C2B	KOG0905 / COG5032
XRCC5	KOG2326	FEN1	*KOG2519 / COG0258	PMS2L3	KOG1978 / COG0323	PIK3C3	KOG0906 / COG5032
XRCC6	KOG2327	GTF2H2	KOG2807 / COG5151	PNKP	KOG2134 / COG0241	PIK3CB	KOG0904 / COG5032
		GTF2H3	KOG2487 / COG5242	POLD1	KOG0969 / COG0417	PIK3CD	KOG0904 / COG5032
		GTF2H4	KOG3471 / COG5144	POLD2	KOG2732 / COG1311	PIK3CG	KOG0904 / COG5032

**Supplementary Figure S8.** KOG-to-COG correspondence. Orthology information of the Clusters of Orthologous Groups (COGs) was retrieved through the KOG-to-COG assignments in the STRING database. Red stripes highlight the eukaryotic groups that match the universally conserved prokaryotic groups in the KOG-to-COG correspondence according to Harris *et al.* (2003).



**Supplementary Figure S9.** Subsample of the data where *C. elegans* isn't at the root of metazoa. This figure presents a comparison between the original analysis (see **Figure 2** ) and an alternative construction where *Nematostella vectensis* is manually placed in the position of *C.elegans*, according to the same species-tree nodes (STNs). Orthology data derive from *Nematostella vectensis* genome web site (<http://genome.jgi-psf.org/Nemve1/Nemve1.home.html>). (Data available at **Supplementary Table S5**).

```

#Homo_sapiens      MVNVLADALK SINNAEKRKG RQVLIRPCSK VIVRFLTVMM KHGYIGEFEI IDHRAGKIV VNLTRGLNKC GVISPRFDVQ LKDLEKWQNN LLPSRQFGFI VLTTASAGIMD HEEARRKHTG GKILGFF
#Pan_troglodytes  .....
#Macaca_mulatta   .....
#Rattus_norvegicus .....
#Mus_musculus     .....
#Canis_familiaris ..... V..... I..... D..... I.....
#Bos_taurus       .....
#Monodelphis_domestica .....
#Gallus_gallus    ..... Y.....
#Xenopus_tropicalis ..... KCLTAR. KALCFISET. FFFSPVDHS. .... Y.....
#Takifugu_rubripes ..... L..... Y.....
#Tetraodon_nigroviridis ..... L..... Y.....
#Danio_rerio      ..... L.....
#Ciona_intestinalis ..... CLC..... L..... V.K..... A.N.T.V..... E.N..... A IR... TT..... Y..... C.....
#Drosophila_melanogaster ..... C..... L..... IK..... V... S..... P IN.I... T..... YV..... G..... L.....
#Anopheles_gambiae ..... S..... C..... N..... VIK..... V... S.V..... A I..... A N.I.R.T..... YV..... G..... L..... Y.....
#Apis_mellifera   ..... S..... L..... IK..... RK..... V... S.V..... S..... P IN.I... T..... YV..... G..... L.....
#Caenorhabditis_elegans ..... N A..... A..... A..... V..... A S..... LNIR. N... Y.T. T..... YL I..... L.....
#Kluyveromyces_lactis ..... TS..... N A..... T..... S..... IK. Q. Q..... Y..... S..... Q.N..... N.K IA.V... TA... A... YV I..... H... VS..... V.....
#Saccharomyces_cerevisiae ..... TS..... N A..... T..... S..... IK. Q. Q..... Y..... S..... Q.N..... N.K IG.I... TA... A... YV I..... VS..... V.....
#Candida_glabrata ..... TS..... N A..... T..... S..... IK. Q. Q R..... Y..... S..... Q.N..... N.K IG.I... TA... A... YV I..... VS..... V.....
#Eremothecium_gossypii ..... TS..... N A..... T..... S..... IK. Q. Q..... Y..... S..... Q..... N.K IN.V... TA... A... YV I..... H... VA..... V.....
#Debaryomyces_hansenii ..... TS..... N..... T..... T..... S..... K..... Q..... Y..... S..... Q.N..... N.K VG.M.R.TD... A... YV I..... VA..... V.....
#Yarrowia_lipolytica ..... TS..... Q... T... A..... S..... IK..... Q. Q..... Y V... S..... Q.N..... N.K ID.V... IQ..... Y. I..... R.VA.....
#Aspergillus_fumigatus ..... S.N... N A..... A..... S..... K. S. Q..... E V... S..... IQ.N..... N. YP... P.I.Q.AVQ..... YV..... VA... L.....
#Schizosaccharomyces_pombe ..... S... C.N N.V... R... R..... S..... K..... Q..... D. TE..... S..... IQ.N... I..... N.K... I... V.Q..... V.V..... R... S... N... A.DA.....
#Filobasidiella_neoformans ..... S.N... N N.V... R..... S..... VIKV.S.Q..... G.V..... IQ.N..... N.P VDSI.N.VAQ... A.S.K.I... V.N.V... A.V.....
#Encephalitozoon_cuniculi ..... SAA.NMC.A... SRA... L.L.FSTE EMRM.LQ.L R... SG.SY.H.K.S.SI ID.N... R... A... NYI.K RDGI.DFRTR... A... HV LFN.K.L.K.CLTENV... Y.....
#Dictyostelium_discoideum ..... S.N.C.N V... RQ... V... S... K.E... KR... V... S... ID.I... I... T.DEI... ASY... H... L... N.KTR... L...
#Arabidopsis_thaliana ..... S.N... MY... M... S... IK.I.Q... Y V... S... E.N... G.V.EI.G.TAR... Y... NV... V...
#Cyanidioschyzon_merolae ..... T... S... T... M... A... I.V.RLLQ RE... D.TF... S... Q.I... V A... Y.A IR... Q.CDA... KL.GL IC... S.N DA.M.R.V... LI... Y...
#Plasmodium_falciparum ..... S... C... T... R... S... VIK.QY.Q K... S... V... S... L.I... A... Y.K.DEI.IITS I... L.HL I... PY... V...
#Cryptosporidium_hominis ..... S.S.C... A.L... M.R... S... IK.QC.Q RR... V V.R... C IE.L... Y.P.S.I.QISTD... Y... S.Y... Q...
#Thalassiosira_pseudonana ..... S... C... T.T... S... K.QC.Q QN... V... SQ... IE.N... VHEI... VV... I.S.TF... T... K... V...
#Giardia_lamblia ..... R... C... QRI... K.IV.S... IE... QL.Q N... SD.AV V.N.SNR... I... A... IP AN.I... VV... L.H.I... Q... I.QHRQI... VI.Y...

```

**Supplementary Figure S10.** Comparison of amino acid sequences of 35 eukaryotic 40S ribosomal proteins obtained from a three-domain orthologous groups (COG0096/KOG1754). Amino acid sequences were aligned using ClustalW. All sites containing alignment gaps and missing-information were removed before computing the distance matrix. Dots represent amino acid residues that are identical to the corresponding sites of the first sequence (*i.e.* *H.sapiens* sequence). Protein IDs: *H.sapiens* (ENSP00000318646); *P.troglodytes* (ENSPTRP00000013350); *M.mulatta* (ENSMMUP00000007395); *R.norvegicus* (ENSRNOP00000024678); *M.musculus* (ENSMUSP00000008827); *C.familiaris* (ENSCAFP00000016520); *B.taurus* (ENSBTAP00000027627); *M.domestica* (ENSMODP00000007323); *G.gallus* (ENSGALP00000010942); *X.tropicalis* (ENSXETP00000006430); *T.rubripes* (NEWSINFRUP00000151157); *T.nigroviridis* (GSTENP00022669001); *D.rerio* (ENSDARP00000007879); *C.intestinalis* (ENSCINP00000008963); *D.melanogaster* (CG2033-PE); *A.gambiae* (ENSANGP00000029176); *A.mellifera* (ENSAPMP00000009198); *C.elegans* (F53A3.3.3); *K.lactis* (KLLA0B07601g); *S.cerevisiae* (YJL190C); *C.glabrata* (CAGL0K04587g); *E.gossypii* (AEL151C); *D.hansenii* (DEHA0B06864g); *Y.lipolytica* (YALI0D05731g); *A.fumigatus* (Afu1g15730); *S.pombe* (SPAC22A12.04c); *F.neoformans* (CNK02900); *E.cuniculi* (ECU09\_1350); *D.discoideum* (DDB0167031); *A.thaliana* (AT5G59850.1); *C.merolae* (CMI202C); *P.falciparum* (MAL3P6.30); *C.hominis* (Chro.80500); *T.pseudonana* (26367); *G.lamblia* (15228).

```

#Homo_sapiens      MSTSKTGKHG HAKVHLVGID IFTGKKYEDI CPSTHNMDVP NIKRNDYQLI CIQDYLSSLT ETGEVRDLKL PEGELGKEIE GKYN AEDVQV SVMCAMSEEE AVAIAK
#Pan_troglodytes  .....
#Macaca_mulatta   .....
#Rattus_norvegicus .....
#Mus_musculus     .....
#Canis_familiaris .....
#Bos_taurus       .....F...G.....Q DS.....R. ...D.....Q..DC.EILI T.LS..T..A .....
#Monodelphis_domestica .....
#Gallus_gallus    .....N.....G.....S.....D.....F.....I ..IS..N..C .....
#Xenopus_tropicalis .....C....G.N.....S DS.D....M ...D..R..L M.HDG.EIL. T.LT..N..C ...L.
#Takifugu_rubripes .....S.....M.....G.....D DN.....I .D.D..... N.NE..EILI T.LN..N..A .ISV.
#Tetraodon_nigroviridis .....S.....M.....G.....D DN.....I .D.D..... S.HE..EILI T.LN..N..A .ISV.
#Danio_rerio      .....N.....V D.SEFF...MD DN.D...RV .D.D..... N.FA..EML. T.LS..G..S ...L.
#Ciona_intestinalis .....M....DS.....S.....Q.. ...TEV.VV DVA.QVTIME .D.DT.E..V E.KPMLDS.H KIIESG.CYI TILA.LGR.K VITA.
#Drosophila_melanogaster .....M....SN......V..E.L... A.D.F.T.M. .S.DL...V .....EQLR LDFDSK.LLC T.LK.CG..C VI...
#Anopheles_gambiae .....M.....V..E...T D.D.F.V.MS DN.DL...I .D....TQLR SEFDSKEIVC T.LKSCG..T VI...
#Apis_mellifera   .....S.....FV..E...A D.D..C.MA DN..L...I .D....AQLR ADHE.KELLC T.LK.CG..V VI...
#Caenorhabditis_elegans .....M.A...S..L... ..VV..RE.L.M A.D..C..MD PEC.QK... .DT...QQ.R DA.EK.S.L. Q.VS.IG..A ILGW.
#Kluyveromyces_lactis .....A...N..L..L S.....E.. VV..TEF..L D.D.F..M. MD..TK.VRA .....DN.Q AAFDEG.LM. TIIA..G..A .ISF.
#Saccharomyces_cerevisiae .....TL...L..L S.....LE.. FV..SE...L D.D...M. MD..TK.V.A .....DSMQ AAFDEG.LM. TIIS..G..A .ISF.
#Candida_glabrata .....A...L..L S.....E.. VV..TEF..L D.D.F..M. MD.DTK.VRH .....DQM AAFDEG.LM. TIIA..G..A .ISF.
#Eremothecium_gossypii .....A...L..L S.....E.. VV..TE...L D.D.F..M. MD..TK.VRA .D....ETMQ ASFDEG.LM. T.I.TS.G..A .ISF.
#Debaryomyces_hansenii .....A...L..L S.....E.. .VG..REF..L D.D.F..M. ND.DTK.V.V .....DK.Q SEFDEG.LL. TIVS..G..A .ISF.
#Yarrowia_lipolytica .....AT...N..L..L S.....E.. .V..EEF... D.D.F...F. AD.STK.VP. ....I.DK.Q TEFDEG.LI. T.IS..G..A VISY.
#Aspergillus_fumigatus .....I.AL...L..L S.....E.. .VT..RE...L DVD.F...MD .A.NTK.V.. ....V.ER.Q KMF.EG.CN. TILT..G..QA CMDV.
#Schizosaccharomyces_pombe .....I.AL...N.R..M S.....E.. VV..DE...V N.D...N.M. TD.TTK.VR. ....N... EGFEEG.LII T.VS..G..I .L.CR
#Filobasidiella_neoformans .....A...L..L S.....E.. .VR..QEF..L D...F.N.MD SD.GSK.V.V ..T.I.QQ.M ADFE.G.LM. TIIS..D..Q .ISY.
#Encephalitozoon_cuniculi T.SV.N...A..TTISSKI LS..SNHKG. YTANDSII.C RPEKVQLK.. D.TSTFTDSS GS.DS.DAGRM SSEDKIVQTV EGS.SS.LSL R.LPDFYKLE S.RPS
#Dictyostelium_discoideum .....NITA...S...E.. ....I... .VS..KE.TVM DV.....D AG...K..A. ....DDI...T QMLKEKEPL. ...IS.LGK.G V.SV.
#Arabidopsis_thaliana V.....C.F.A...S..L... V..S..C... HVN.T... D.E..V... DN..STK... .NDD.LQQ.K SGFDDK.LV. ...S..G..Q IN.L.
#Cyanidioschyzon_merolae V.....ANI..L... ..V A.TS.T.YQ. VVV..SE... D.AEFC..MD DK..T.... DPNVHAK.K EDFE.KQLT. V.LK..KQ.K IMQS.
#Plasmodium_falciparum Y.....A.I...R..... .TS..... VV..TEL... D.E.FV...Y DN.DTK..S. .KDTVA.Q.R NLFDNKS.L. .LS.CGQ.K II.A.
#Cryptosporidium_hominis .....A.I..L...N...V .TS...A.. FV..SE... D.DNFT... .N.STK..A .TDAVALQ.Q QM.ADKAIL. G.LA.CGI.K I.A.
#Thalassiosira_pseudonana I.V.....CNFTAV...L..M V..S.GTT.. IVT.TEWEI. D.EEE.T.MD .G.NQK.VN. .GVPMQA..R DAW.D.Q.S. T.QA.VGQ.Q V..Y.
#Giardia_lamblia I.....CSITAV...S...S..IMC. .VS..T..L.. N.D..AFYFD QNA.QQR.P. ENE..VADLR AAFE.DTIMI NIQR..NT.G I..F.

```

**Supplementary Figure S11.** Comparison of amino acid sequences of 35 eukaryotic translation initiation factor 5A proteins obtained from a three-domain orthologous groups (KOG3271/COG0231). Amino acid sequences were aligned using ClustalW. All sites containing alignment gaps and missing information were removed before computing the distance matrix. Dots represent amino acid residues that are identical to the corresponding sites of the first sequence (*i.e.* *H.sapiens* sequence). Protein IDs: *H.sapiens* (ENSP00000295822); *P.troglodytes* (ENSPTRP00000026874); *M.mulatta* (ENSMMP00000030130); *R.norvegicus* (ENSRNOP00000015779); *M.musculus* (ENSMUSP00000050289); *C.familiaris* (ENSCAFP00000022031); *B.taurus* (ENSBTAP00000002616); *M.domestica* (ENSMODP00000020009); *G.gallus* (ENSGALP00000038570); *X.tropicalis* (ENSXETP00000055749); *Trubripes* (NEWSINFRUP00000158888); *T.nigroviridis* (GSTENP00019223001); *D.rerio* (ENSDARP00000027654); *C.nintestinalis* (ENSCINP00000012334); *D.melanogaster* (CG3186-PA); *A.gambiae* (ENSANGP00000015032); *A.mellifera* (ENSAPMP00000024661); *C.elegans* (F54C9); *K.lactis* (KLLA0E22286g); *S.cerevisiae* (YJR047C); *C.glabrata* (CAGL0L01353g); *E.gossypii* (AFR356C); *D.hansenii* (DEHA0F13640g); *Y.ipolytica* (YALIO06886g); *A.umigatus* (Afu1g04070); *S.pombe* (SPBC336); *F.neoformans* (CND05400); *E.cuniculi* (ECU09\_1370); *D.discoideum* (DDB0191442); *A.thaliana* (AT1G13950); *C.merolae* (CMS351C); *P.falciparum* (PFL0210c); *C.hominis* (Chro.70262); *T.pseudonana* (158382); *G.lamblia* (14614).

```

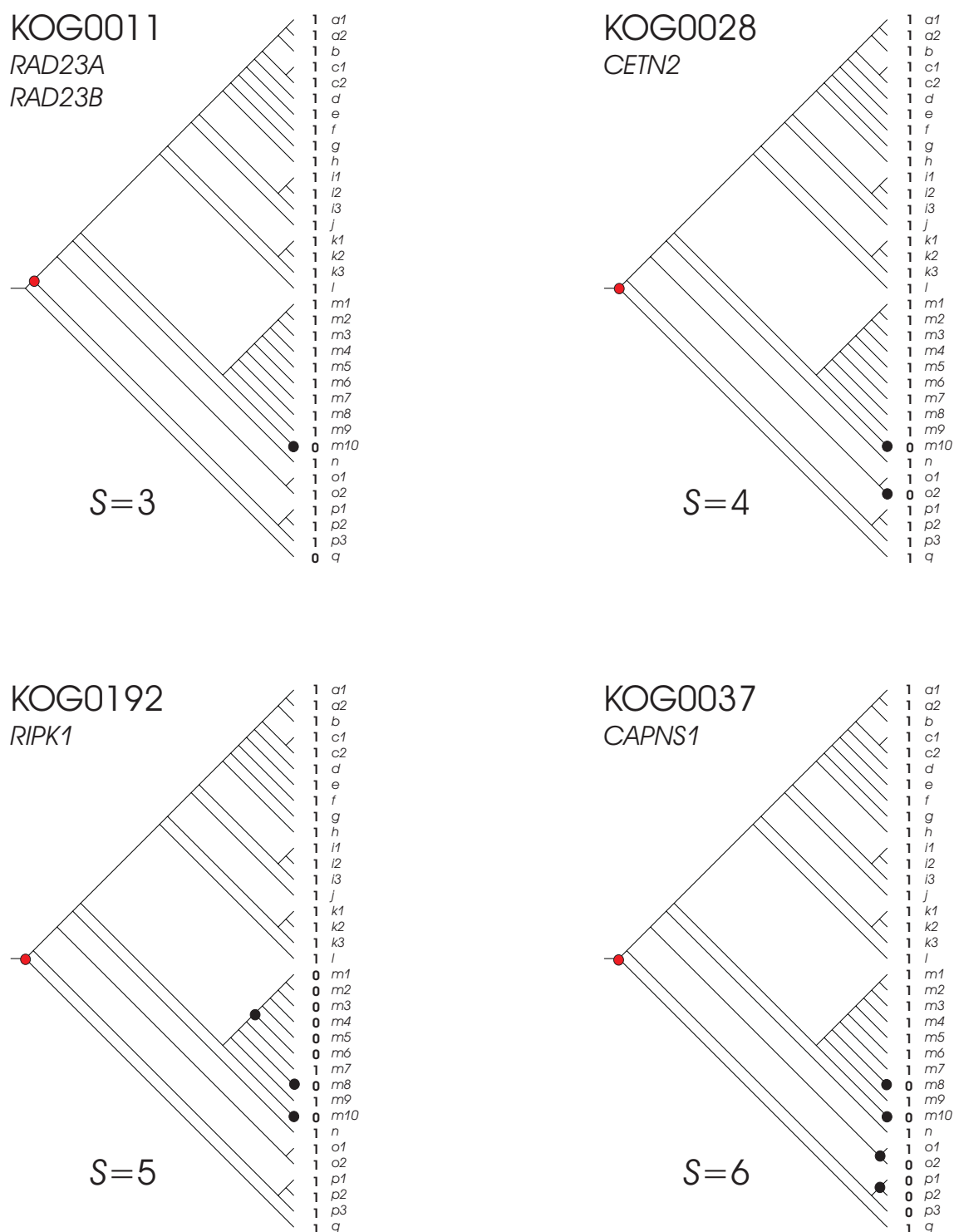
#Homo_sapiens      MSIQYFLIAV  GGLQNEEGET  SHLMGMFYRT  IRMMENGIKP  VYVFDGKPPQ  LKSGEKKIQE  FHLSRILQEL  GLNQEQFVDL  CILLGSDYCE  SIRGIGPKRA  VDLIQKHKSI  EEIVRRLDYP  VPELHKEAHO  LFLEVLDELK  WSEPNEEELI  KFMCGEKQFS  EERISRGVKR  LSKSRQGGQR  LDDFFKVTG
#Pan_troglodytes
#Macaca_mulatta
#Rattus_norvegicus
#Mus_musculus
#Canis_familiaris
#Bos_taurus
#Monodelphis_domestica
#Gallus_gallus
#Xenopus_tropicalis
#Takifugu_rubripes
#Tetraodon_nigroviridis
#Danio_rerio
#Ciona_intestinalis
#Drosophila_melanogaster
#Anopheles_gambiae
#Apis_mellifera
#Caenorhabditis_elegans
#Kluyveromyces_lactis
#Saccharomyces_cerevisiae
#Candida_glabrata
#Eremothecium_gossypii
#Debaryomyces_hansenii
#Yarrowia_lipolytica
#Aspergillus_fumigatus
#Schizosaccharomyces_pombe
#Filobasidiella_neoformans
#Encephalitozoon_cuniculi
#Dictyostelium_discoideum
#Arabidopsis_thaliana
#Cyanidioschyzon_merolae
#Plasmodium_falciparum
#Cryptosporidium_hominis
#Thalassiosira_pseudonana
#Giardia_lamblia

```

**Supplementary Figure S12.** Comparison of amino acid sequences of 35 eukaryotic 5-3 exonucleases (Flap Structure-specific Endonuclease 1 like proteins - FEN1) obtained from a three-domain orthologous groups (KOG2519 / COG0258). Amino acid sequences were aligned using ClustalW. All sites containing alignment gaps and missing-information were removed before computing the distance matrix. Dots represent amino acid residues that are identical to the corresponding sites of the first sequence (*i.e.* *H.sapiens* sequence). Protein IDs: *H.sapiens* (ENSP00000305480); *P.troglodytes* (ENSPTRP00000006454); *M.mulatta* (ENSMMUP00000008754); *R.norvegicus* (ENSRNOP00000027842); *M.musculus* (ENSMUSP00000025651); *C.familiaris* (ENSCAFP00000023634); *B.taurus* (ENSBTAP00000000071); *M.domestica* (ENSMODP00000027077); *G.gallus* (ENSGALP00000005724); *X.tropicalis* (ENSXETP00000014663); *T.rubripes* (ENSTRUP00000032042); *T.nigroviridis* (GSTENP00004156001); *D.rerio* (ENSDARP00000004016); *C.ntestinalis* (ENSCINP000000004910); *D.melanogaster* (FBpp0086223); *A.gambiae* (AGAP011448-PA); *A.mellifera* (ENSAPMP00000010885); *C.elegans* (Y47G6A.8); *K.lactis* (KLLA0F02992g); *S.cerevisiae* (YKL113C); *C.glabrata* (CAGL0K11506g); *E.gossypii* (ABL052C); *D.hansenii* (DEHA0F15059g); *Y.ipolytica* (YALIOF20042g); *A.umigatus* (Afu3g06060); *S.pombe* (SPAC3G6.06c); *F.neoformans* (CND01190); *E.cuniculi* (ECU03\_1080); *D.discoideum* (DDB0186301); *A.thaliana* (AT5G26680.1); *C.merolae* (CMG106C); *P.falciparum* (PFD0420c); *C.hominis* (Chro.70245); *T.pseudonana* (132436); *G.lamblia* (16953).

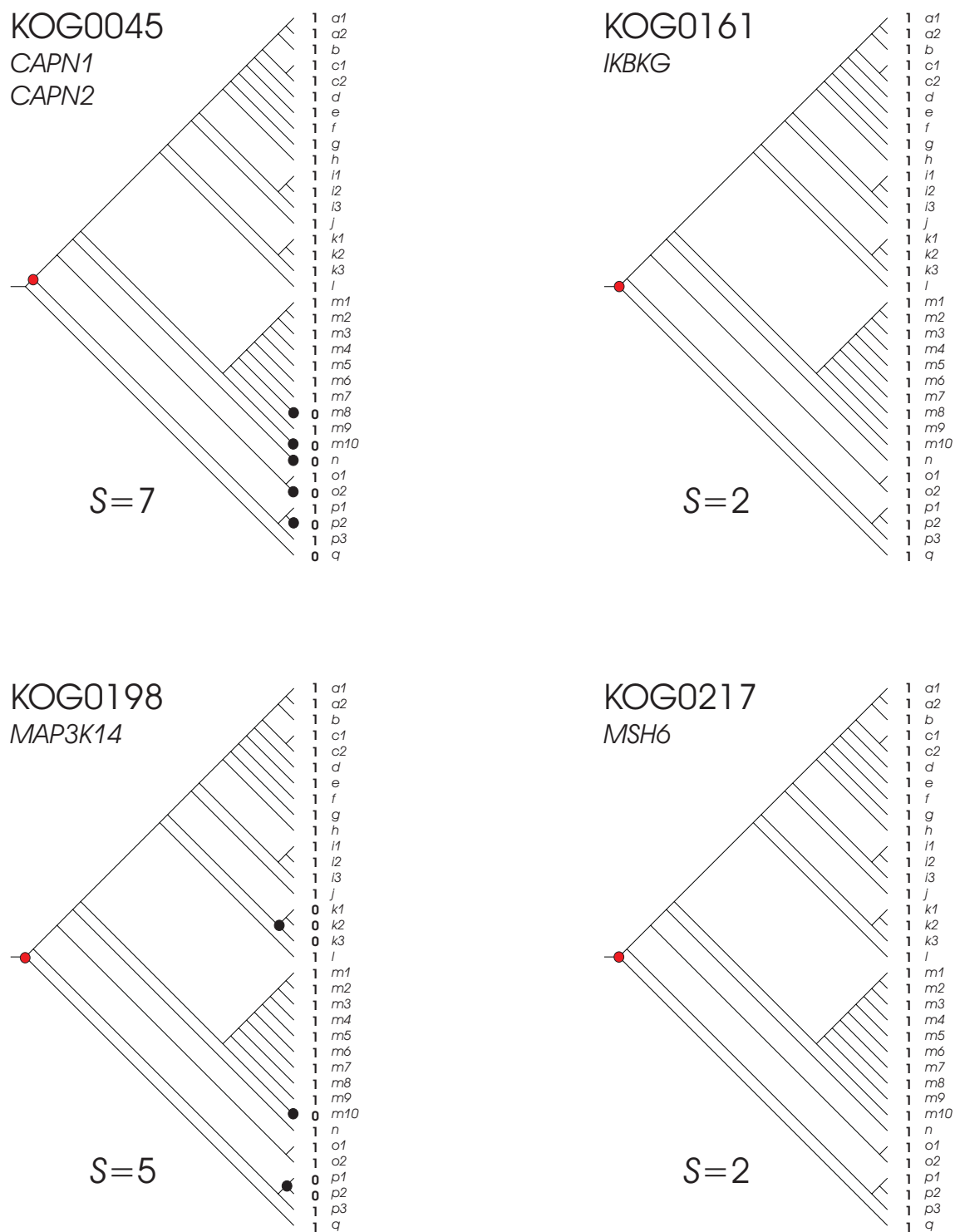




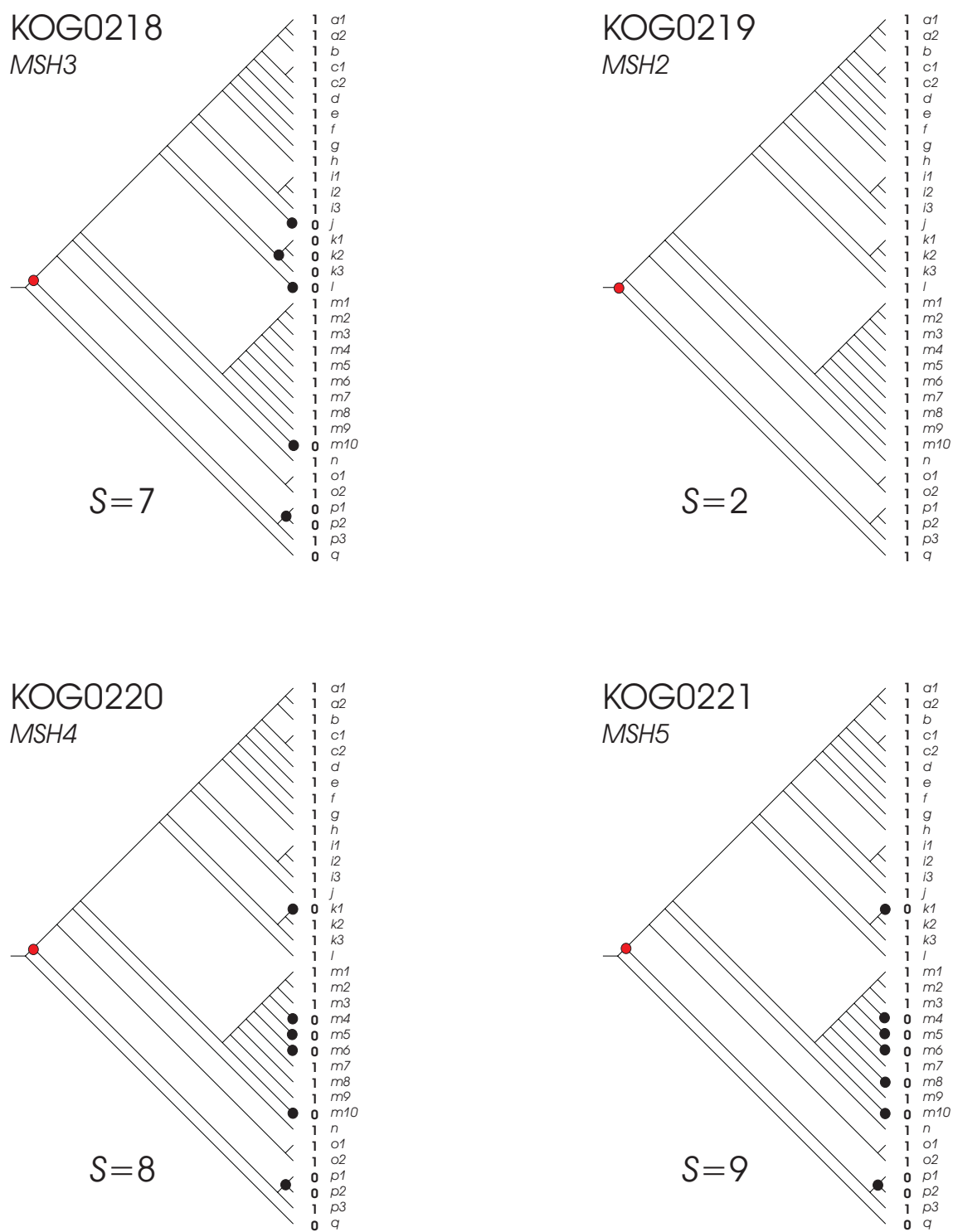


**Supplementary Figure S14.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0011, KOG0028, KOG0192 and KOG0037. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes: **a1** (*Homo sapiens*); **a2** (*Pan troglodytes*); **b** (*Macaca mulatta*); **c1** (*Rattus norvegicus*); **c2** (*Mus musculus*); **d** (*Canis familiaris*); **e1** (*Bos Taurus*); **f** (*Monodelphis domestica*); **g** (*Gallus gallus*); **h** (*Xenopus tropicalis*); **i1** (*Takifugu rubripes*); **i2** (*Tetraodon nigroviridis*); **i3** (*Danio rerio*); **j** (*Ciona intestinalis*); **k1** (*Drosophila melanogaster*); **k2** (*Anopheles gambiae*); **k3** (*Apis mellifera*); **l** (*Caenorhabditis elegans*); **m1** (*Kluyveromyces lactis*); **m2** (*Saccharomyces cerevisiae*); **m3** (*Candida glabrata*); **m4** (*Eremothecium gossypii*); **m5** (*Debaryomyces hansenii*); **m6** (*Yarrowia lipolytica*); **m7** (*Aspergillus fumigatus*); **m8** (*Schizosaccharomyces pombe*); **m9** (*Filobasidiella neoformans*); **m10** (*Encephalitozoon cuniculi*); **n** (*Dictyostelium discoideum*); **o1** (*Arabidopsis thaliana*); **o2** (*Cyanidioschyzon merolae*); **p1** (*Plasmodium falciparum*); **p2** (*Cryptosporidium hominis*); **p3** (*Thalassiosira pseudonana* CCMP1335); **q** (*Giardia lamblia* ATCC 50803).

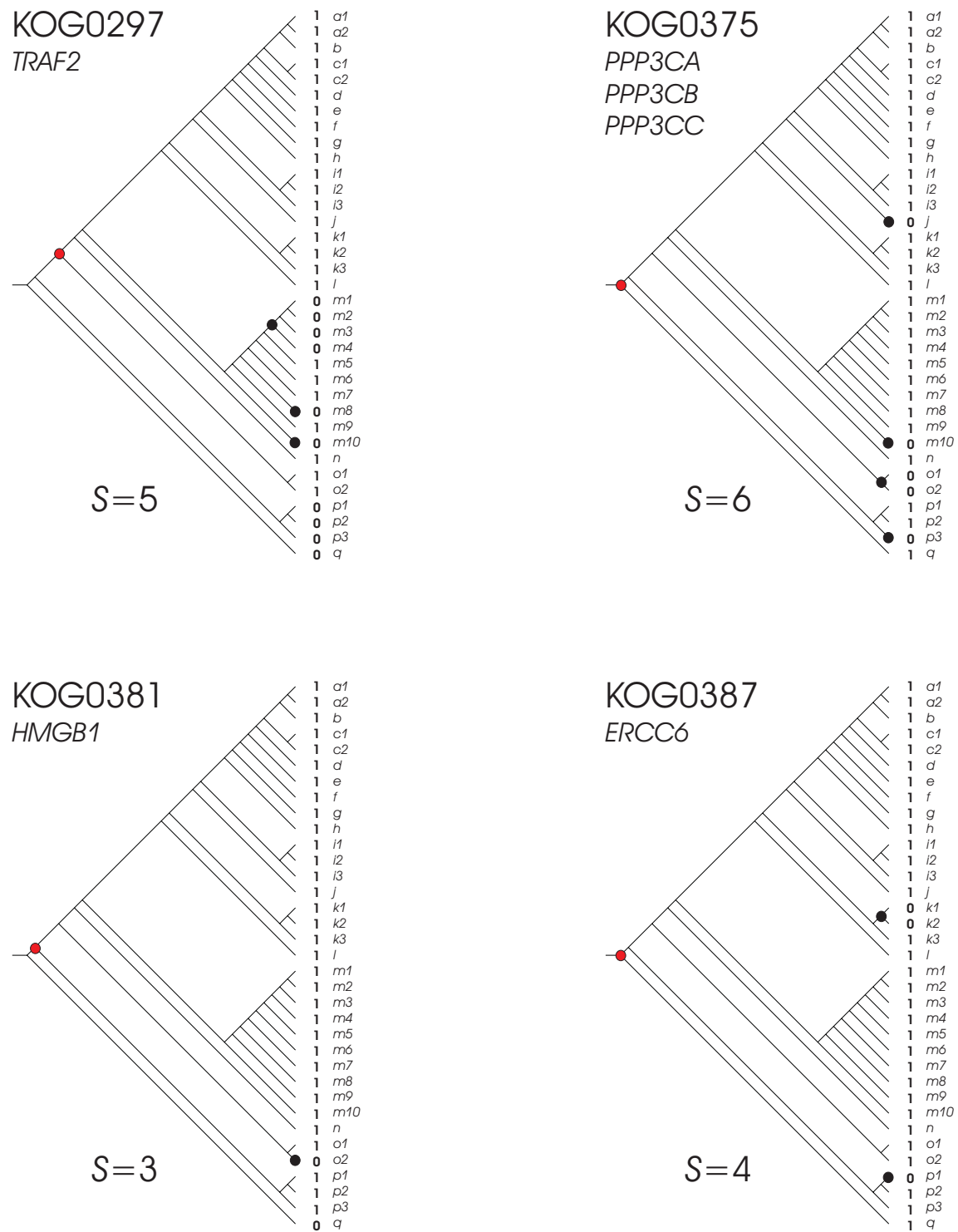




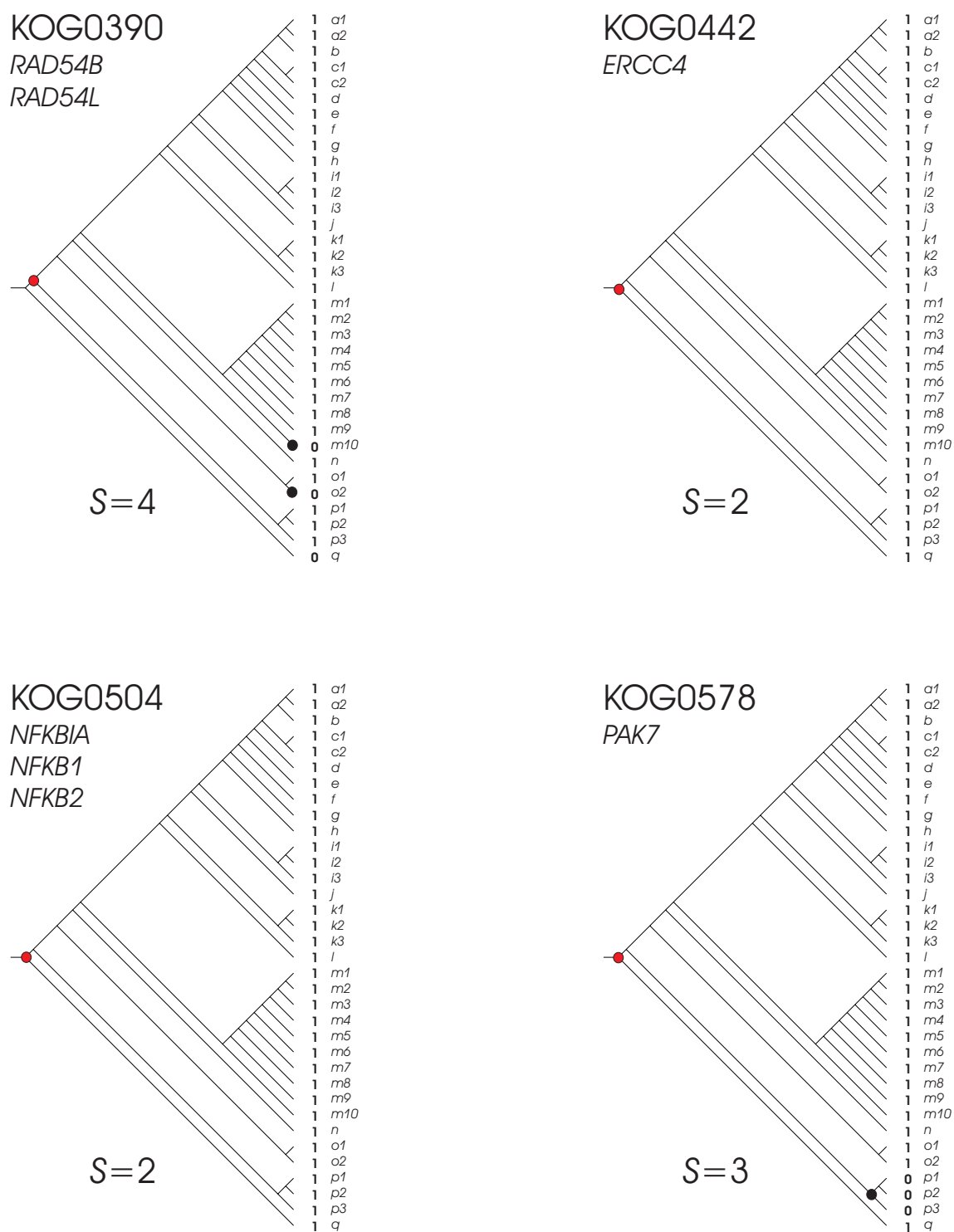
**Supplementary Figure S15.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0045, KOG0161, KOG0198 and KOG0217, as in Supplementary Figure S14.



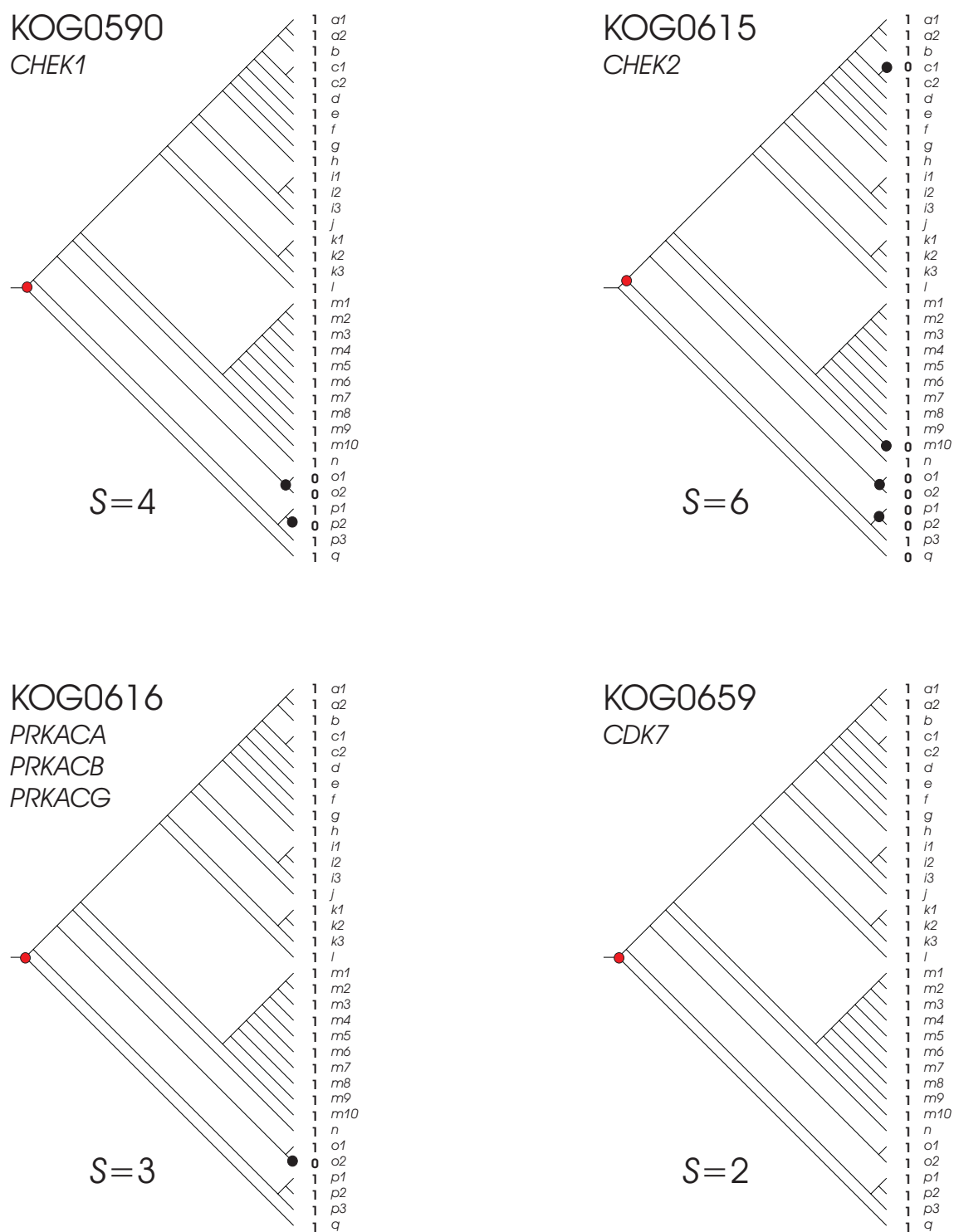
**Supplementary Figure S16.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0218, KOG0219, KOG0220 and KOG0221, as in Supplementary Figure S14.



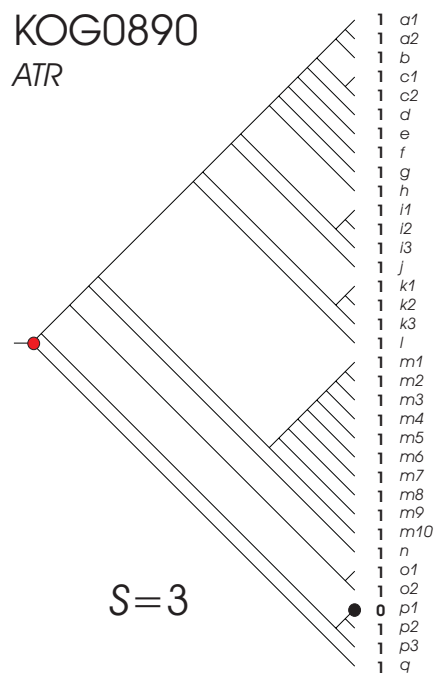
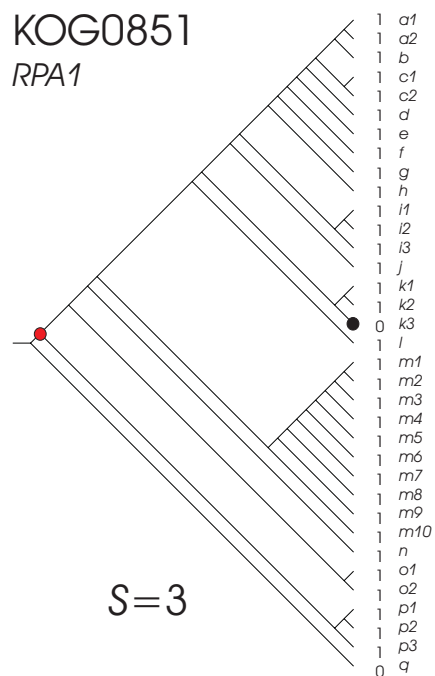
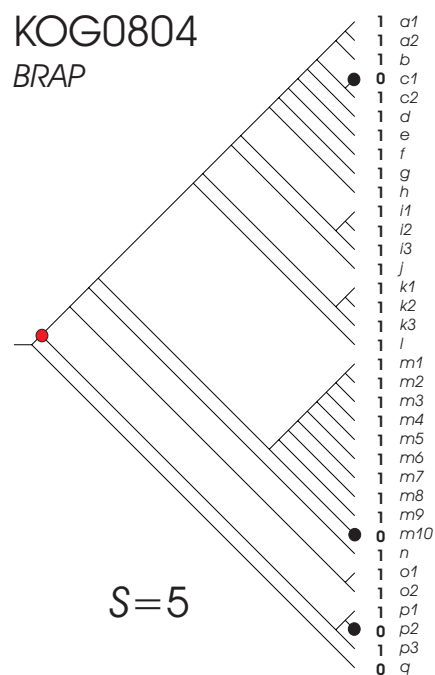
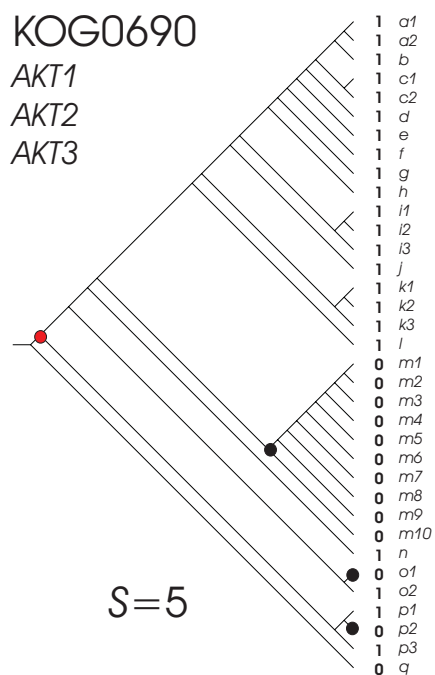
**Supplementary Figure S17.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0297, KOG0375, KOG0381 and KOG0387, as in Supplementary Figure S14.



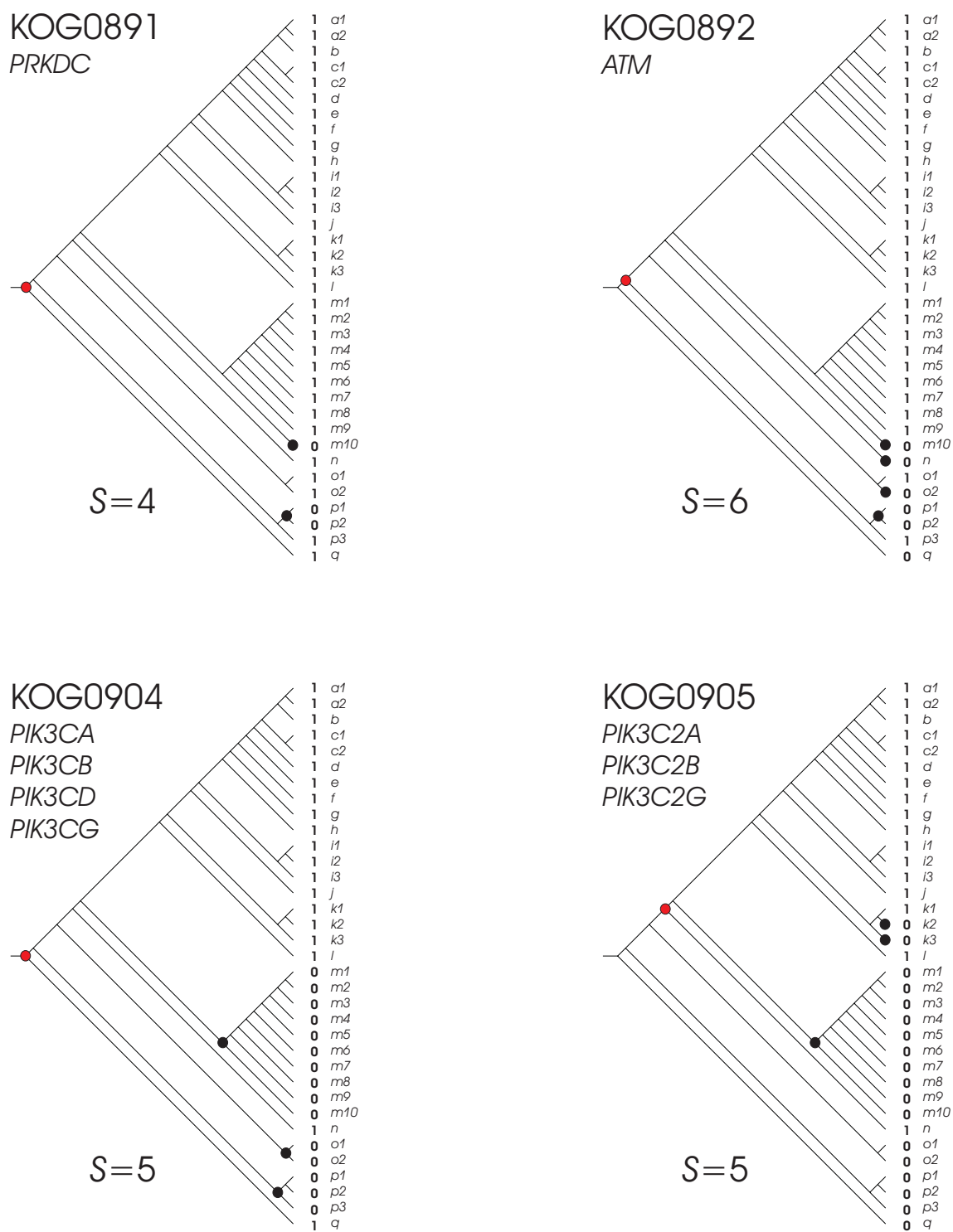
**Supplementary Figure S18.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0390, KOG0442, KOG0504 and KOG0578, as in Supplementary Figure S14.



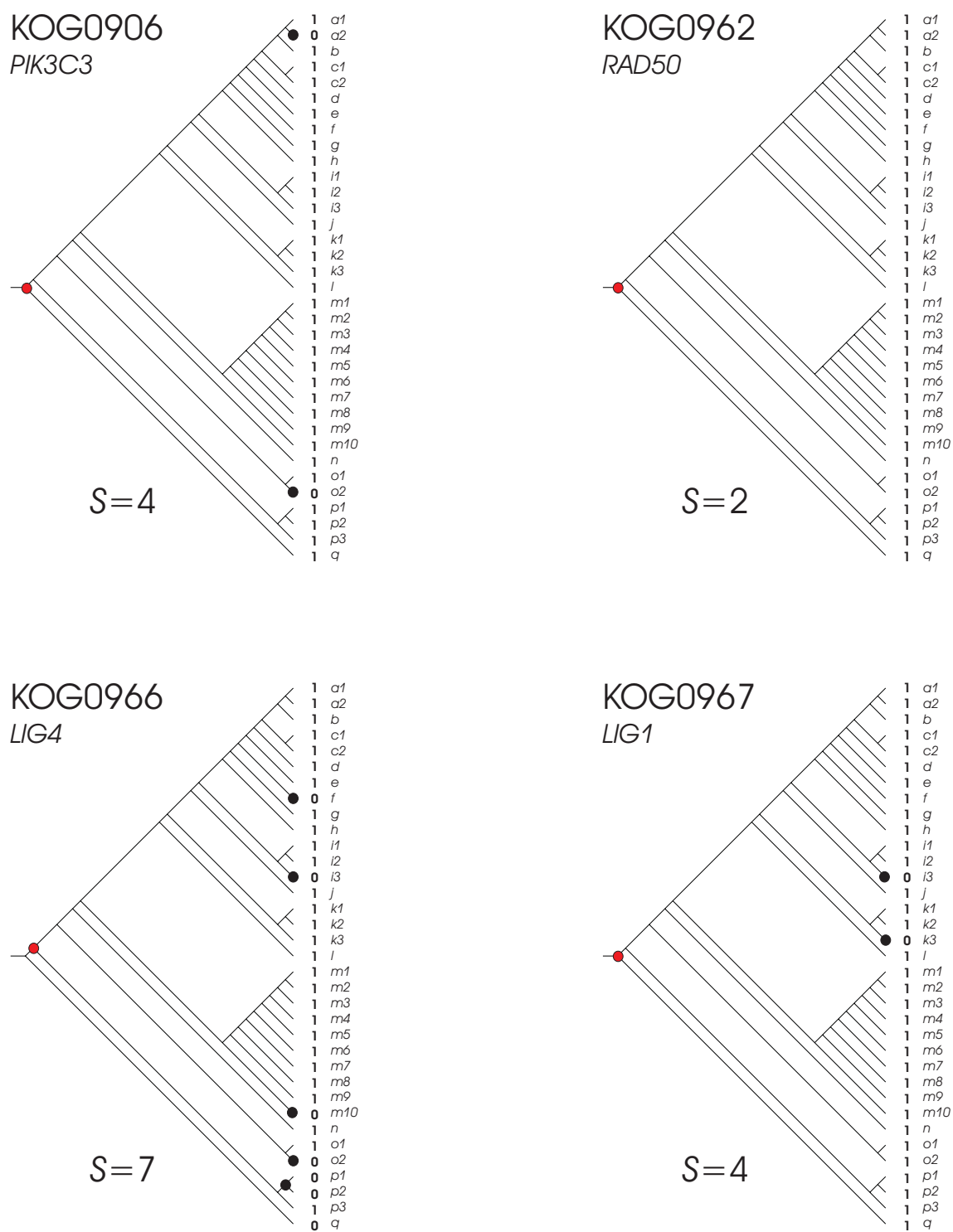
**Supplementary Figure S19.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0590, KOG0615, KOG0616 and KOG0659, as in Supplementary Figure S14.



**Supplementary Figure S20.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0690, KOG0804, KOG0851 and KOG0890, as in Supplementary Figure S14.

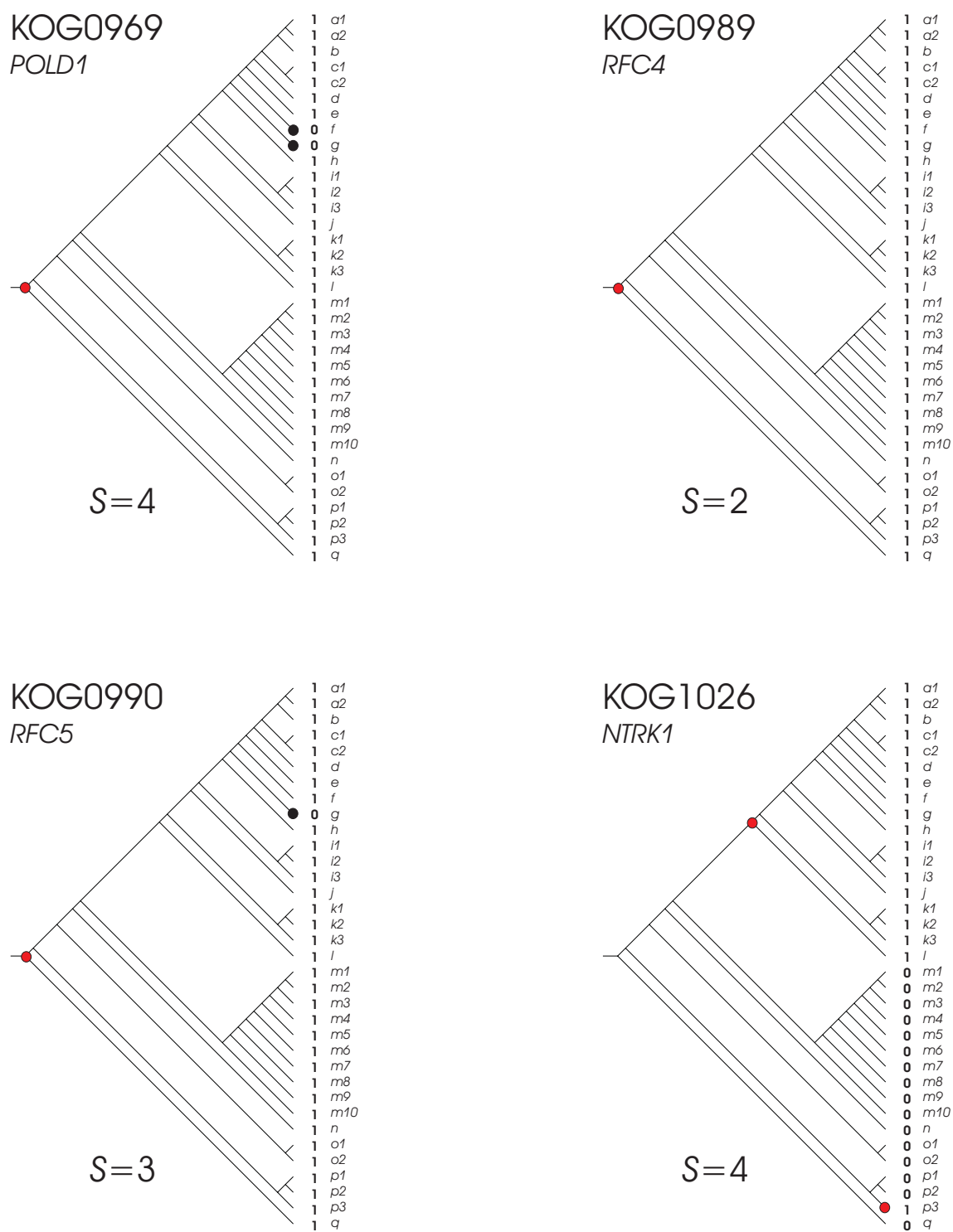


**Supplementary Figure S21.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0891, KOG0892, KOG0904 and KOG0905, as in Supplementary Figure S14.

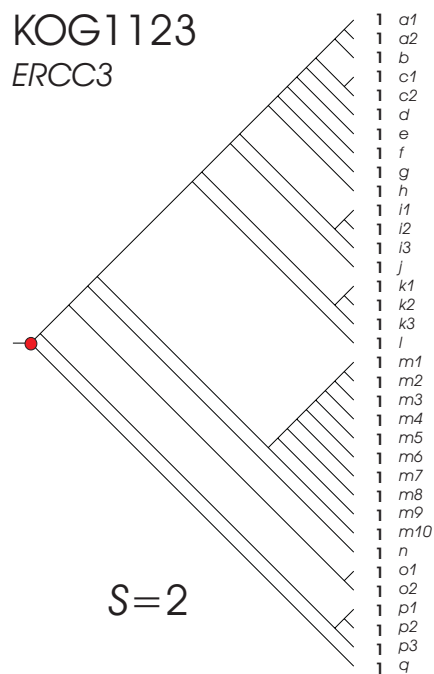
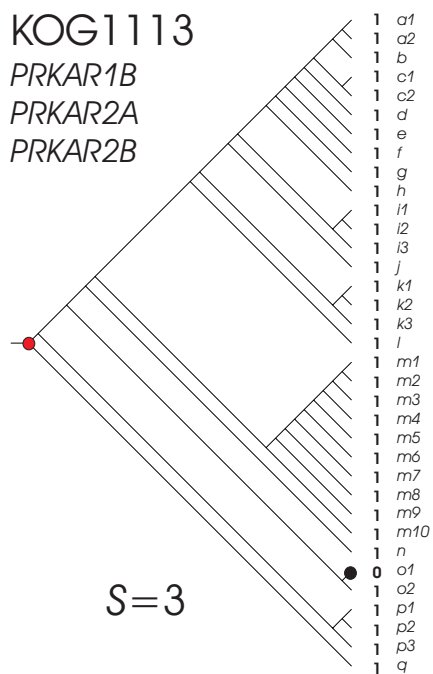
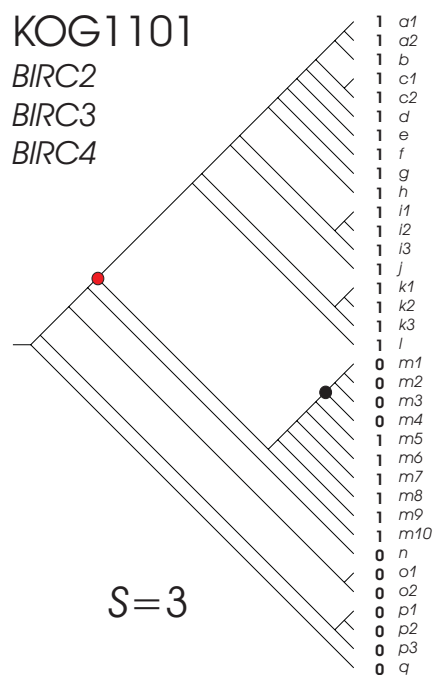
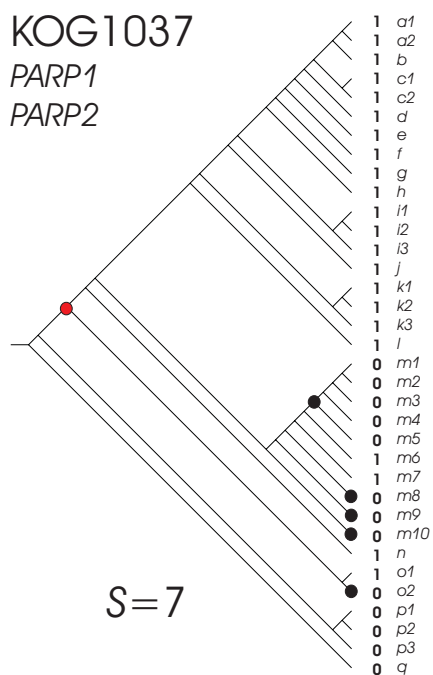


**Supplementary Figure S22.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0906, KOG0962, KOG0966 and KOG0967, as in Supplementary Figure S14.

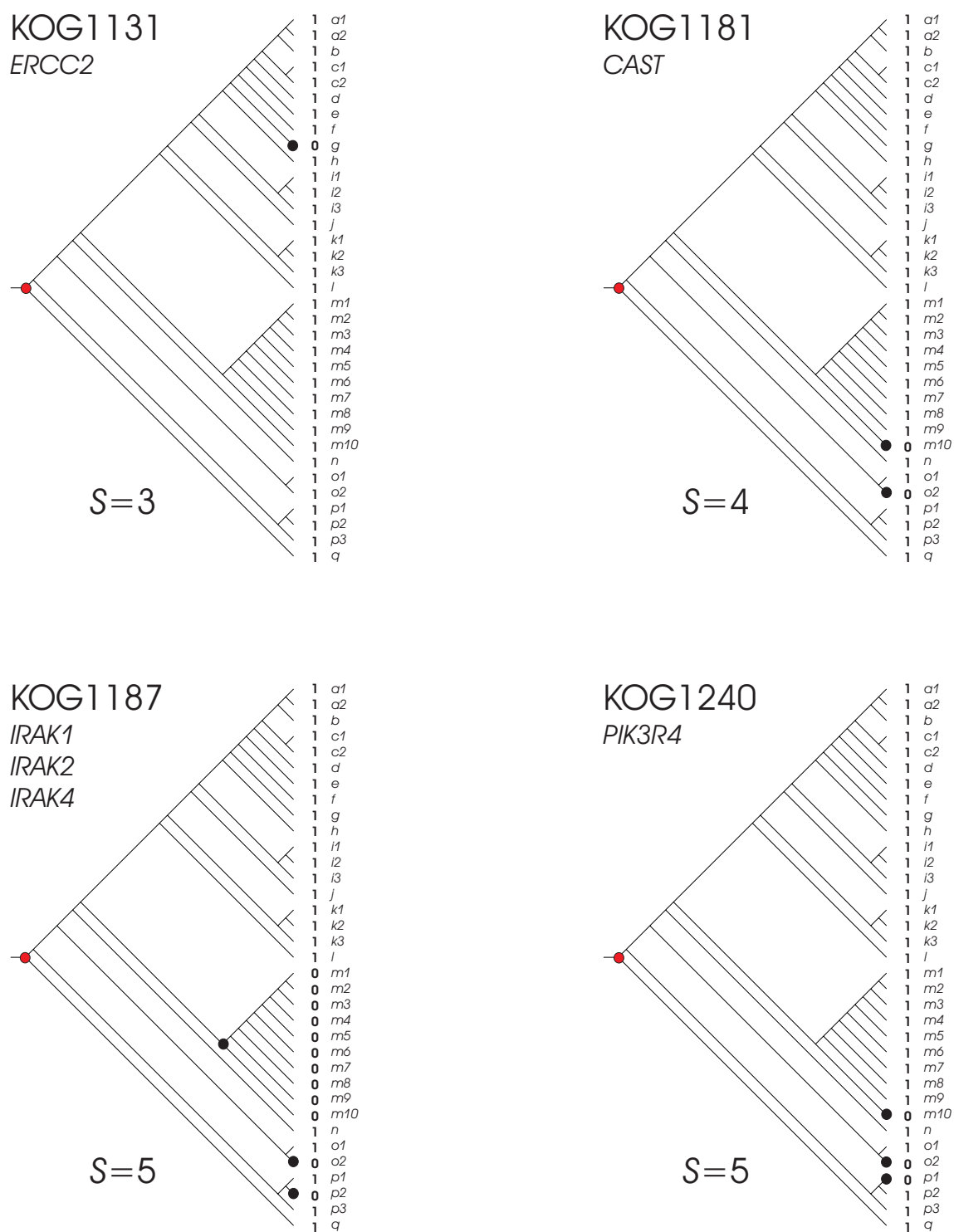




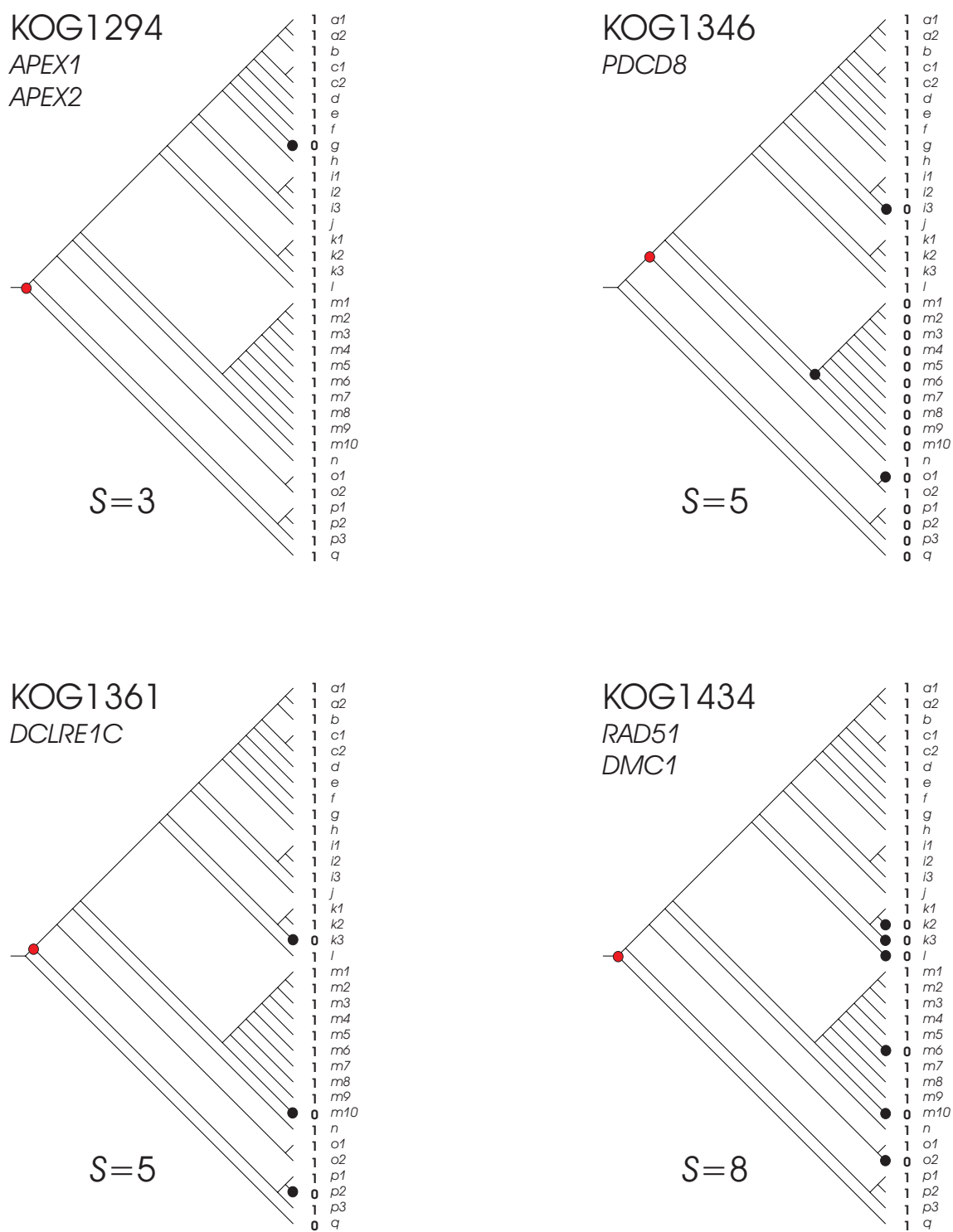
**Supplementary Figure S23.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG0969, KOG0989, KOG0990 and KOG1026, as in Supplementary Figure S14.



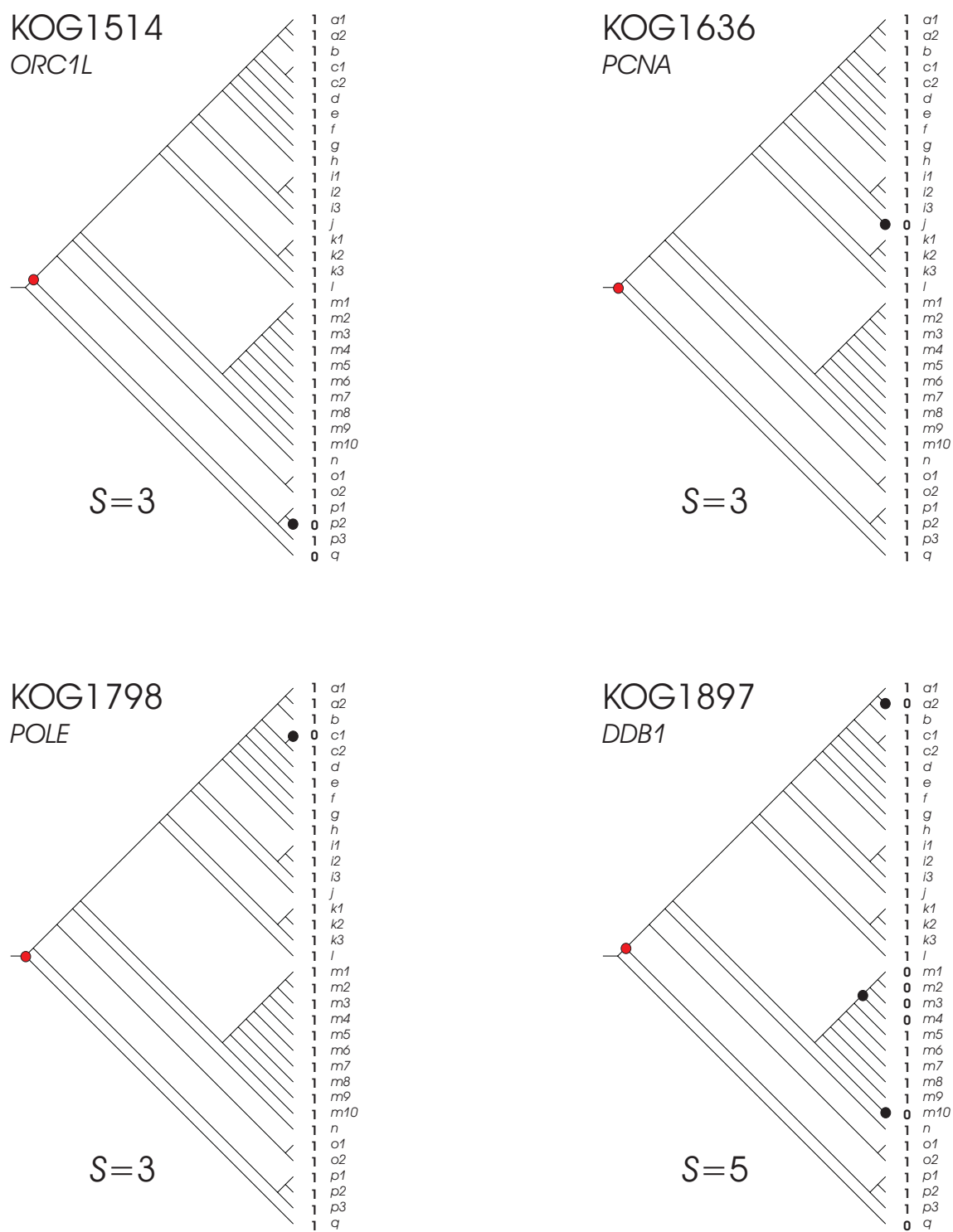
**Supplementary Figure S24.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG1037, KOG1101, KOG1113 and KOG1123, as in Supplementary Figure S14.



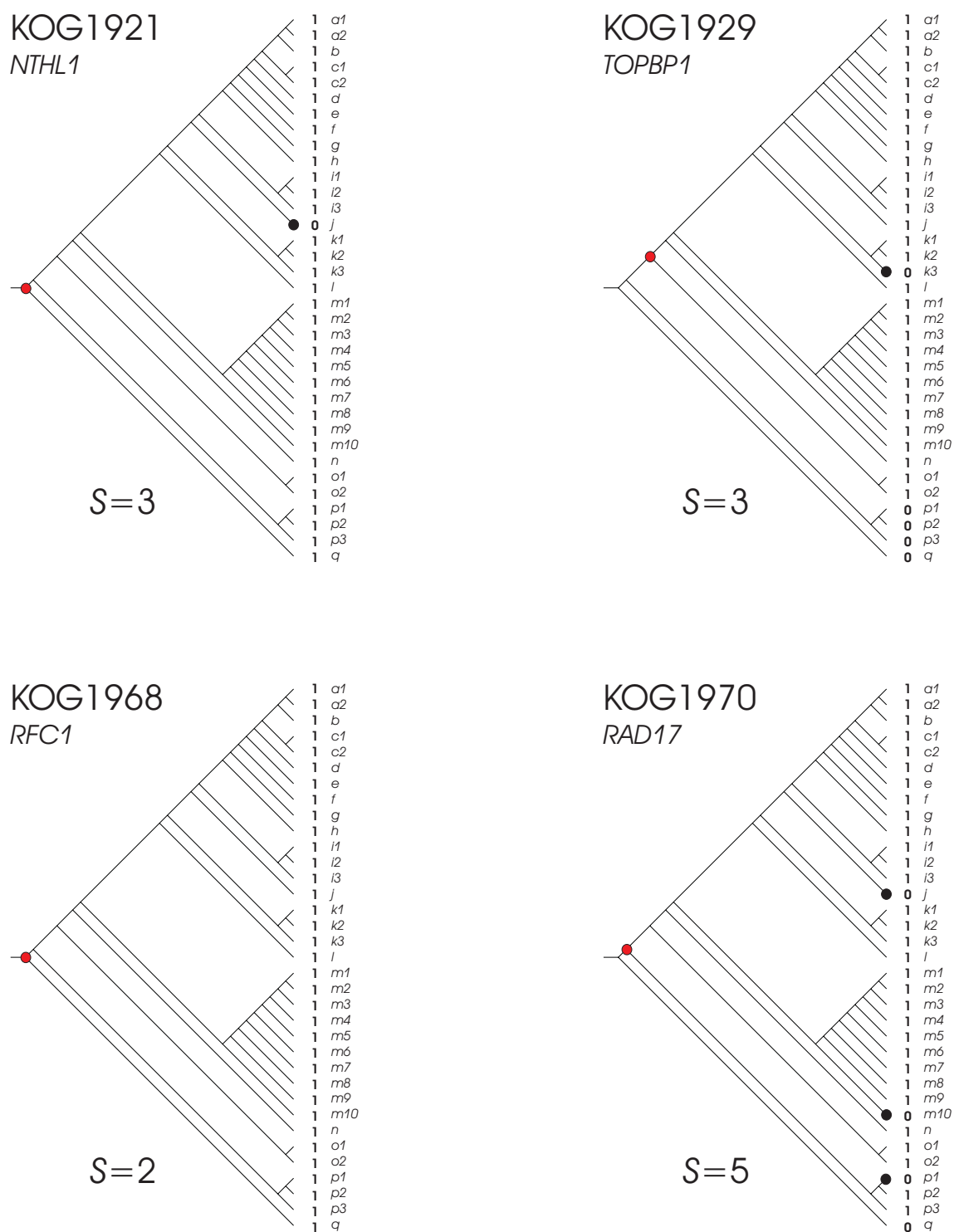
**Supplementary Figure S25.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG1131, KOG1181, KOG1187 and KOG1240, as in Supplementary Figure S14.



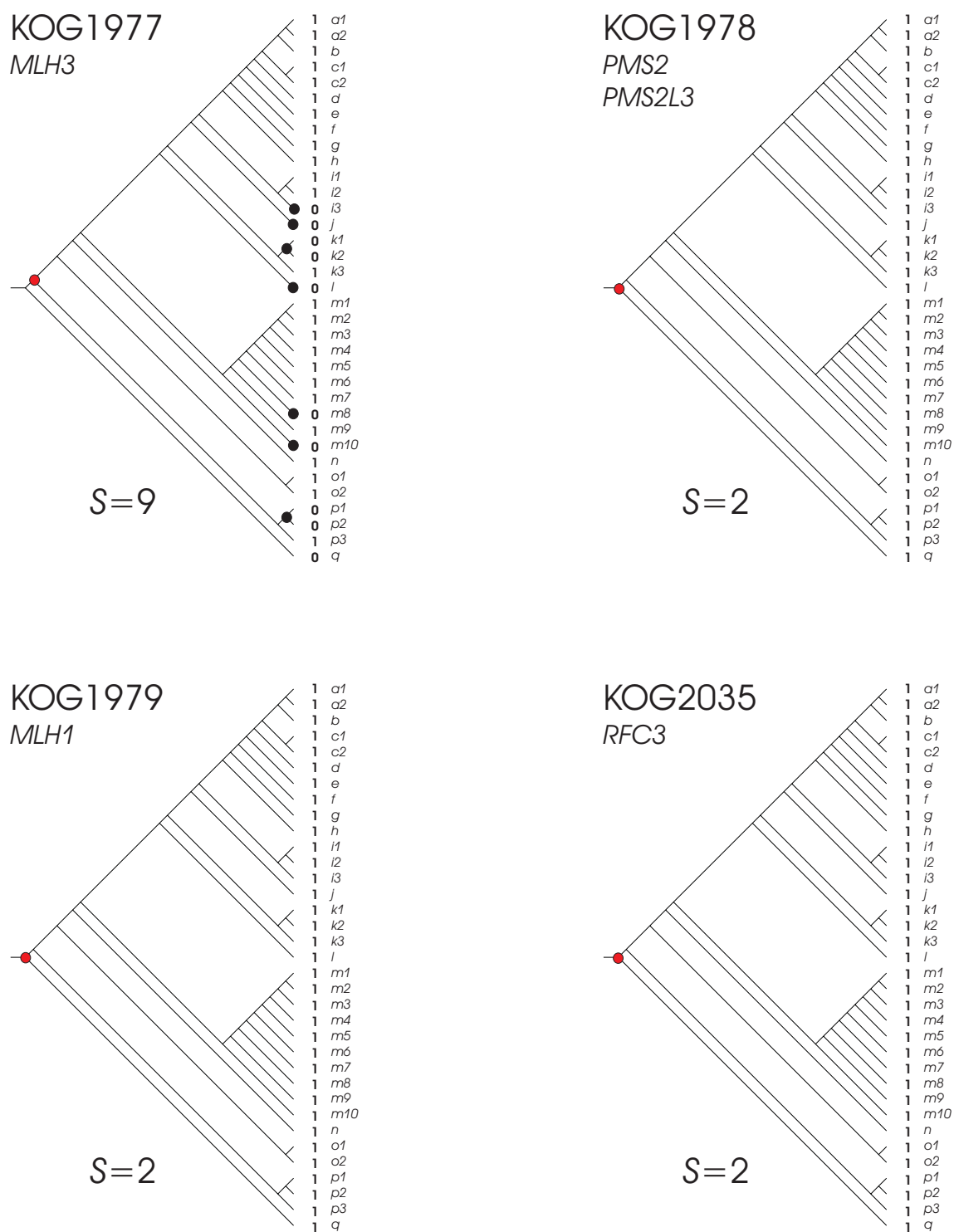
**Supplementary Figure S26.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG1294, KOG1346, KOG1361 and KOG1434, as in Supplementary Figure S14.



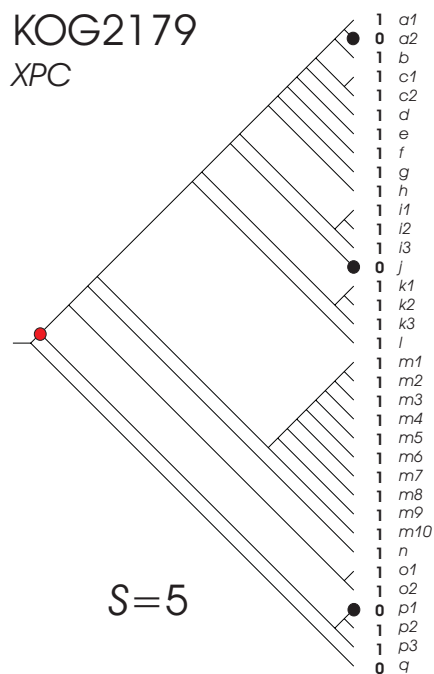
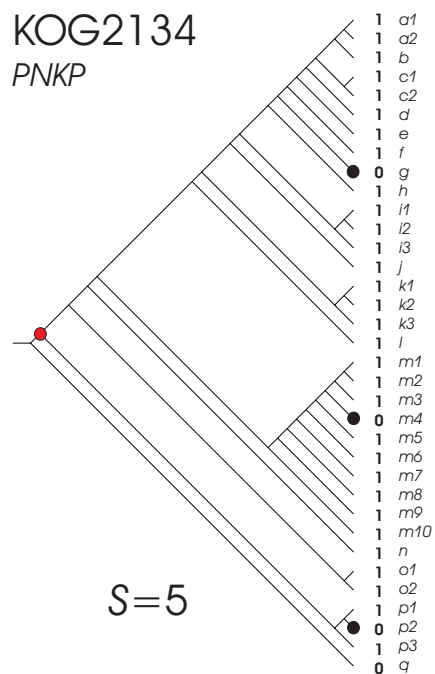
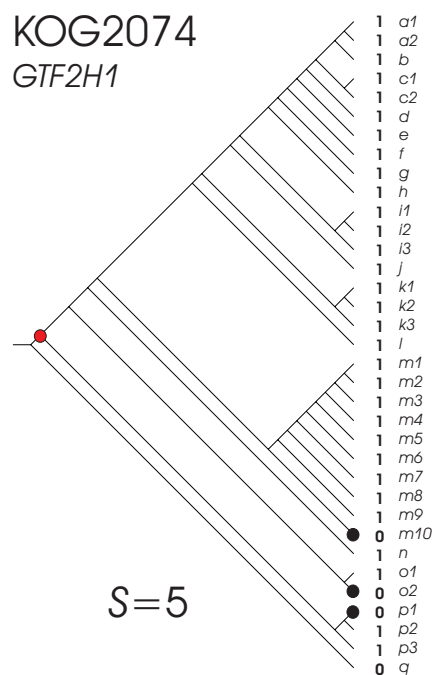
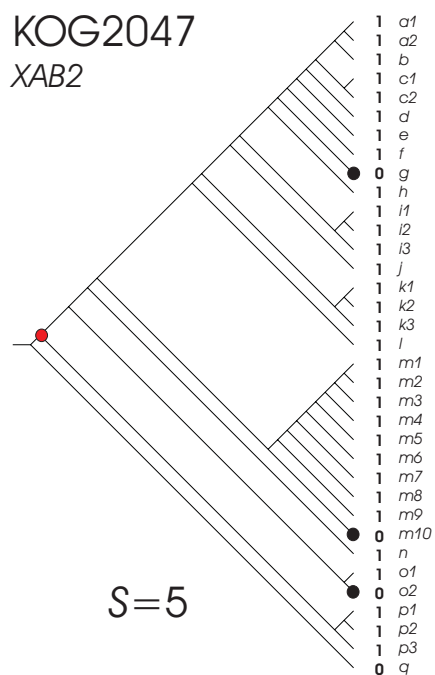
**Supplementary Figure S27.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG1514, KOG1636, KOG1798 and KOG1897, as in Supplementary Figure S14.



**Supplementary Figure S28.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG1921, KOG1929, KOG1968 and KOG1970, as in Supplementary Figure S14.

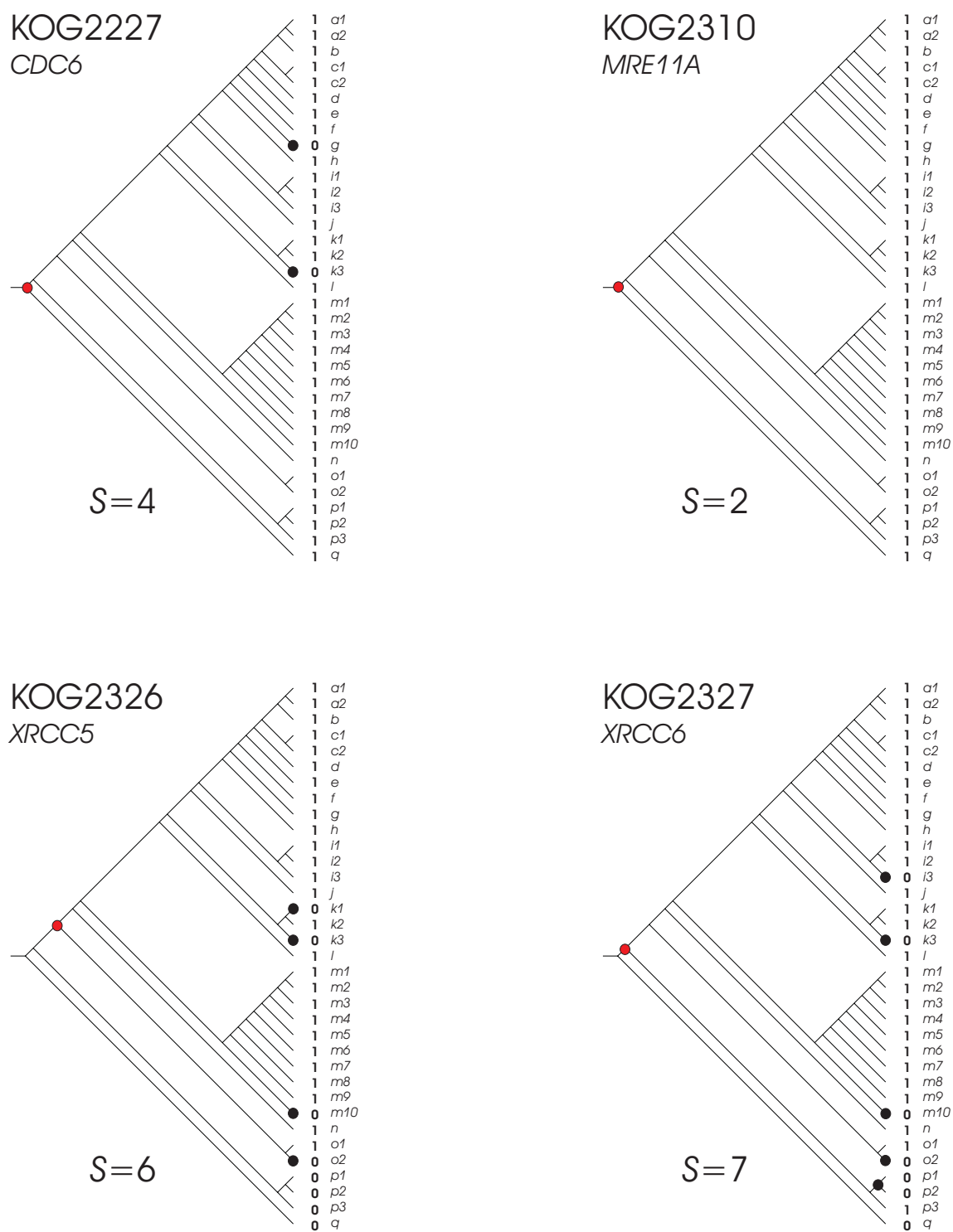


**Supplementary Figure S29.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG1977, KOG1978, KOG1979 and KOG2035, as in Supplementary Figure S14.

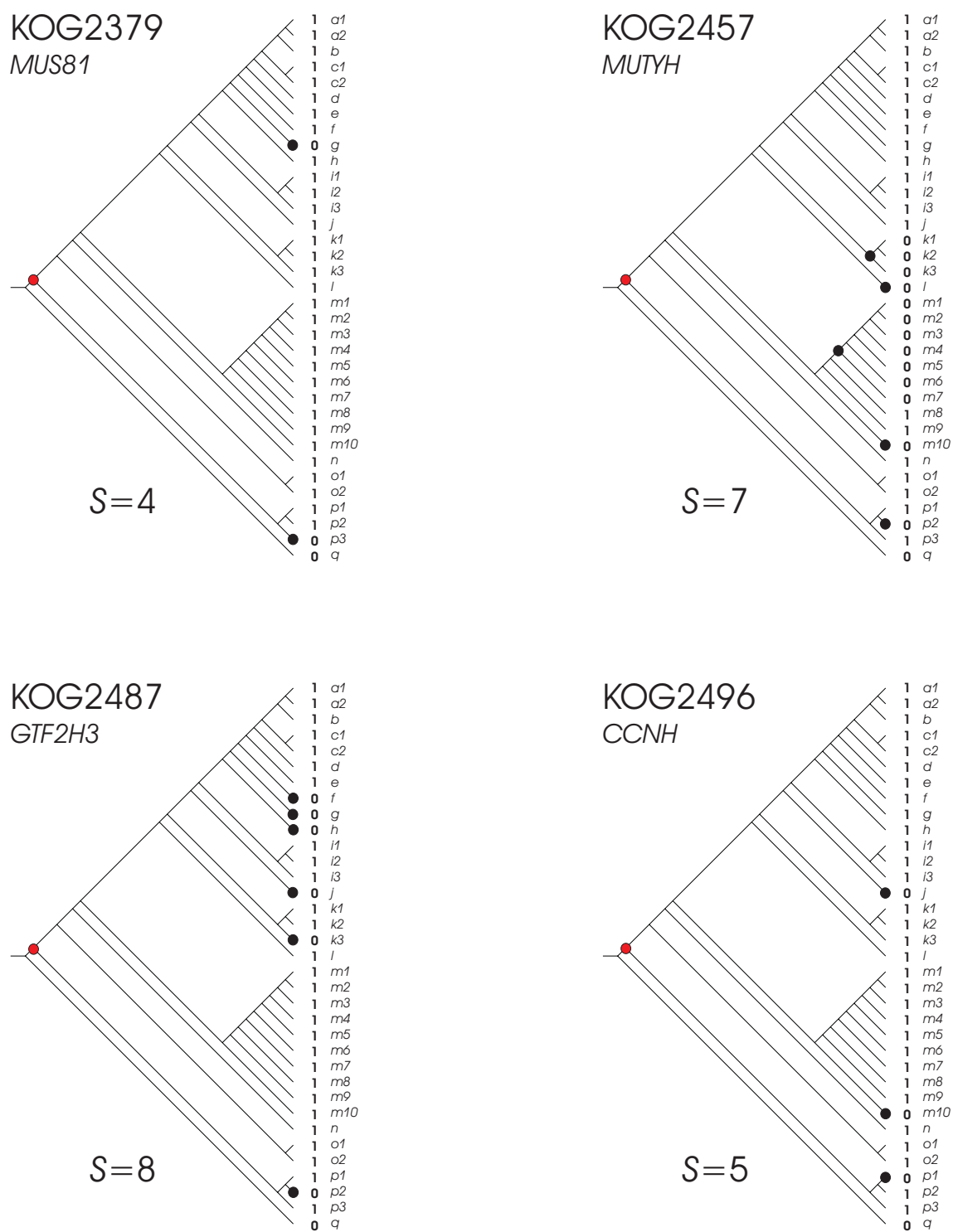


**Supplementary Figure S30.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG2047, KOG2074, KOG2134 and KOG2179, as in Supplementary Figure S14.

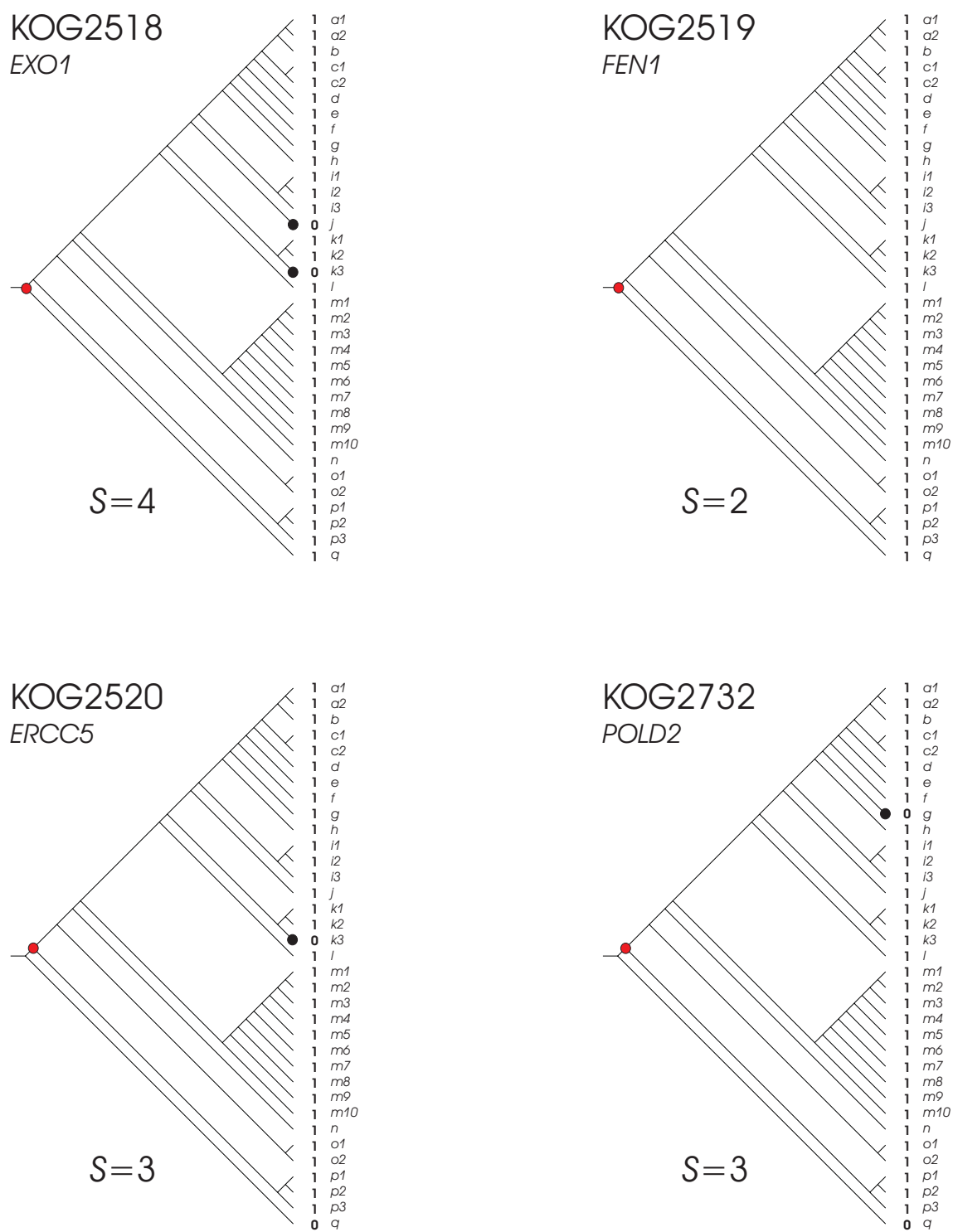




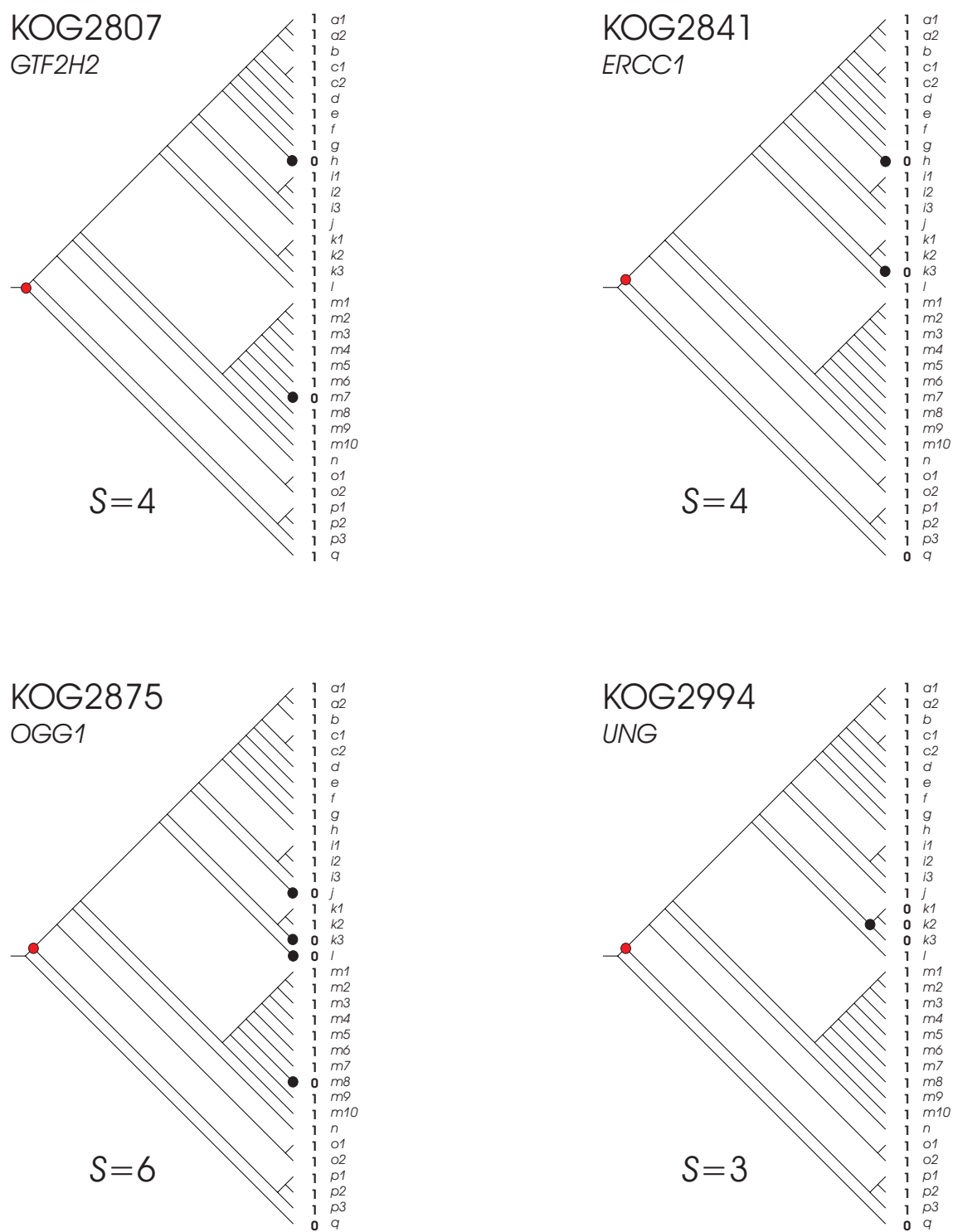
**Supplementary Figure S31.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG2227, KOG2310, KOG2326 and KOG2327, as in Supplementary Figure S14.



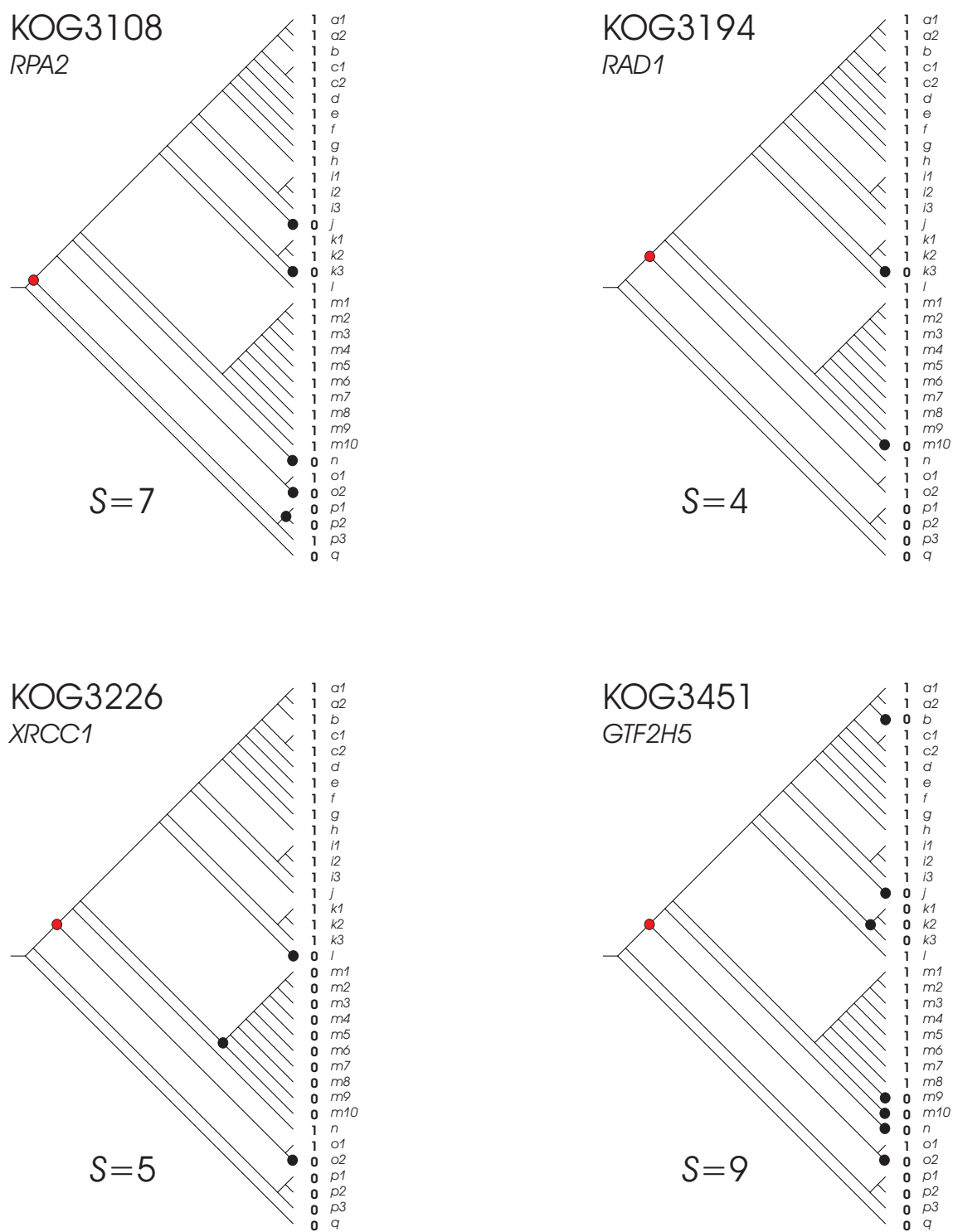
**Supplementary Figure S32.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG2379, KOG2457, KOG2487 and KOG2496, as in Supplementary Figure S14.



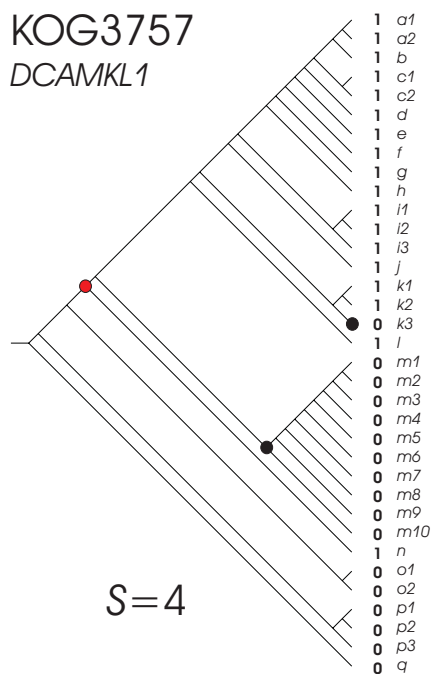
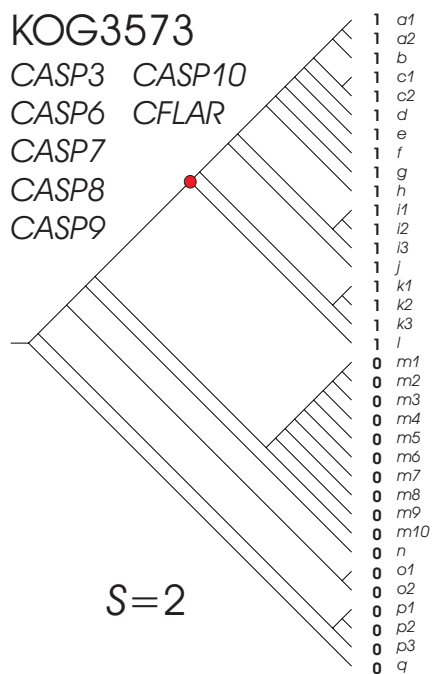
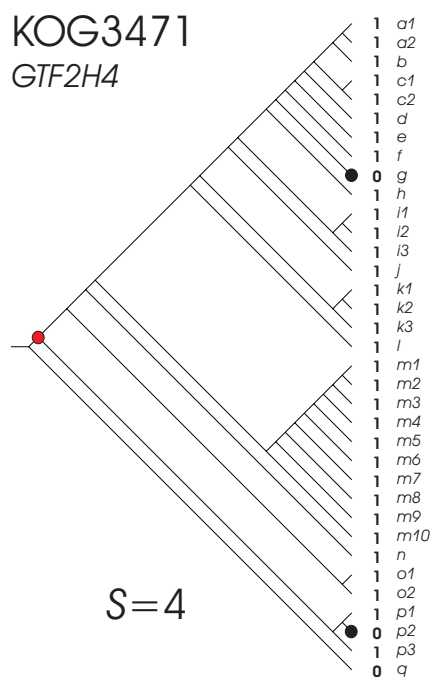
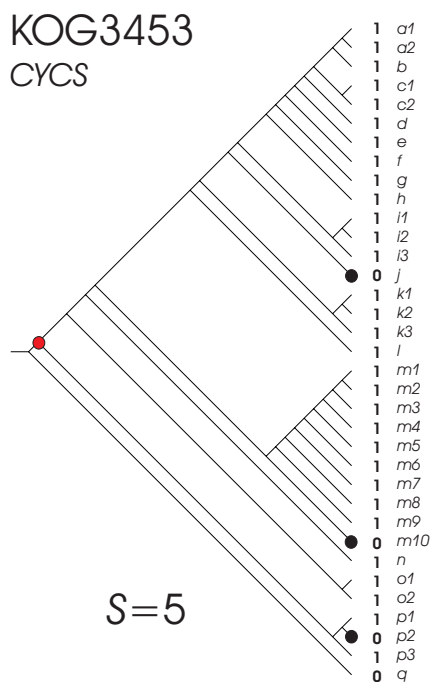
**Supplementary Figure S33.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG2518, KOG2519, KOG2520 and KOG2732, as in Supplementary Figure S14.



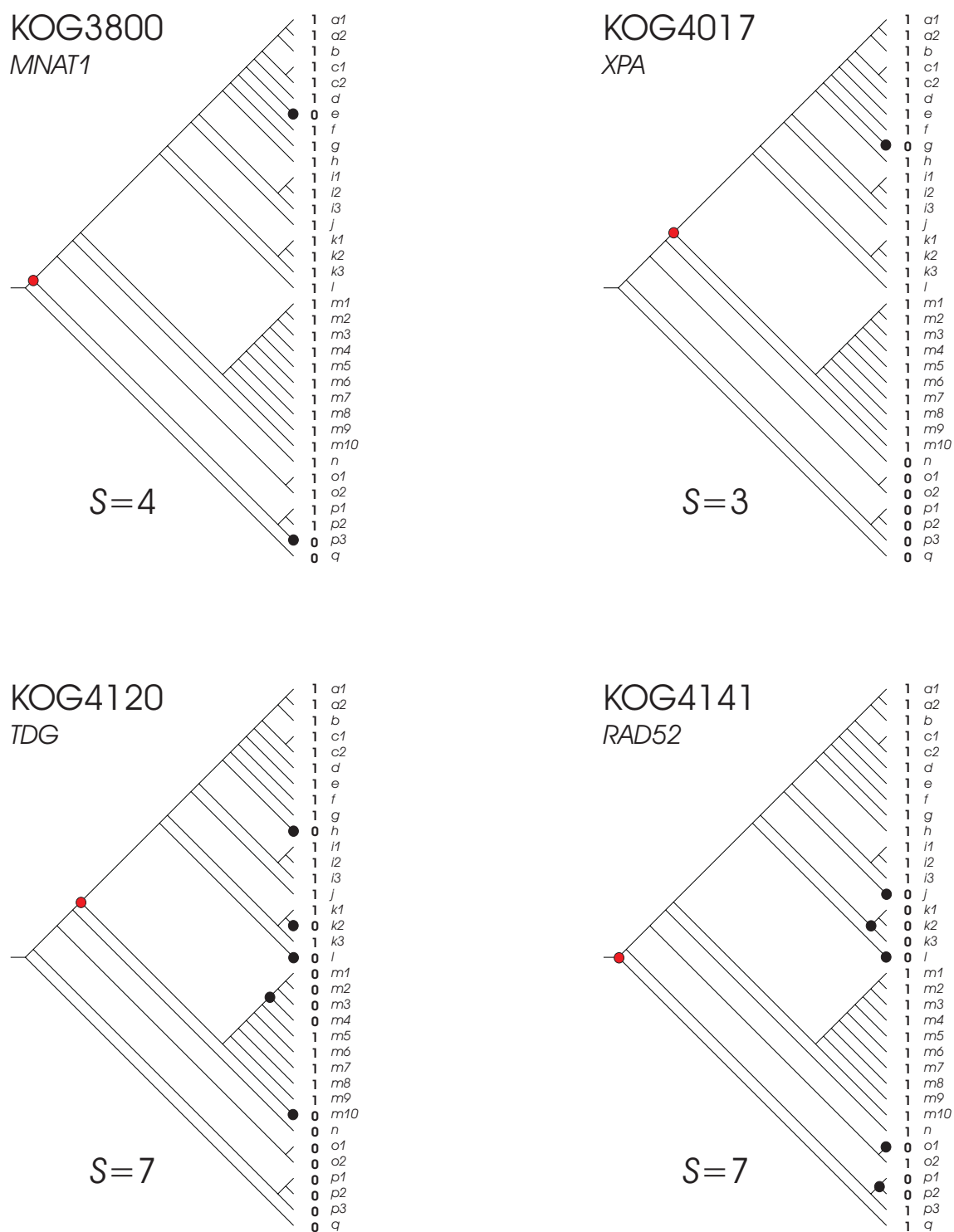
**Supplementary Figure S34.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG2807, KOG2841, KOG2875 and KOG2994, as in Supplementary Figure S14.



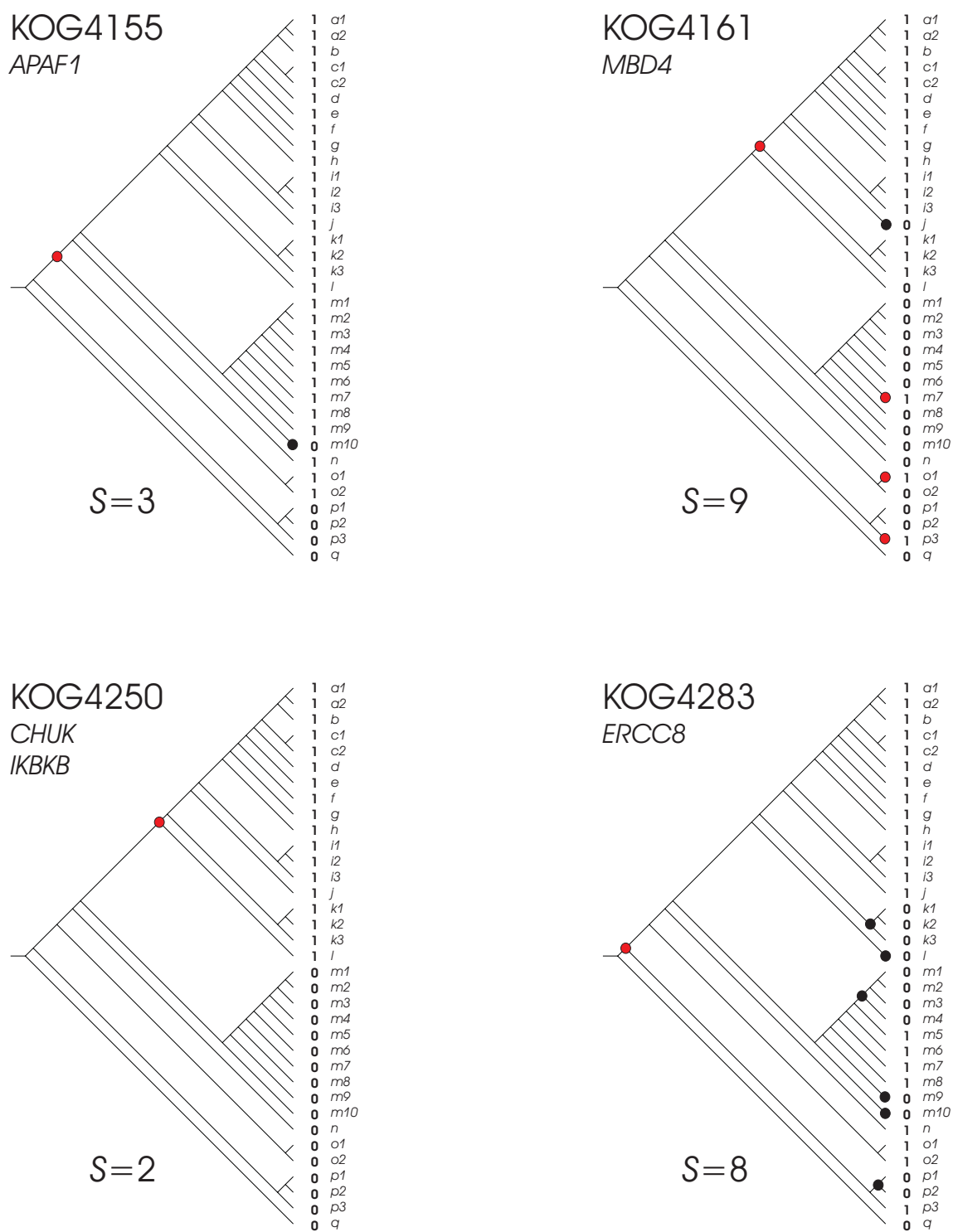
**Supplementary Figure S35.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG3108, KOG3194, KOG3226 and KOG3451, as in Supplementary Figure S14.



**Supplementary Figure S36.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG3453, KOG3471, KOG3573 and KOG3757, as in Supplementary Figure S14.

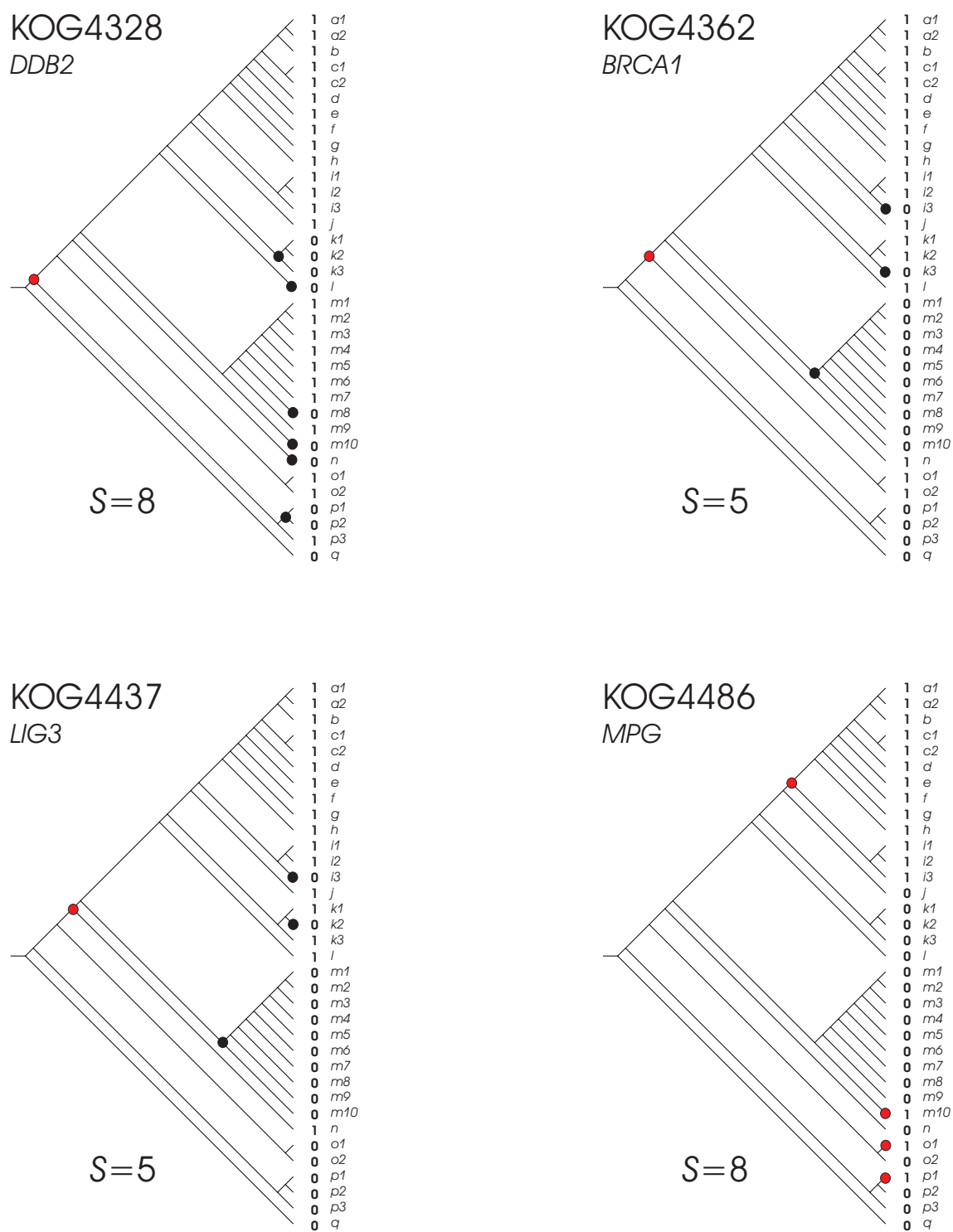


**Supplementary Figure S37.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG3800, KOG4017, KOG4120 and KOG4141, as in Supplementary Figure S14.

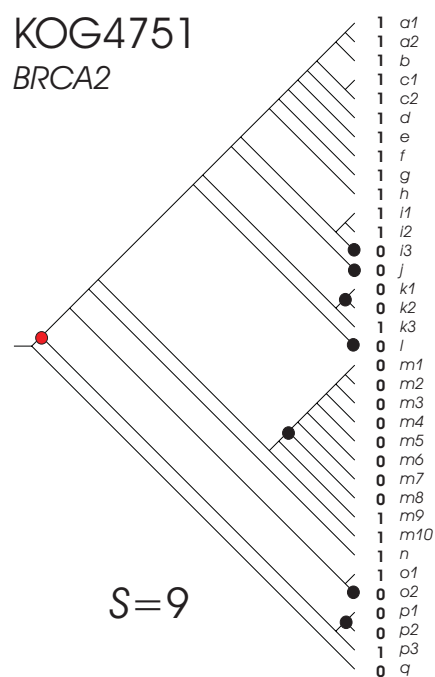
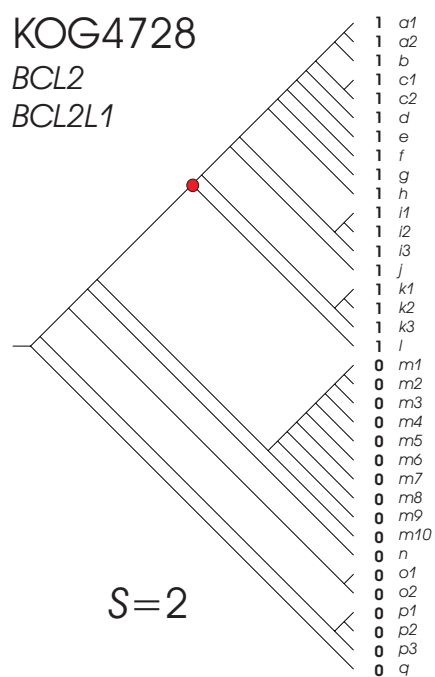
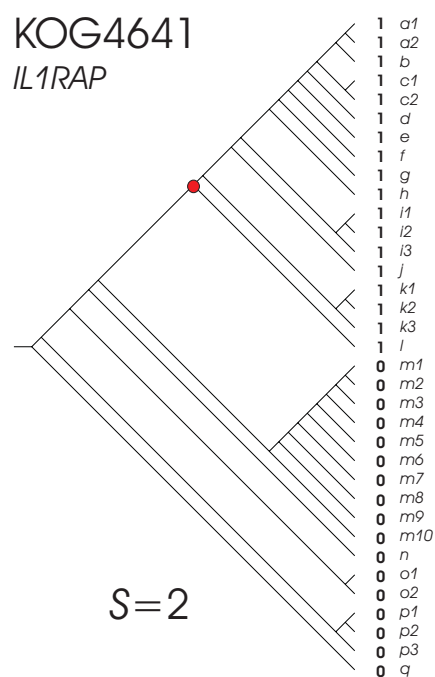
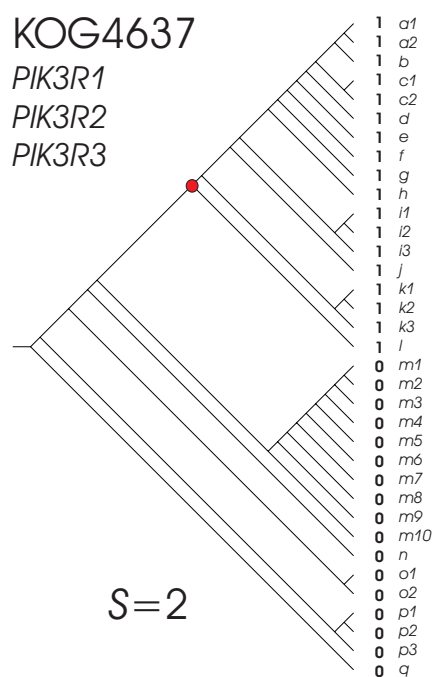


**Supplementary Figure S38.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG4155, KOG4161, KOG4250 and KOG4283, as in Supplementary Figure S14.

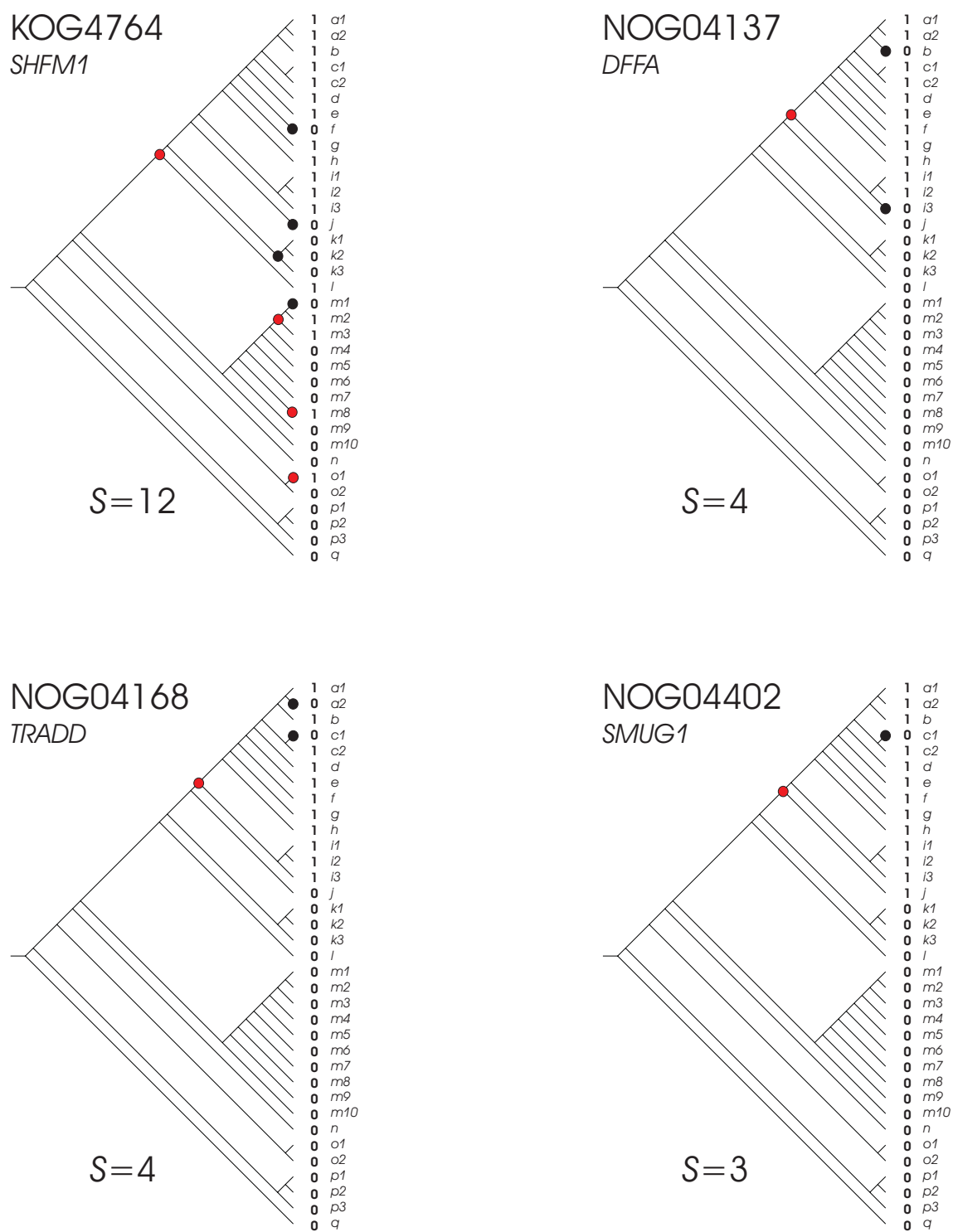




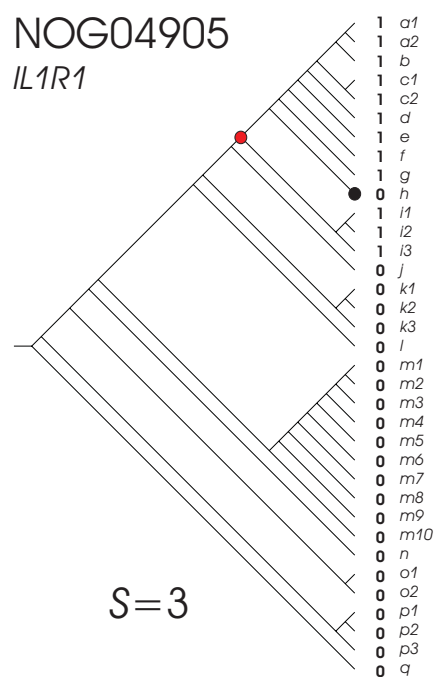
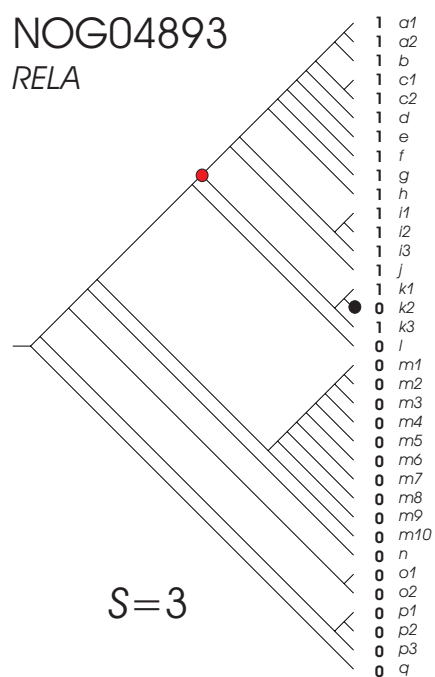
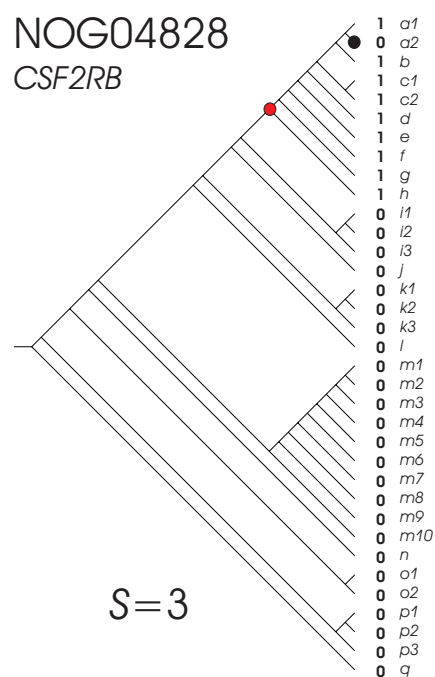
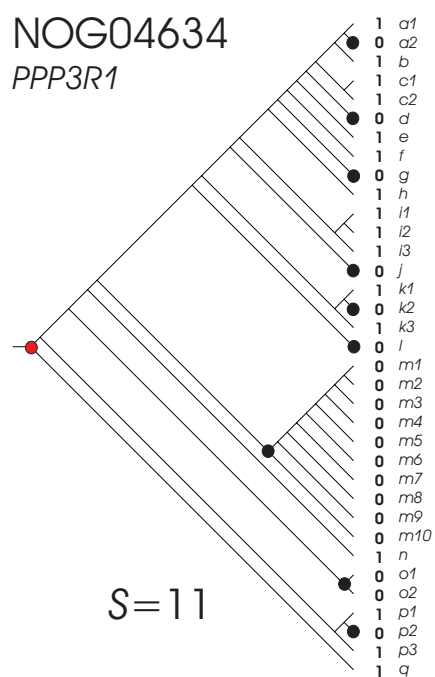
**Supplementary Figure S39.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG4328, KOG4362, KOG4437 and KOG4486, as in Supplementary Figure S14.



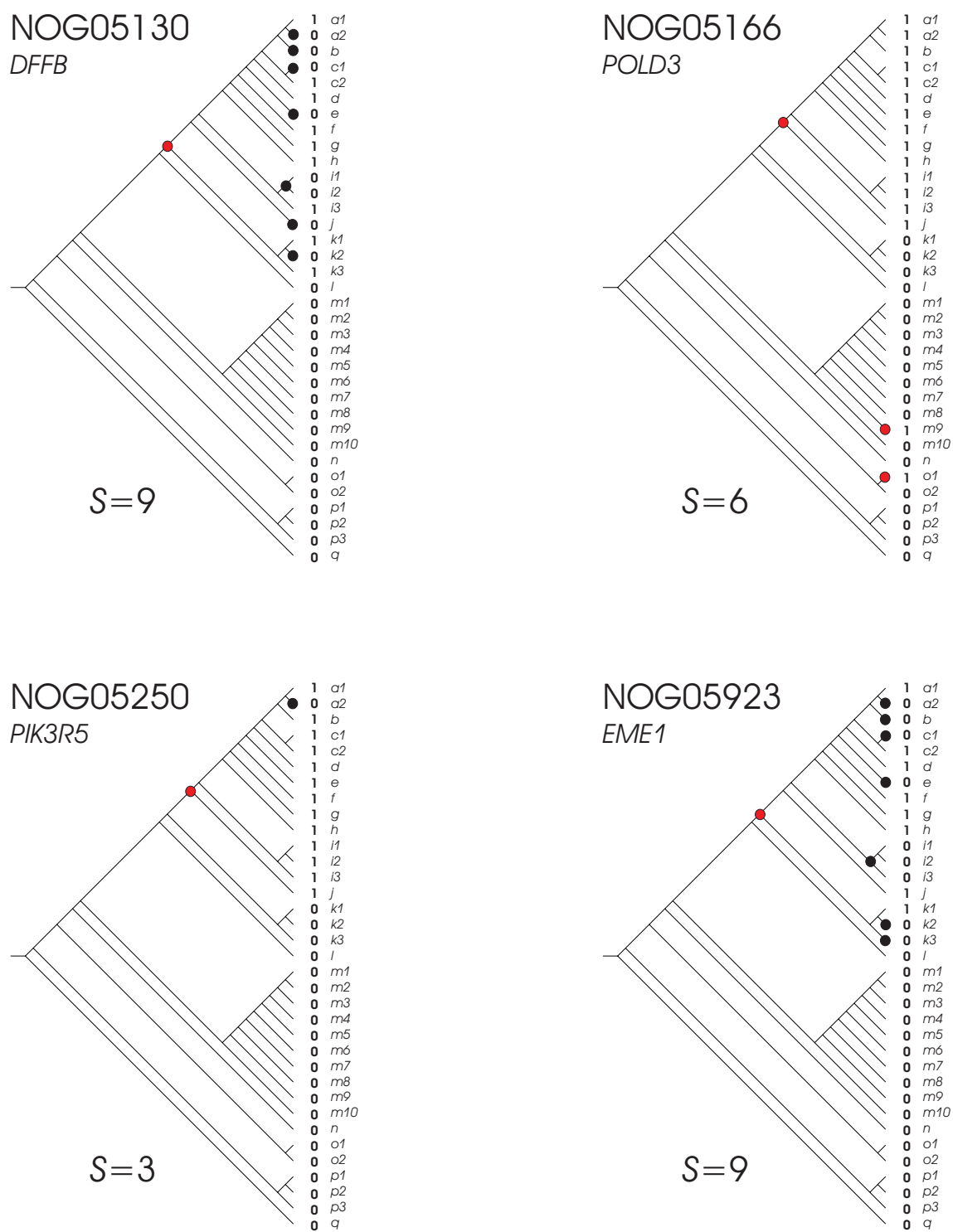
**Supplementary Figure S40.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG4637, KOG4641, KOG4728 and KOG4751, as in Supplementary Figure S14.



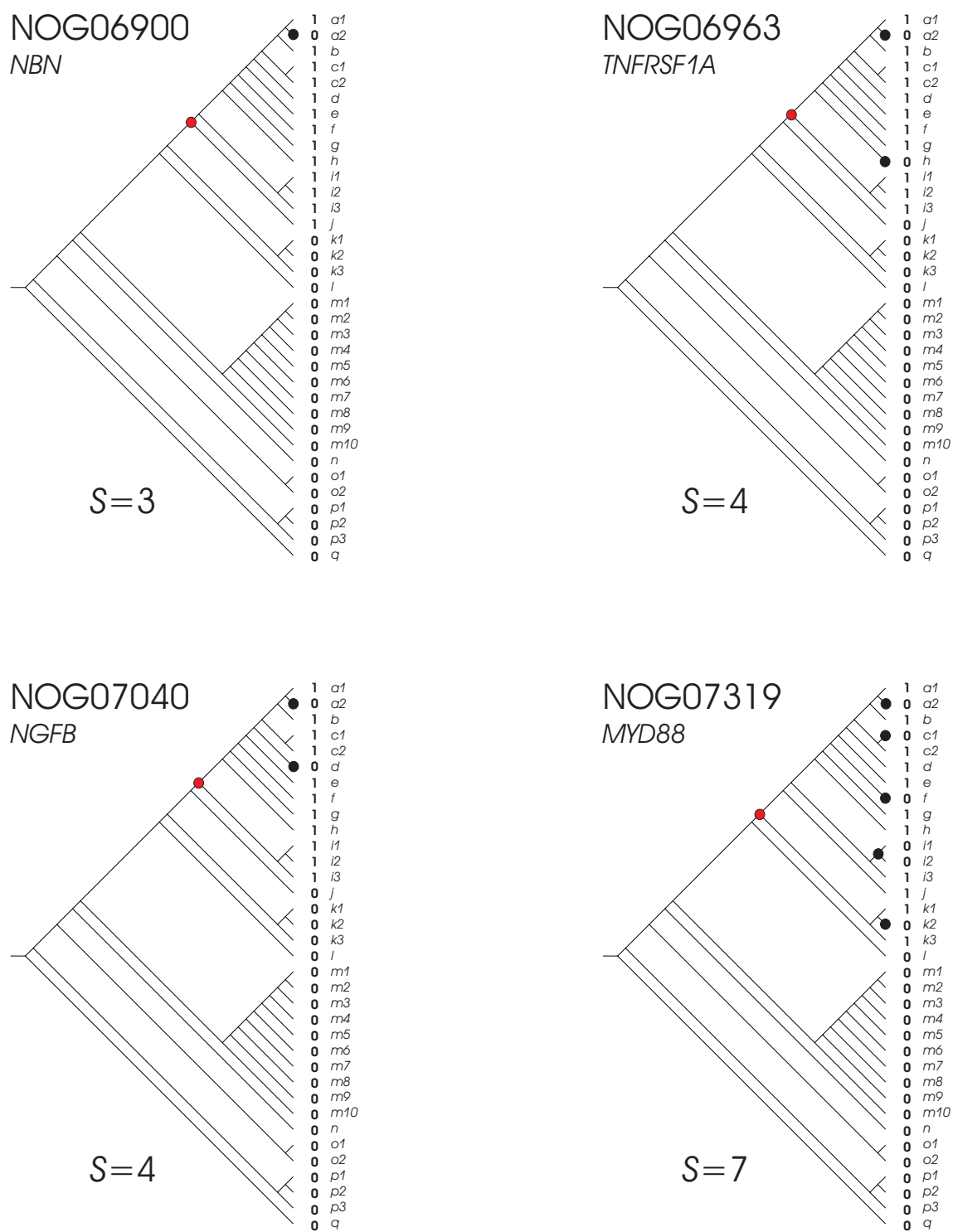
**Supplementary Figure S41.** Parsimony analysis of eukaryotic clusters of orthologous groups: KOG4764, NOG04137, NOG04168 and NOG04402, as in Supplementary Figure S14.



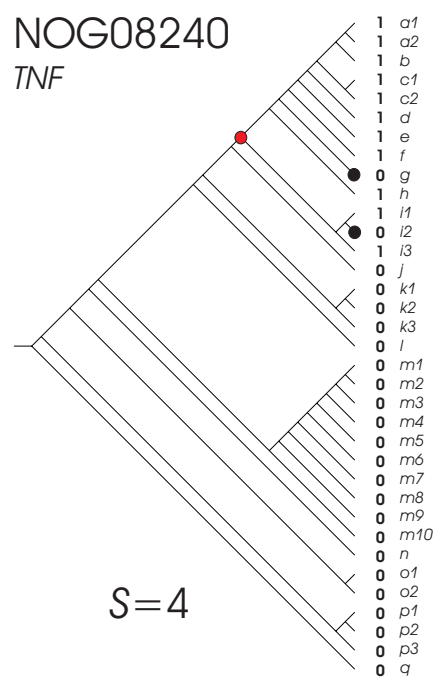
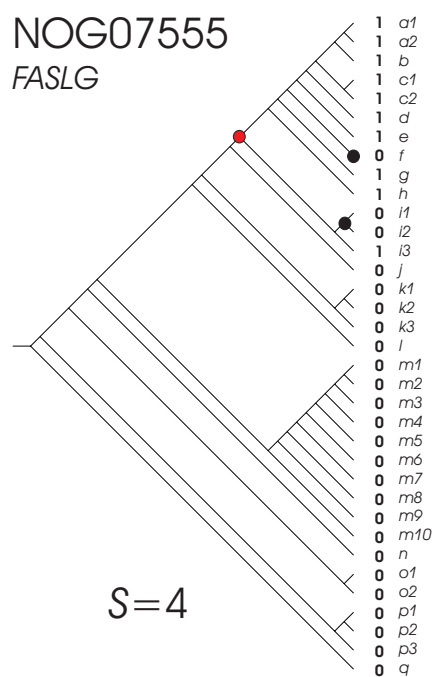
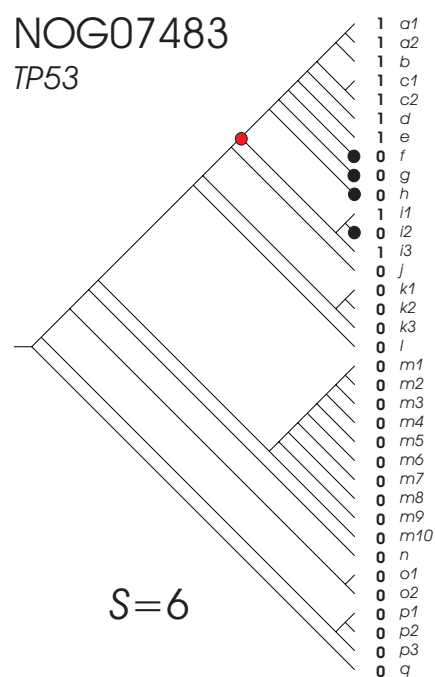
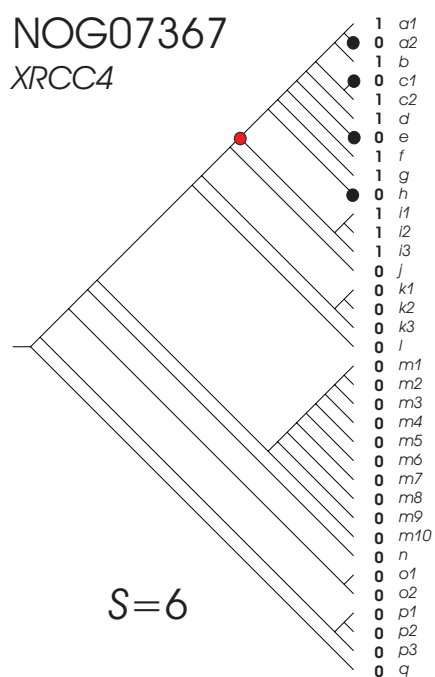
**Supplementary Figure S42.** Parsimony analysis of eukaryotic clusters of orthologous groups: NOG04634, NOG04828, NOG04893 and NOG04905, as in Supplementary Figure S14.



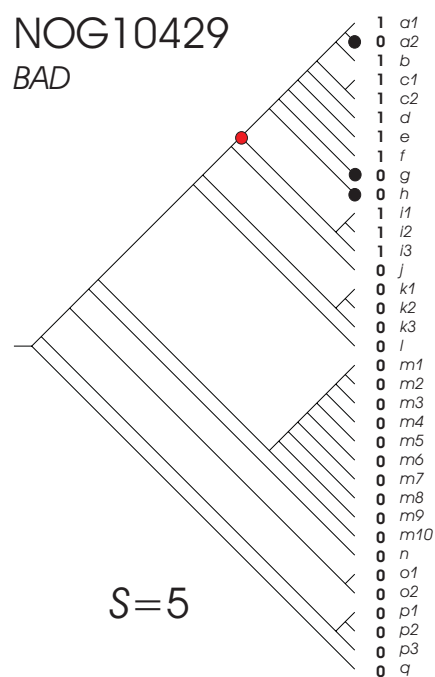
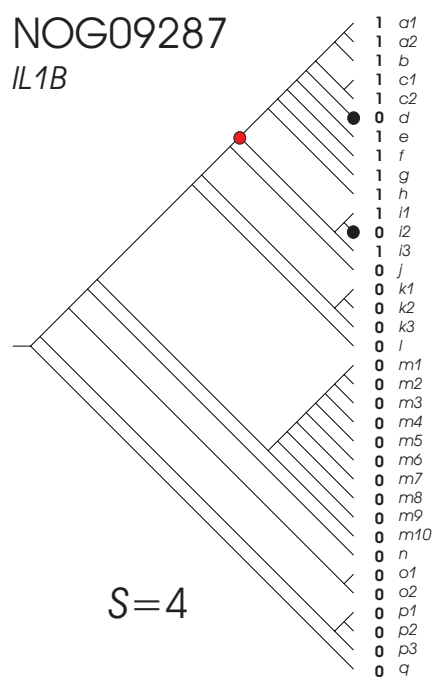
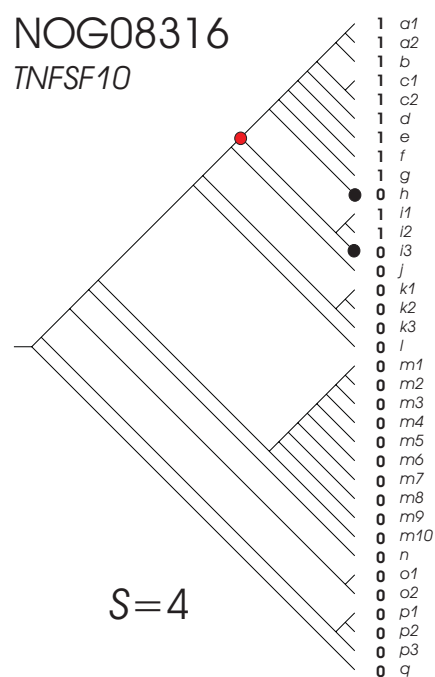
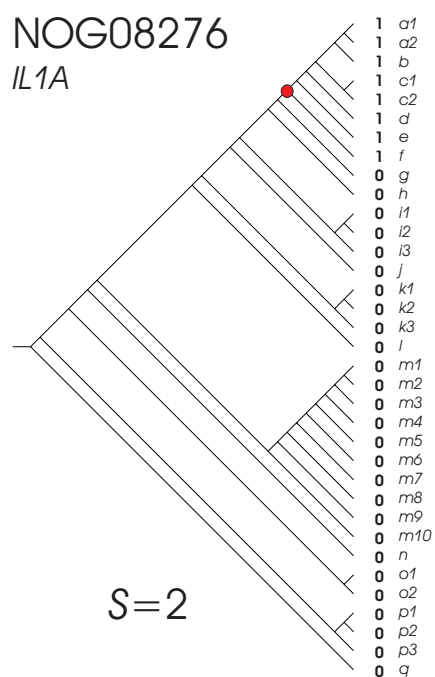
**Supplementary Figure S43.** Parsimony analysis of eukaryotic clusters of orthologous groups: NOG05130, NOG05166, NOG05250 and NOG05923, as in Supplementary Figure S14.



**Supplementary Figure S44.** Parsimony analysis of eukaryotic clusters of orthologous groups: NOG06900, NOG06963, NOG07040 and NOG07319, as in Supplementary Figure S14.

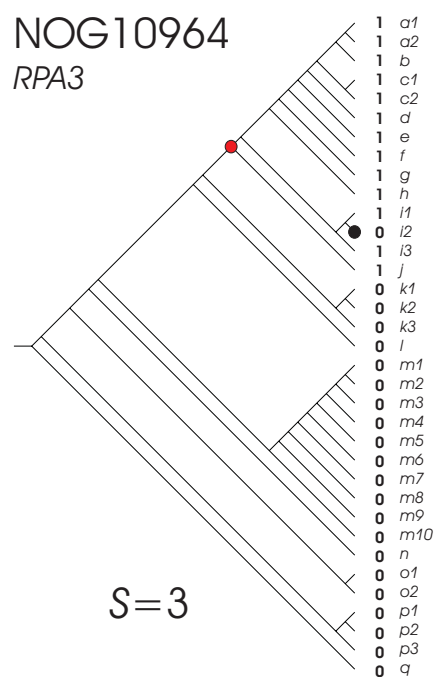
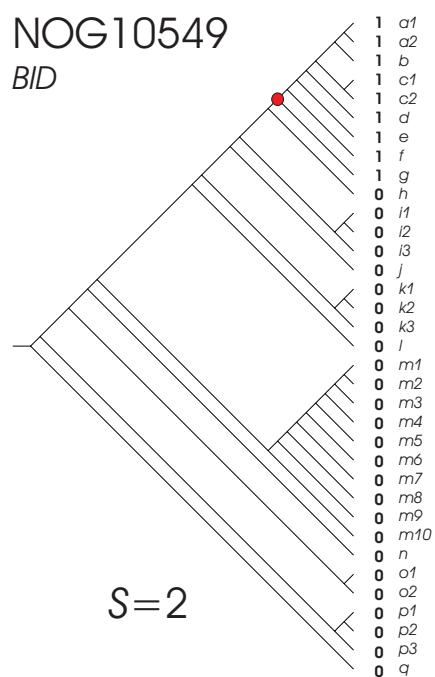
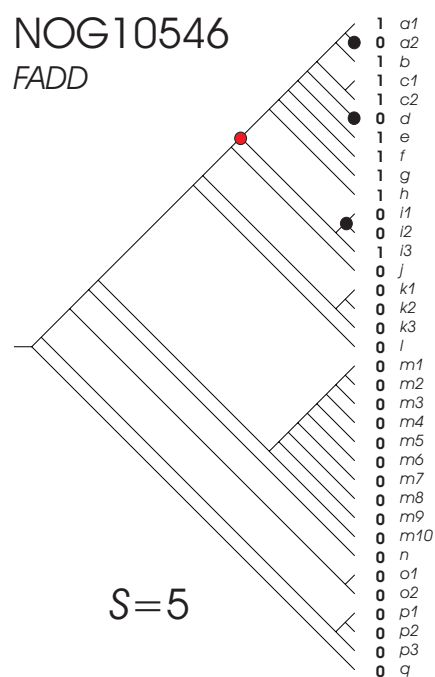
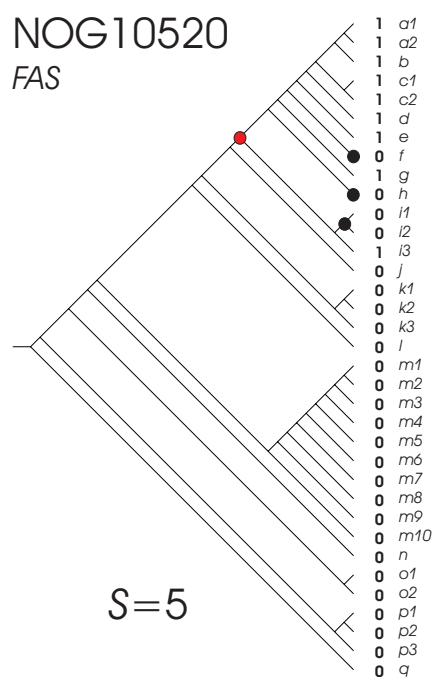


**Supplementary Figure S45.** Parsimony analysis of eukaryotic clusters of orthologous groups: NOG07367, NOG07483, NOG07555 and NOG08240, as in Supplementary Figure S14.

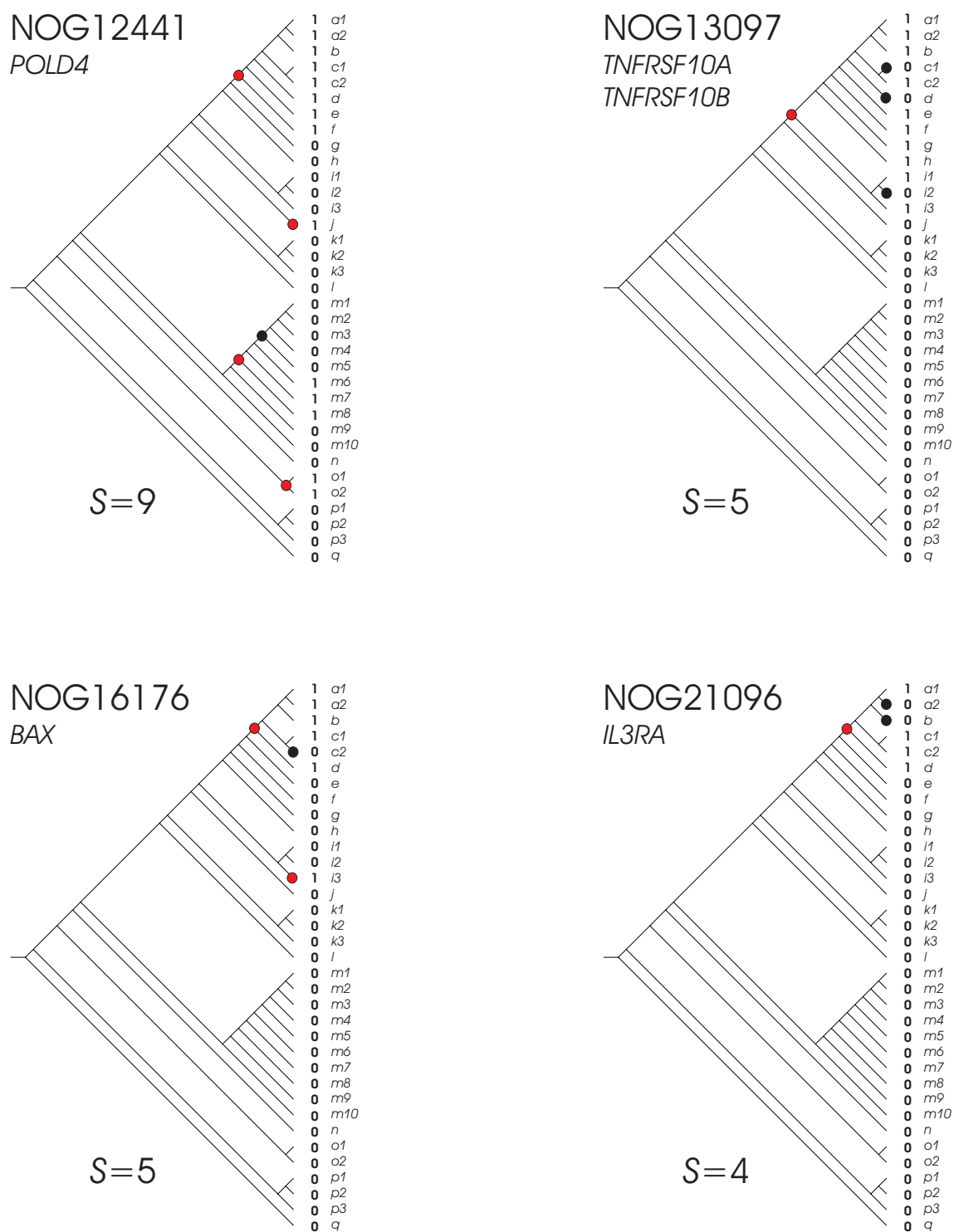


**Supplementary Figure S46.** Parsimony analysis of eukaryotic clusters of orthologous groups: NOG08276, NOG08316, NOG09287 and NOG10429, as in Supplementary Figure S14.

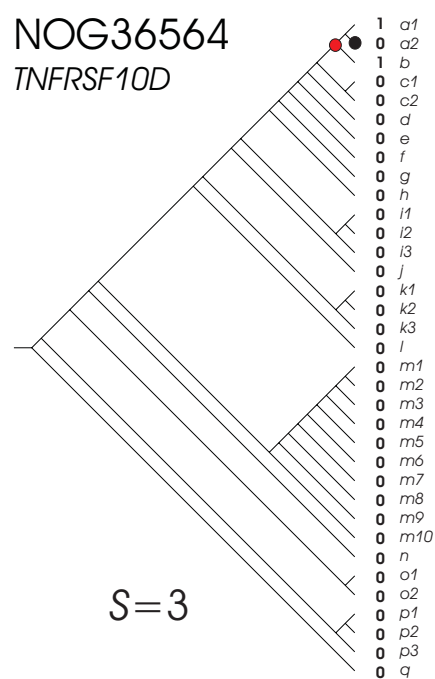
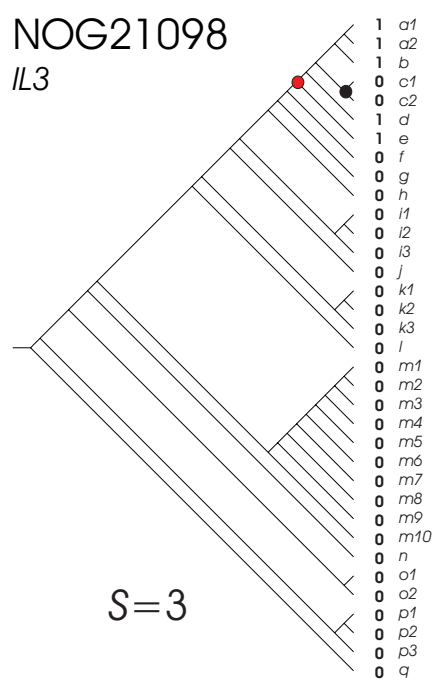




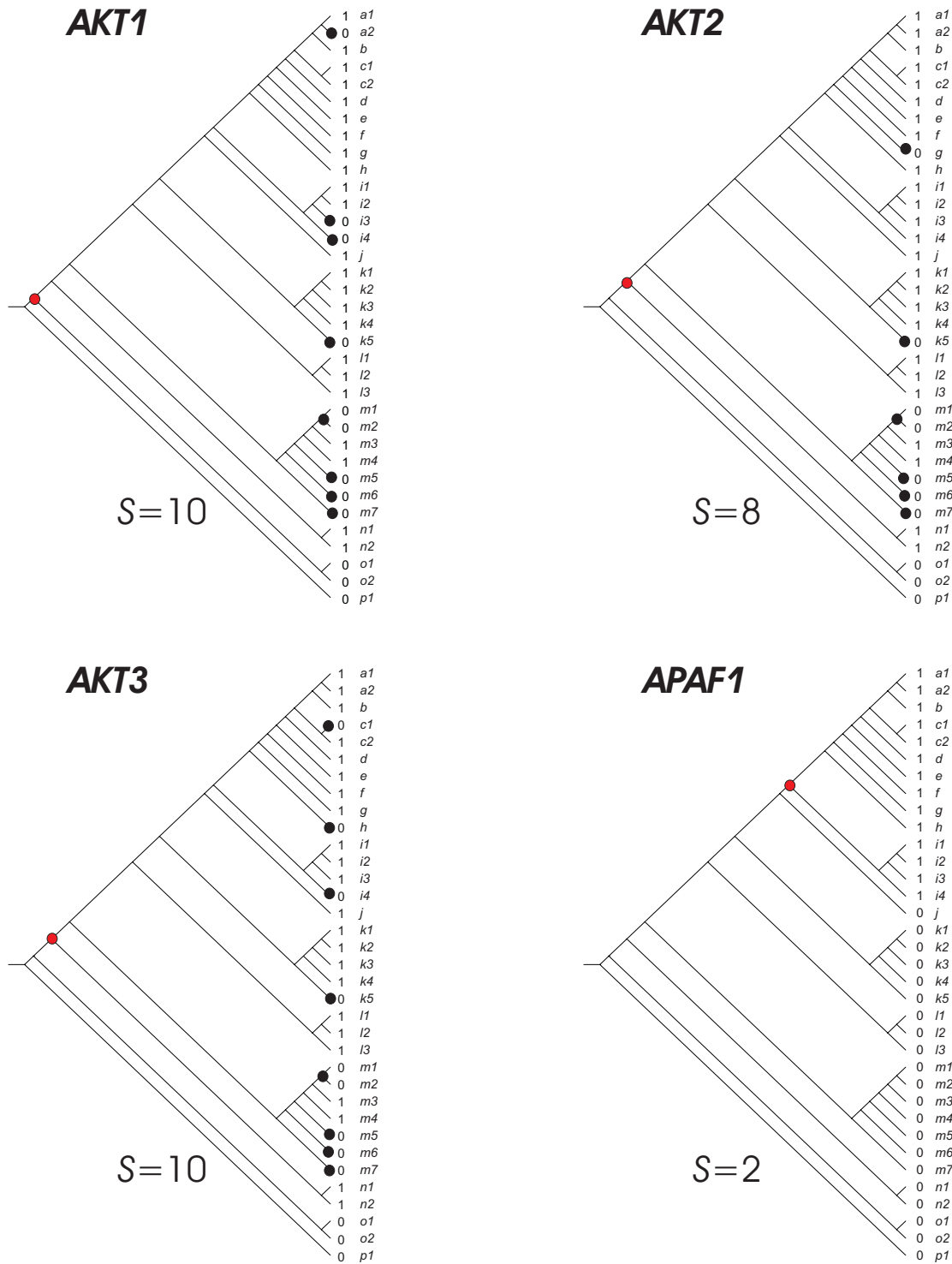
**Supplementary Figure S47.** Parsimony analysis of eukaryotic clusters of orthologous groups: NOG10520, NOG10546, NOG10549 and NOG10964, as in Supplementary Figure S14.



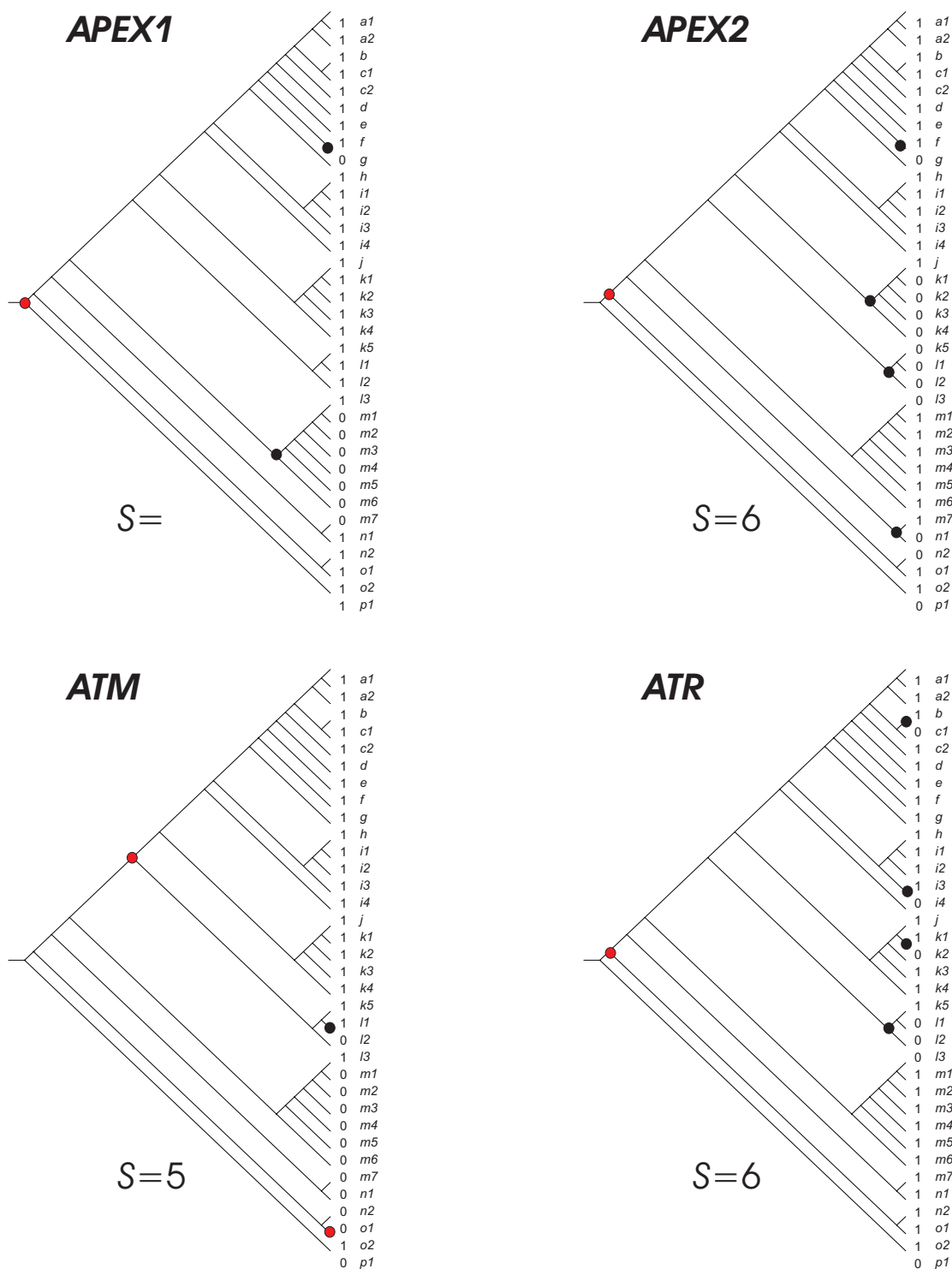
**Supplementary Figure S48.** Parsimony analysis of eukaryotic clusters of orthologous groups: NOG12441, NOG13097, NOG16176 and NOG21096, as in Supplementary Figure S14.



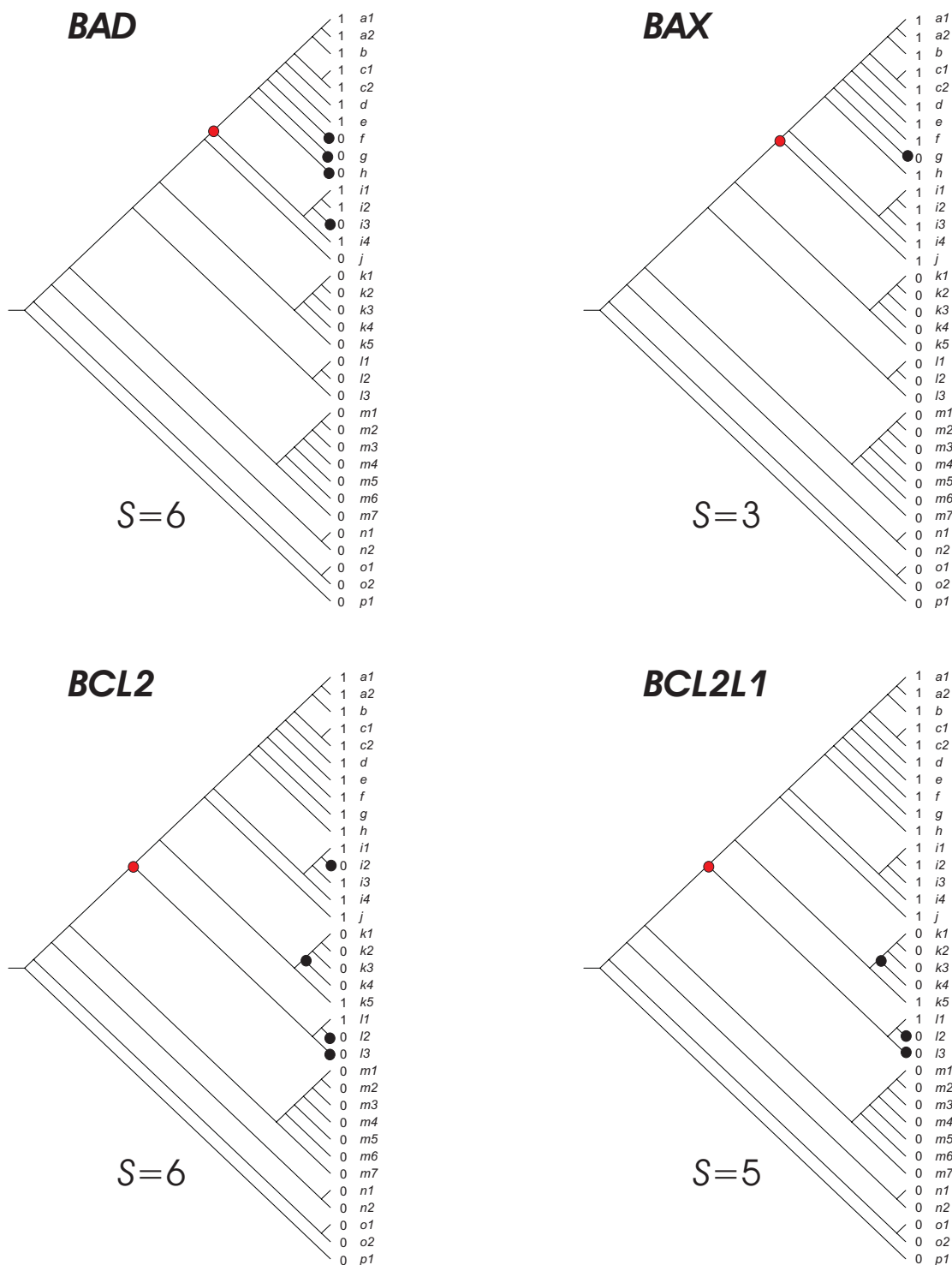
**Supplementary Figure S49.** Parsimony analysis of eukaryotic clusters of orthologous groups: NOG21098 and NOG36564, as in Supplementary Figure S14.



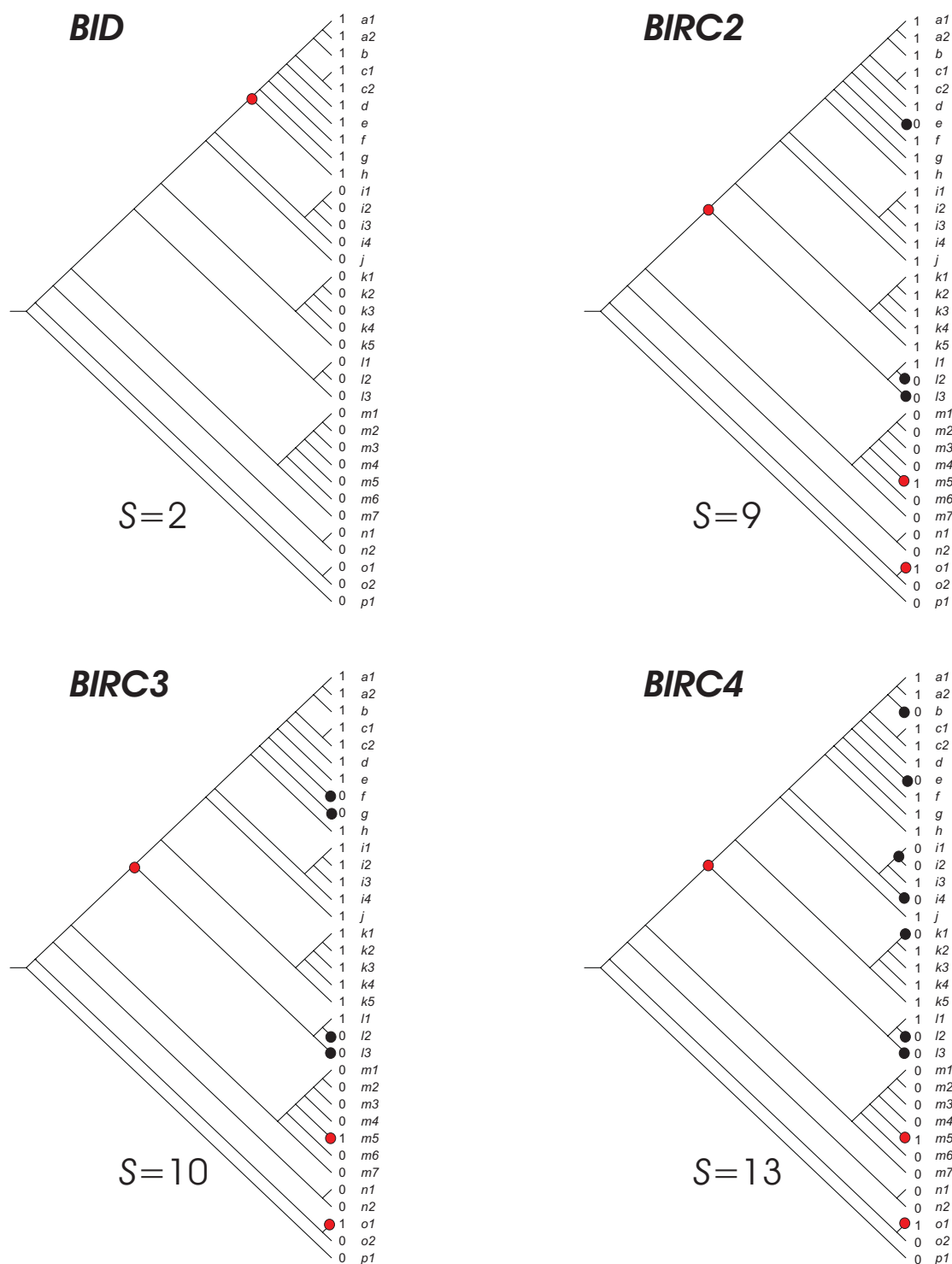
**Supplementary Figure S50.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes: **a1** (*Homo sapiens*); **a2** (*Pan troglodytes*); **b** (*Macaca mulatta*); **c1** (*Rattus norvegicus*); **c2** (*Mus musculus*); **d** (*Canis familiaris*); **e1** (*Bos Taurus*); **f** (*Monodelphis domestica*); **g** (*Gallus gallus*); **h** (*Xenopus tropicalis*); **i1** (*Takifugu rubripes*); **i2** (*Tetraodon nigroviridis*); **i3** (*Gasterosteus aculeatus*); **i4** (*Danio rerio*); **j** (*Ciona intestinalis*); **k1** (*Drosophila melanogaster*); **k2** (*Drosophila pseudoobscura*); **k3** (*Anopheles gambiae*); **k4** (*Aedes aegypti*); **k5** (*Apis mellifera*); **l1** (*Caenorhabditis briggsae*); **l2** (*Caenorhabditis remanei*); **l3** (*Caenorhabditis elegans*); **m1** (*Kluyveromyces lactis*); **m2** (*Saccharomyces cerevisiae*); **m3** (*Candida glabrata*); **m4** (*Debaryomyces hansenii*); **m5** (*Yarrowia lipolytica*); **m6** (*Schizosaccharomyces pombe*); **m7** (*Cryptococcus neoformans*); **n1** (*Entamoeba histolytica*); **n2** (*Dictyostelium discoideum*); **o1** (*Arabidopsis thaliana*); **o2** (*Oryza sativa*); **p1** (*Escherichia coli*K12);



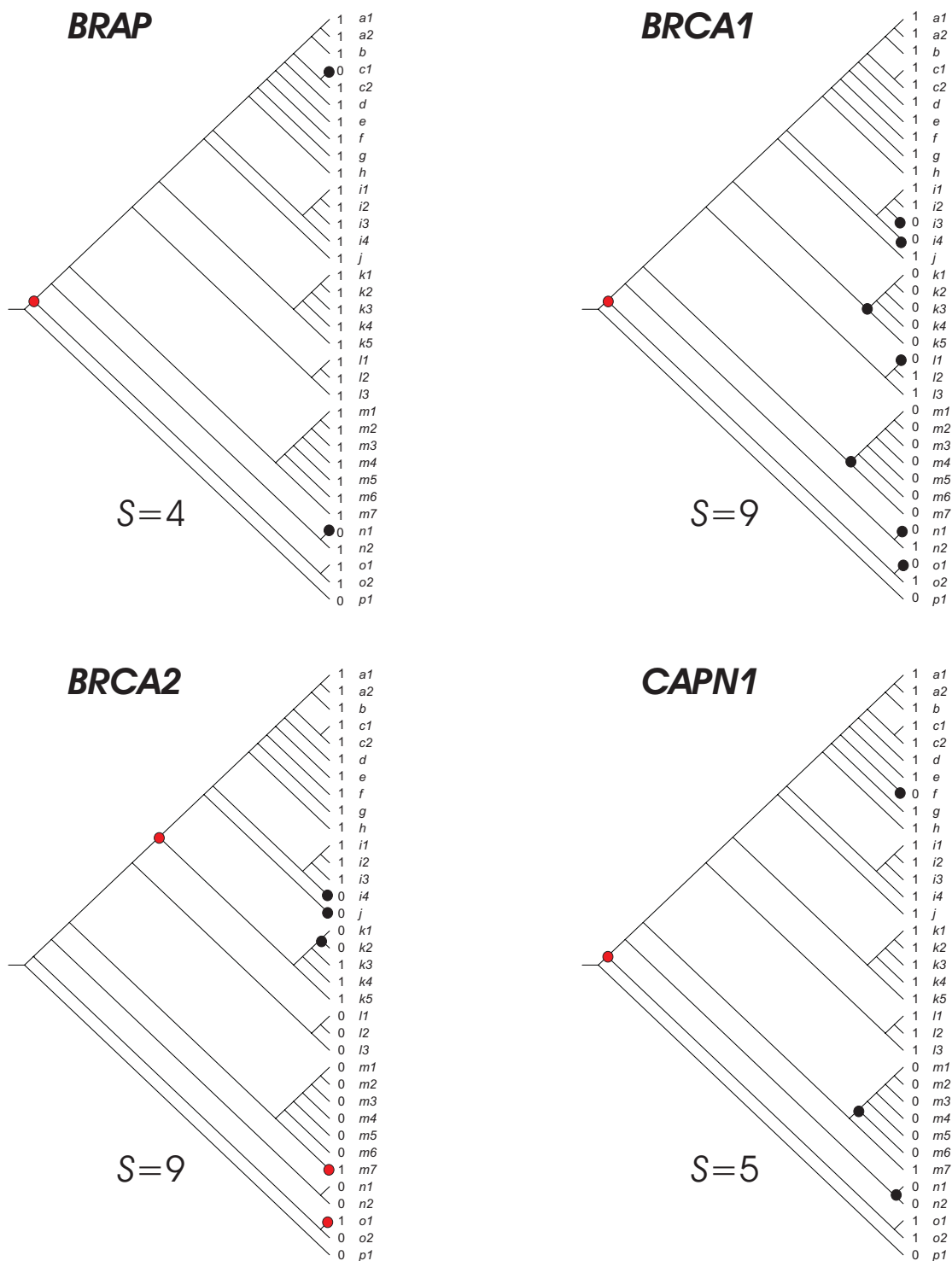
**Supplementary Figure S51.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



**Supplementary Figure S52.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

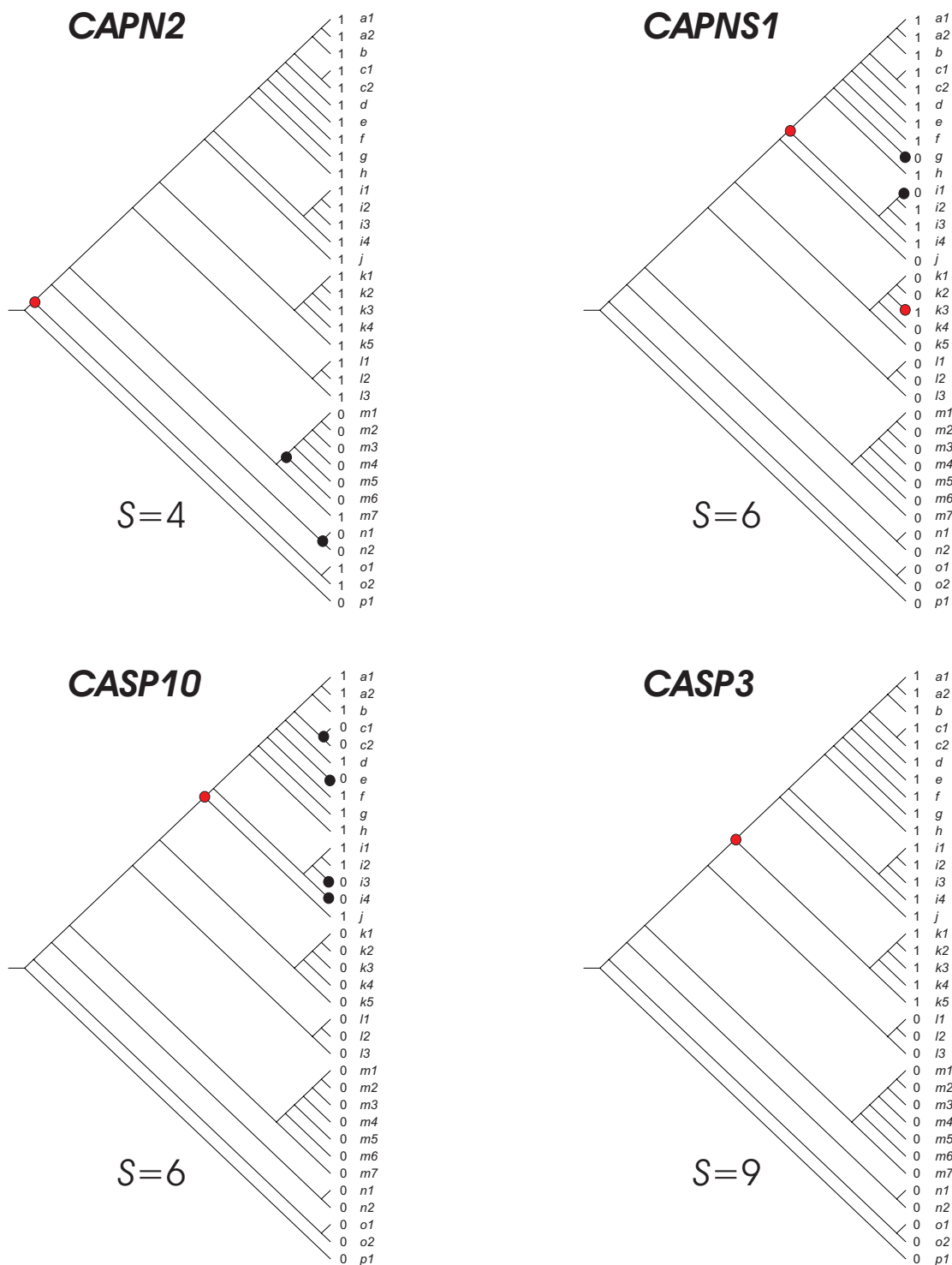


**Supplementary Figure S53.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

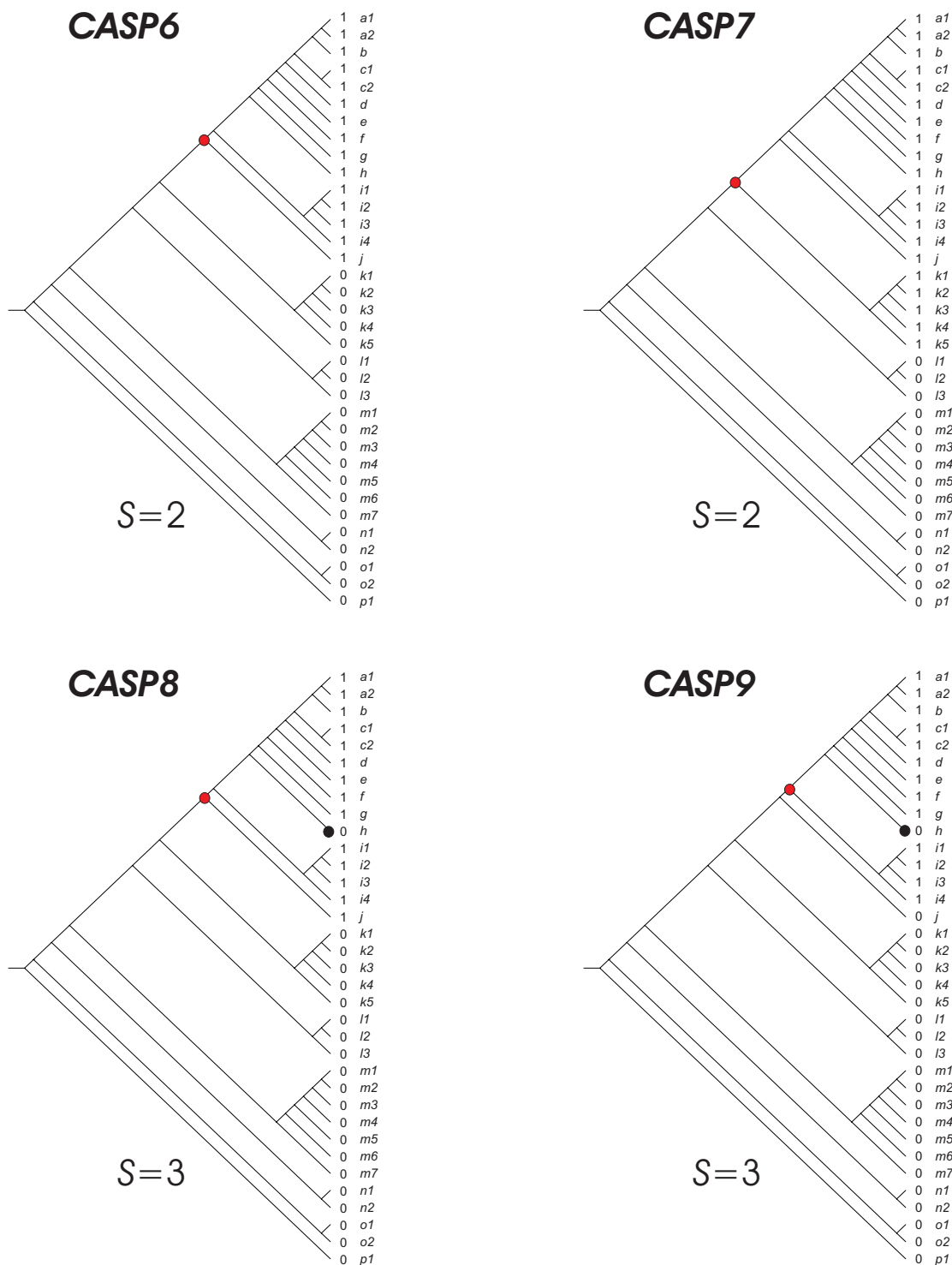


**Supplementary Figure S54.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

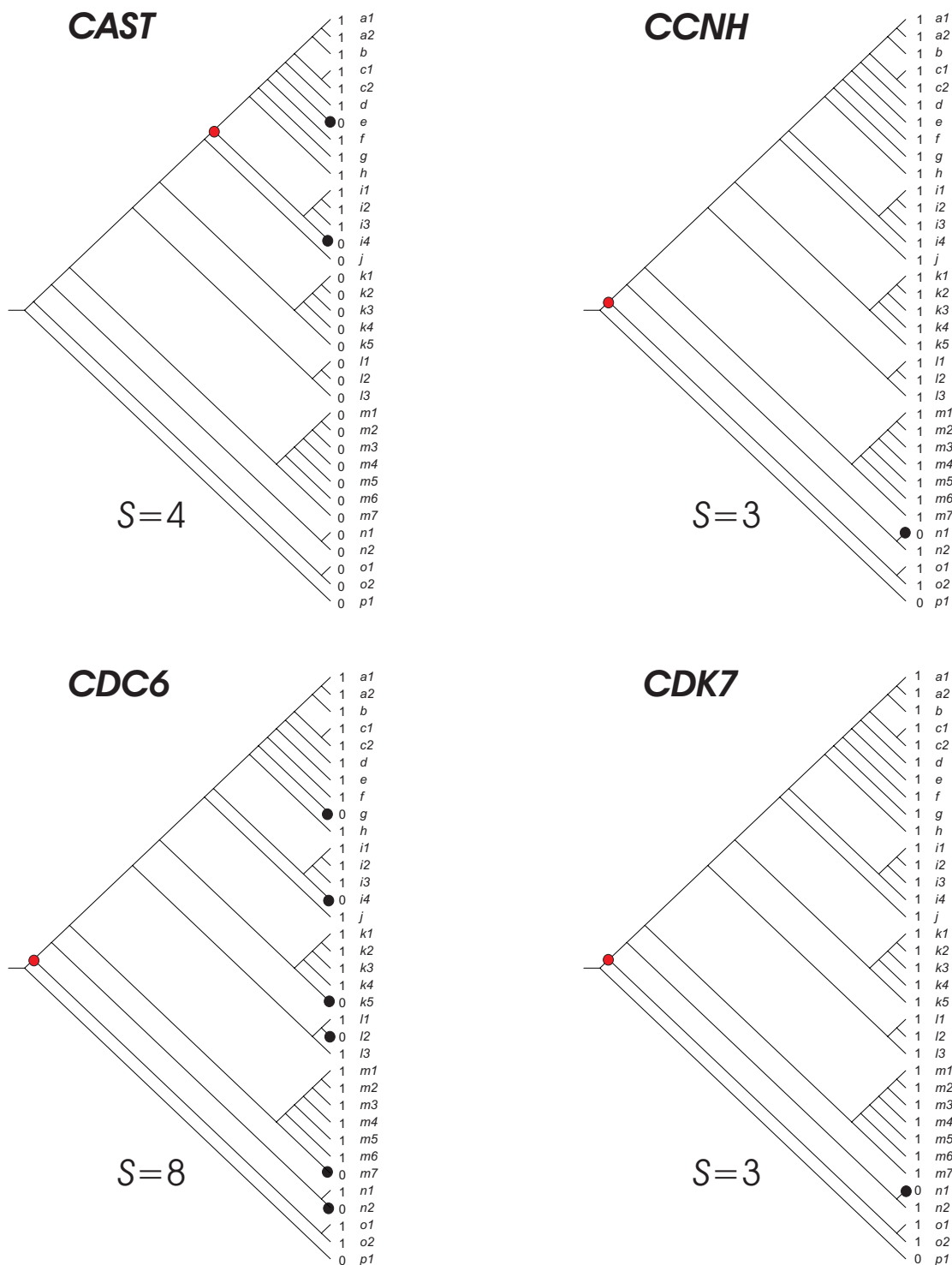




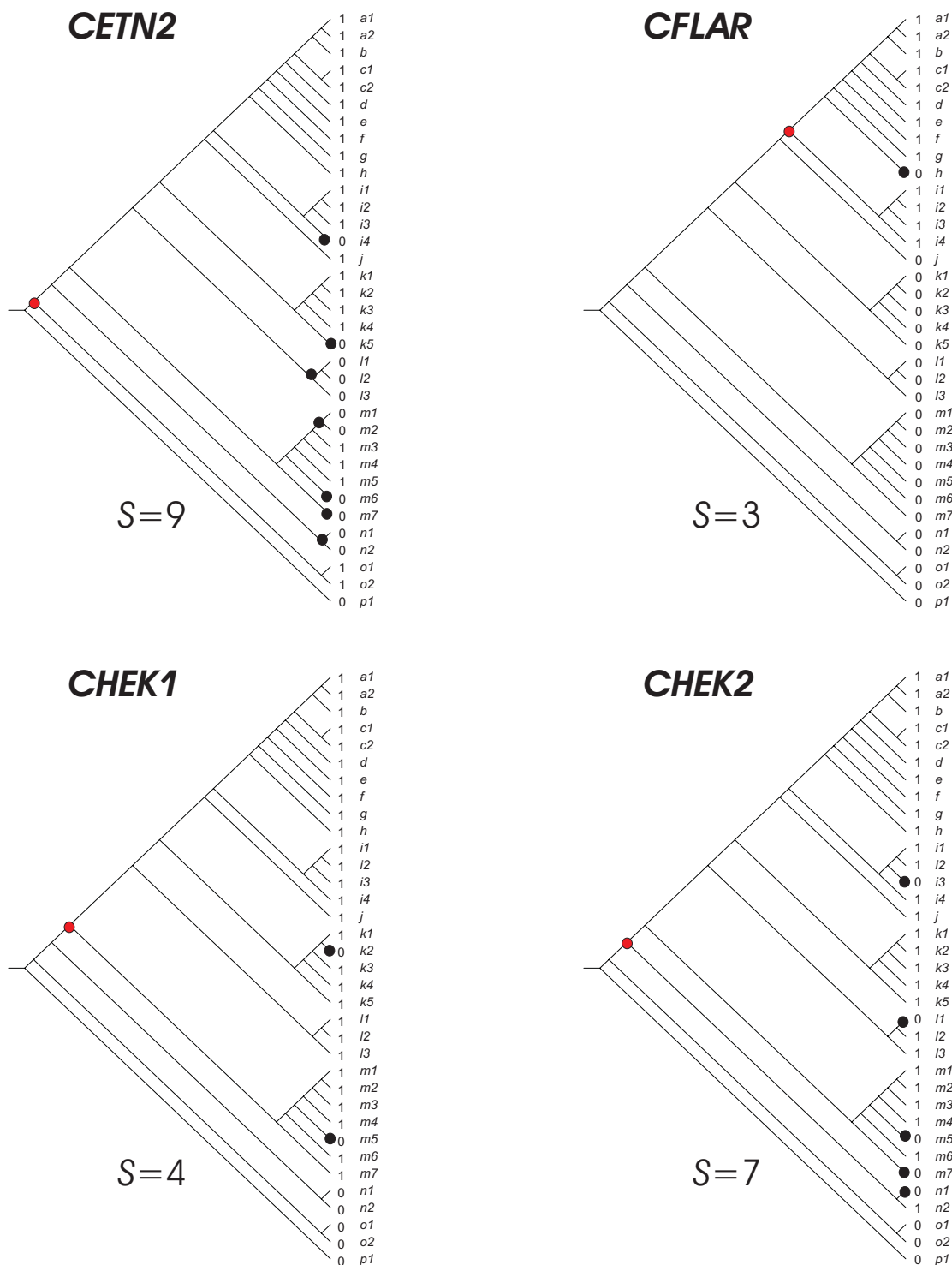
**Supplementary Figure S55.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



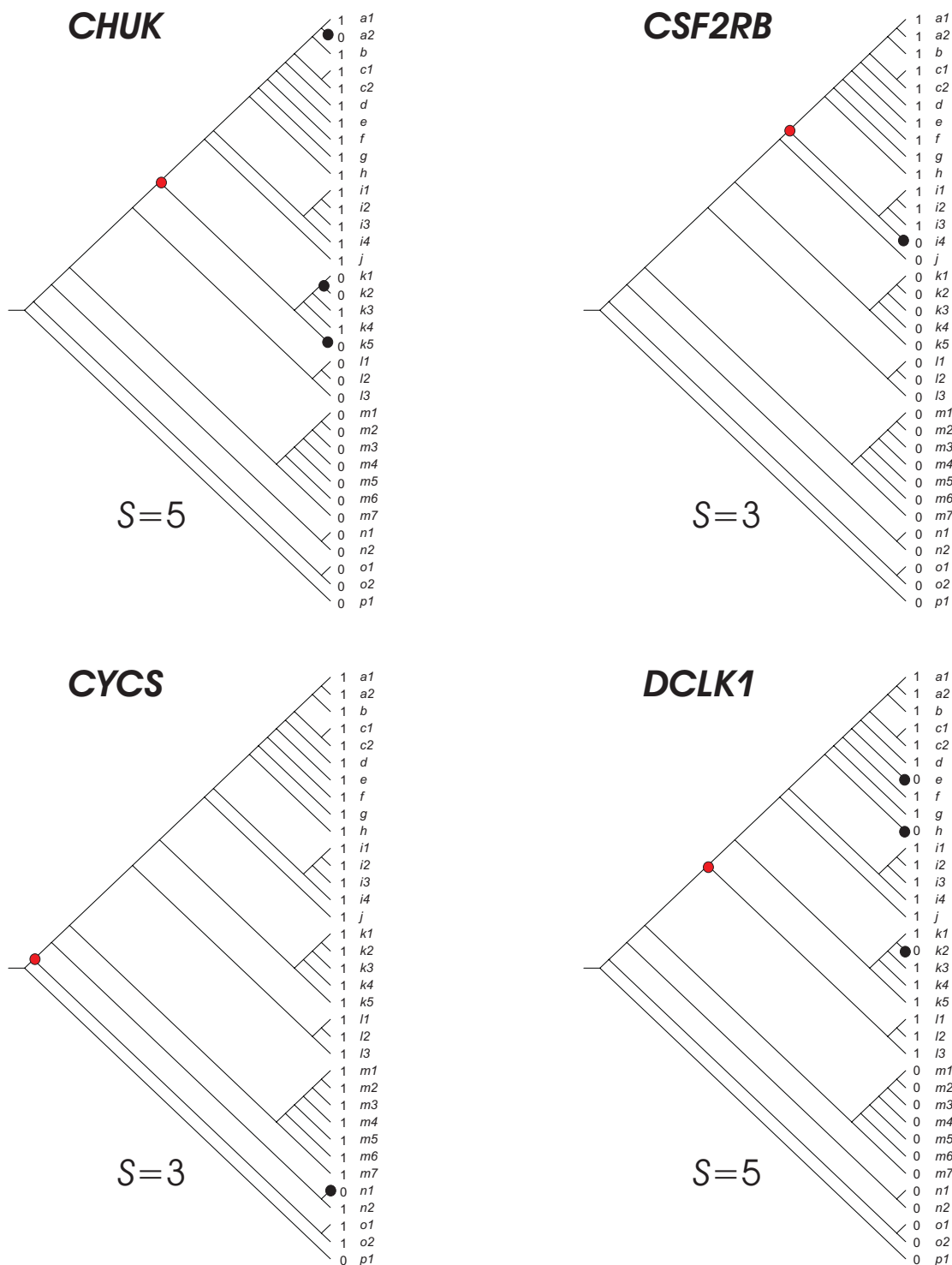
**Supplementary Figure S56.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



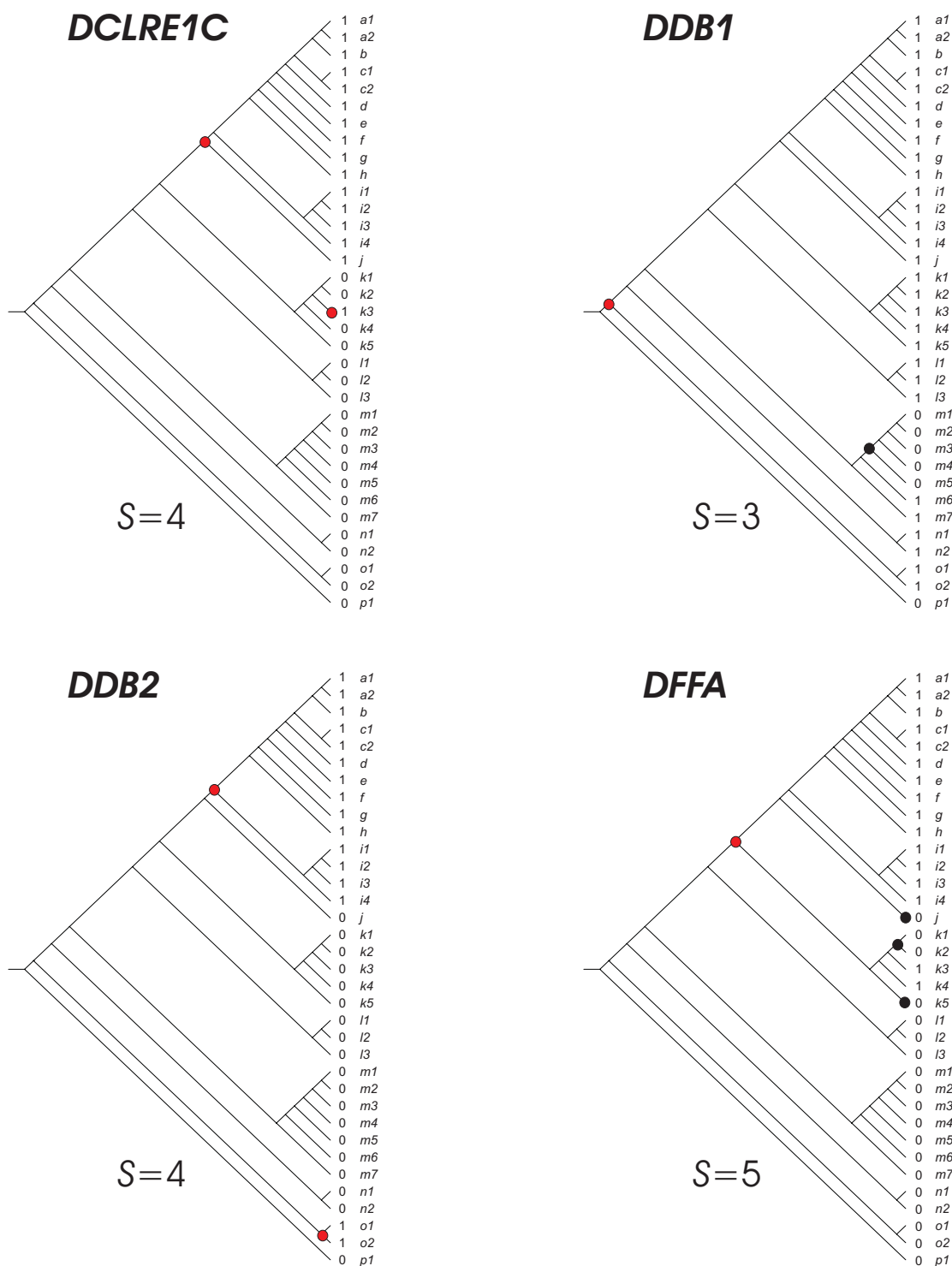
**Supplementary Figure S57.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



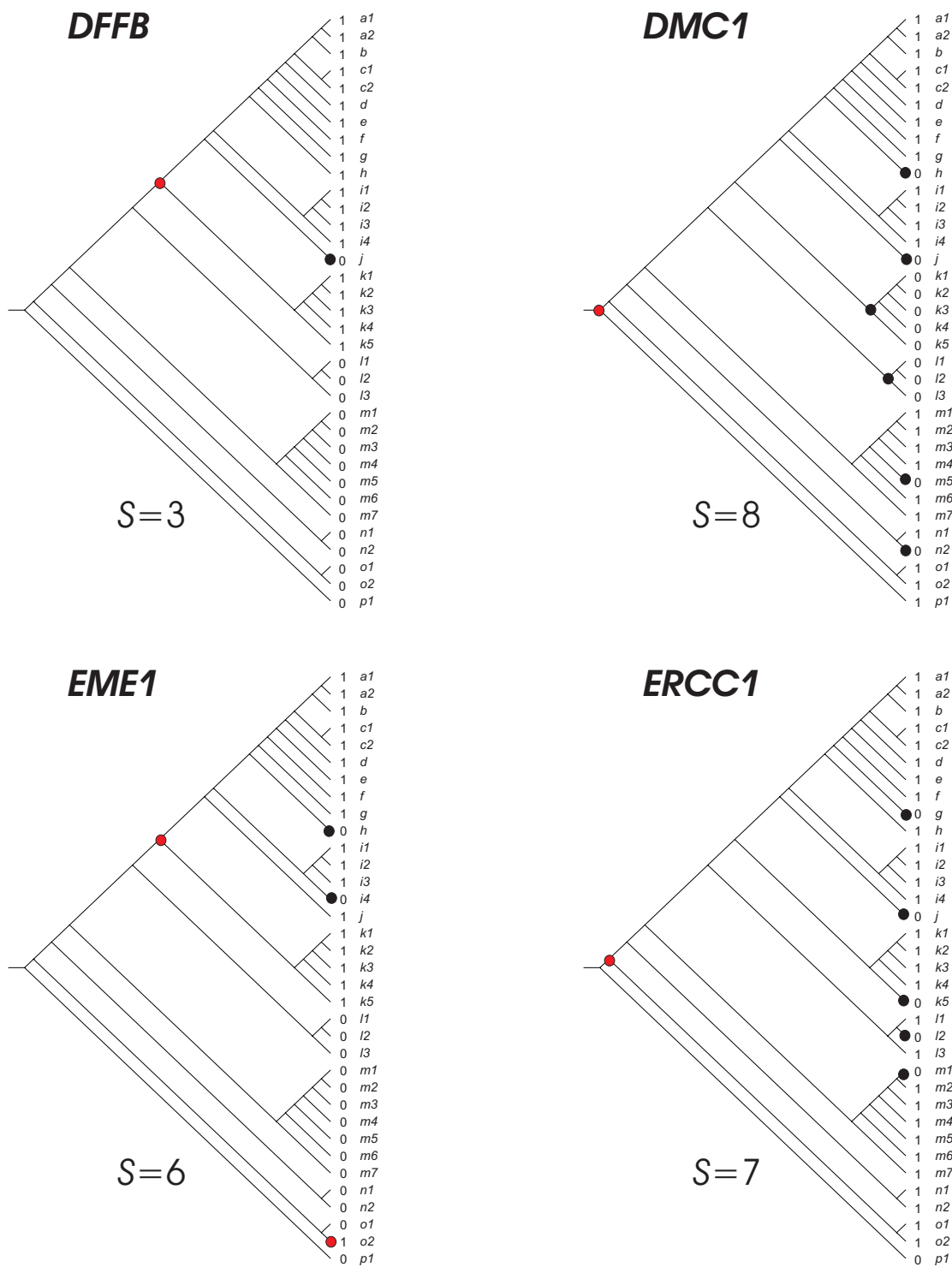
**Supplementary Figure S58.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



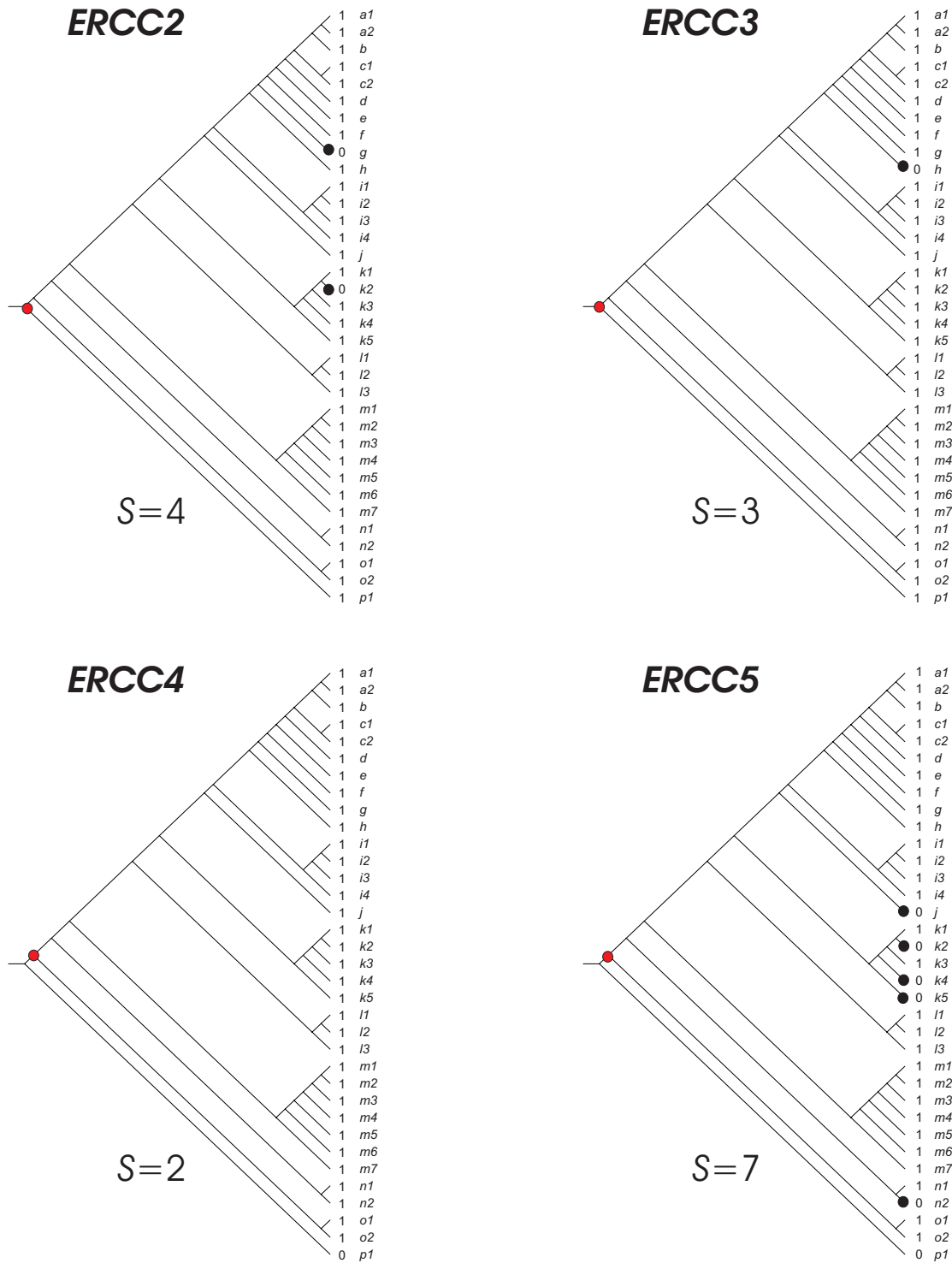
**Supplementary Figure S59.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



**Supplementary Figure S60.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

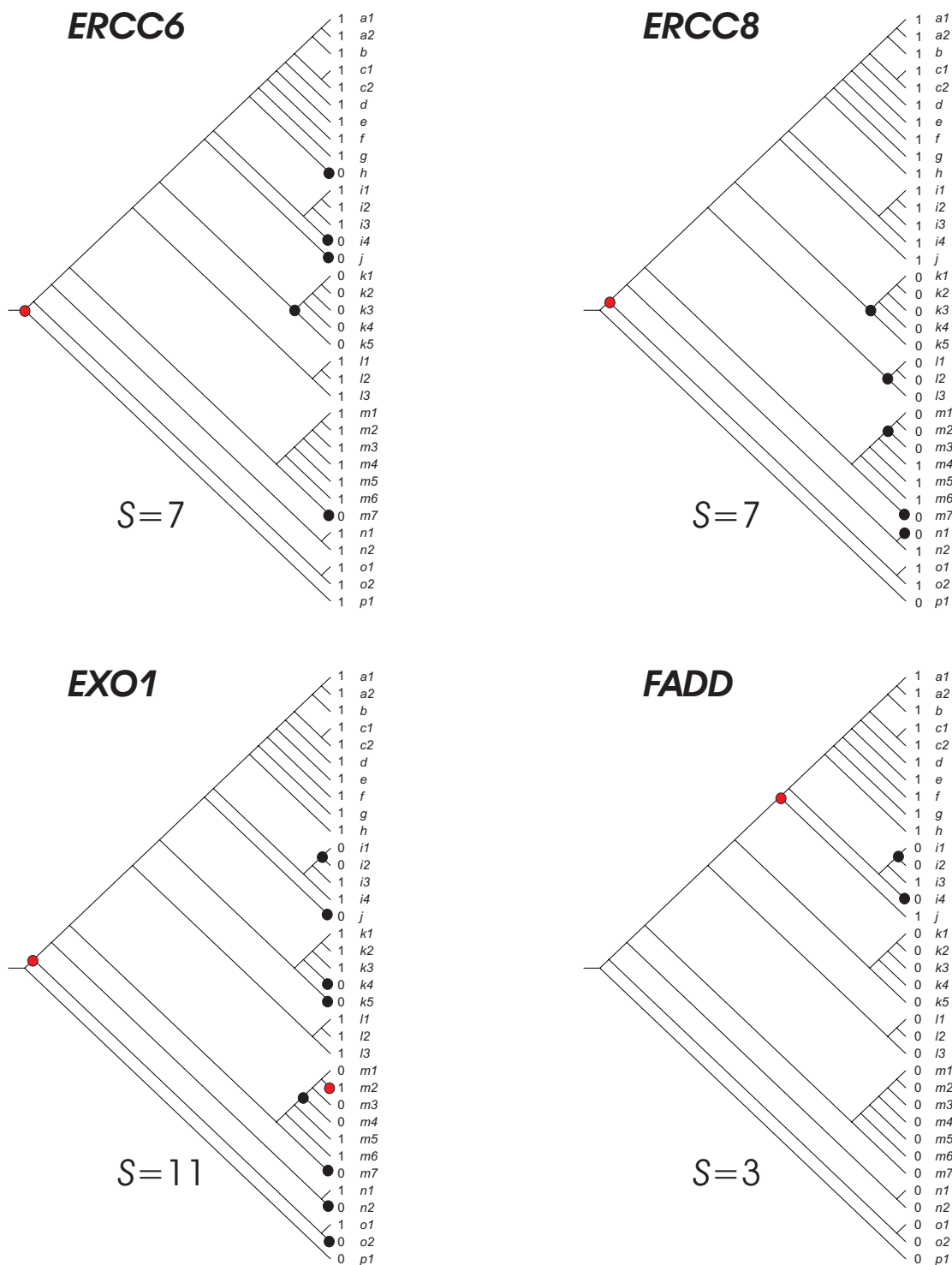


**Supplementary Figure S61.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

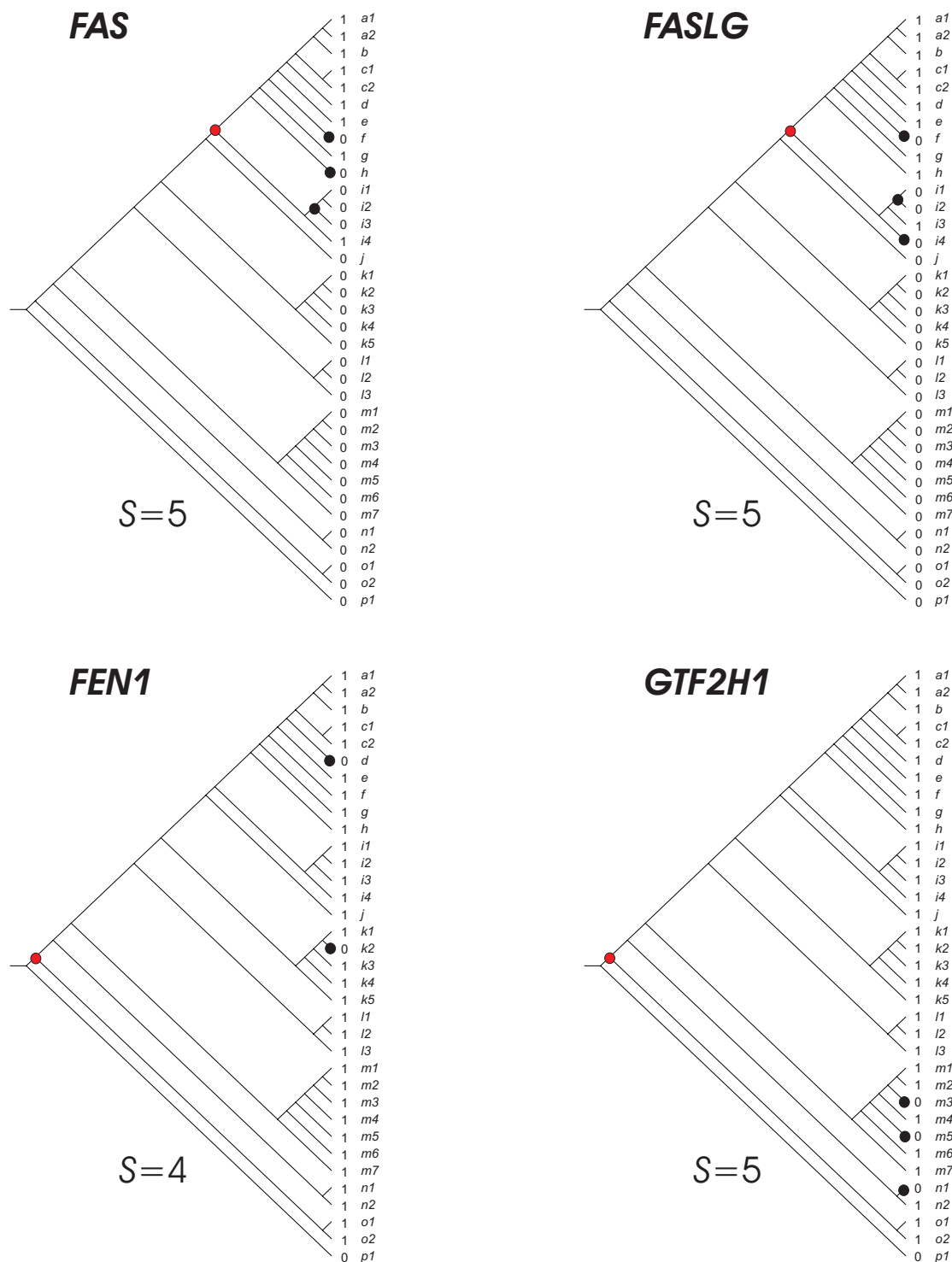


**Supplementary Figure S62.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

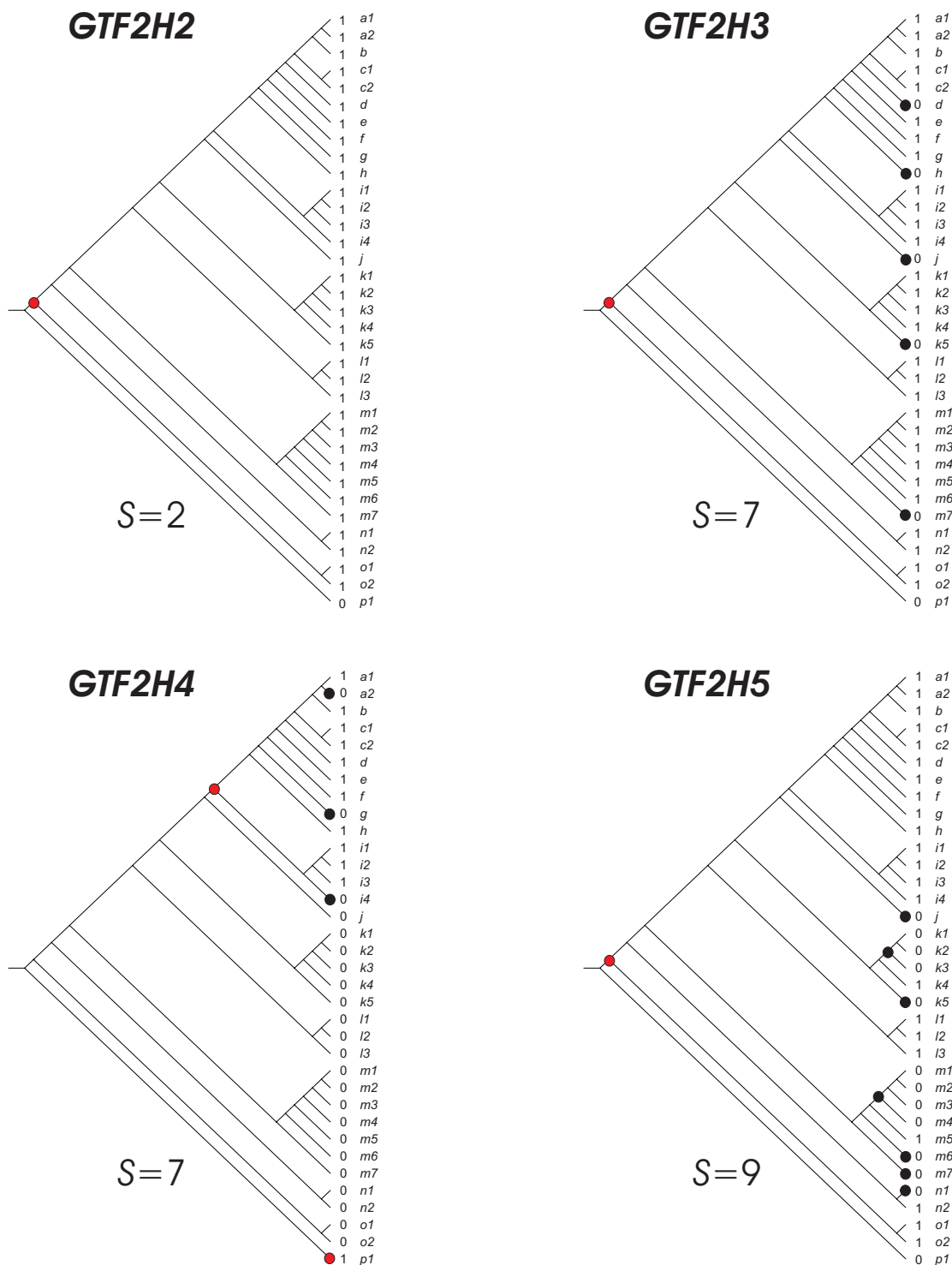




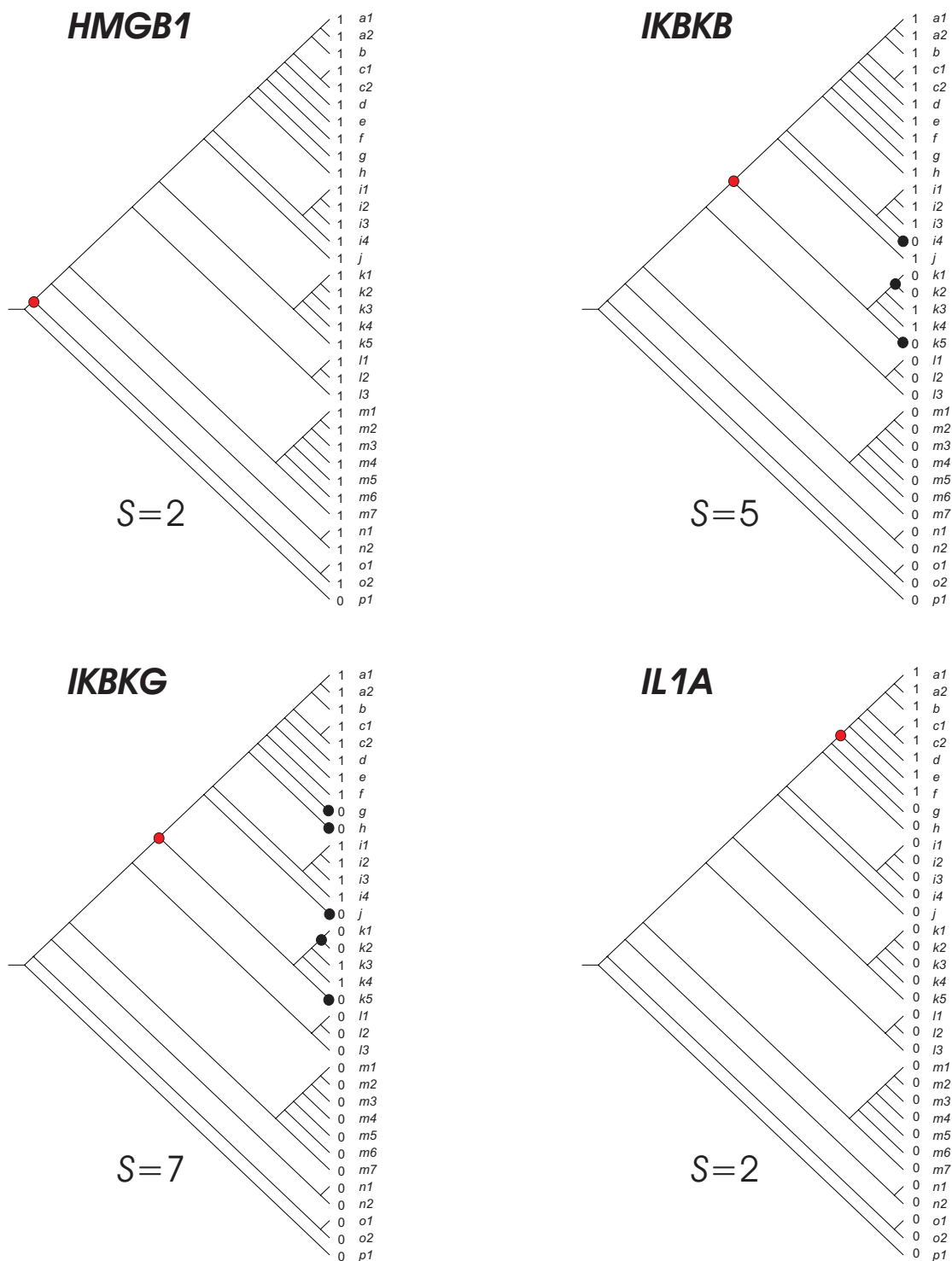
**Supplementary Figure S63.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



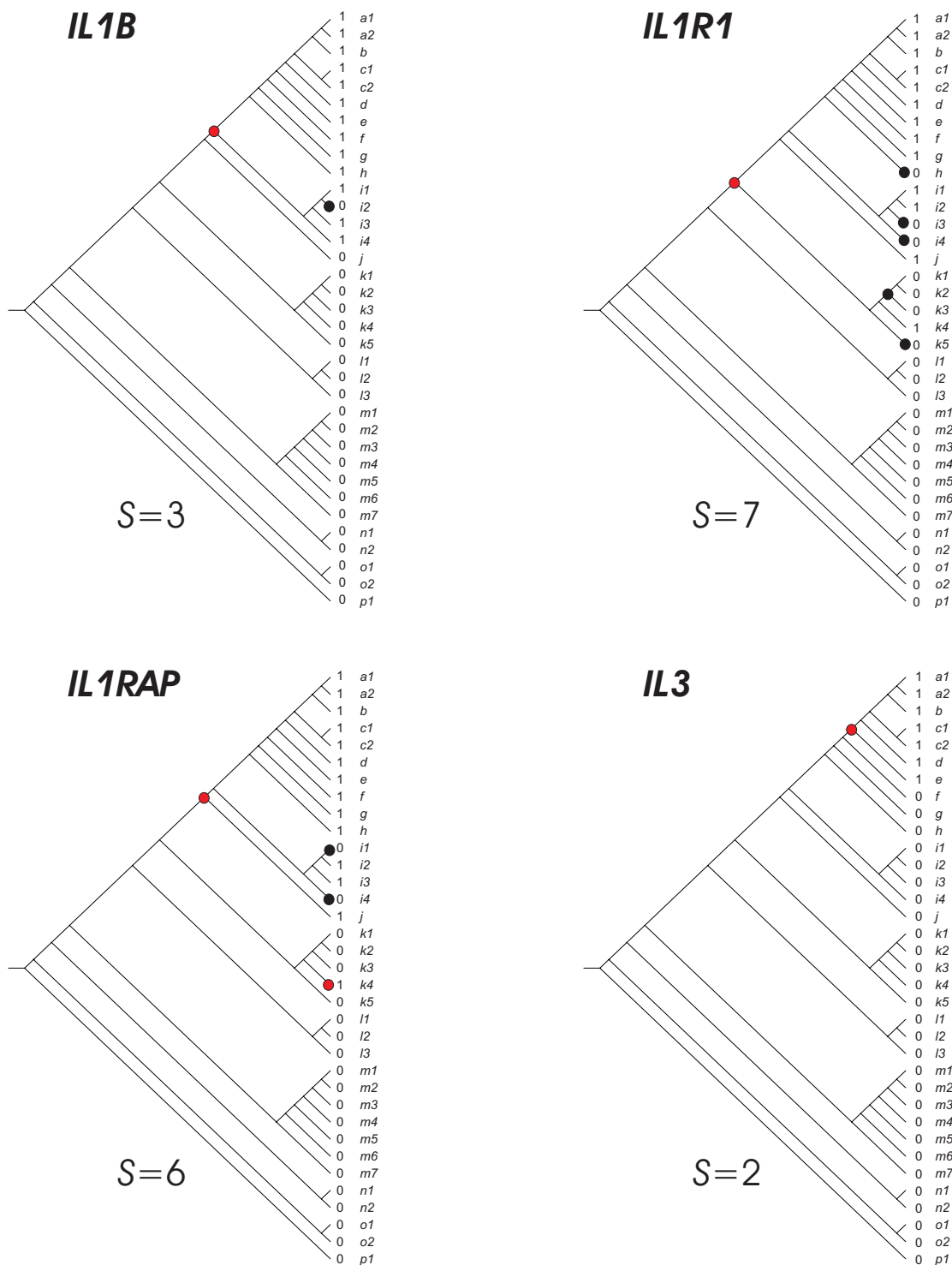
**Supplementary Figure S64.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



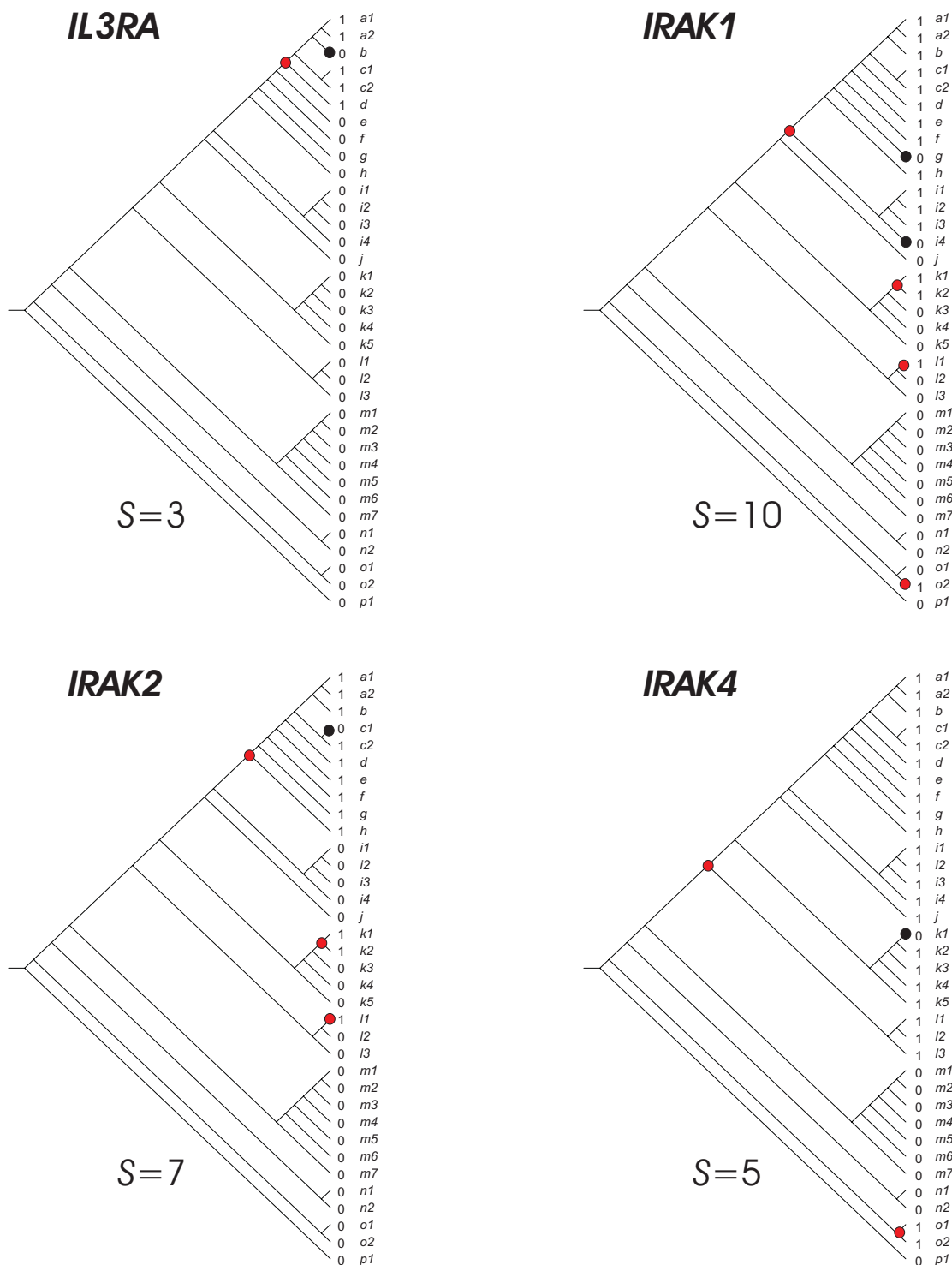
**Supplementary Figure S65.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



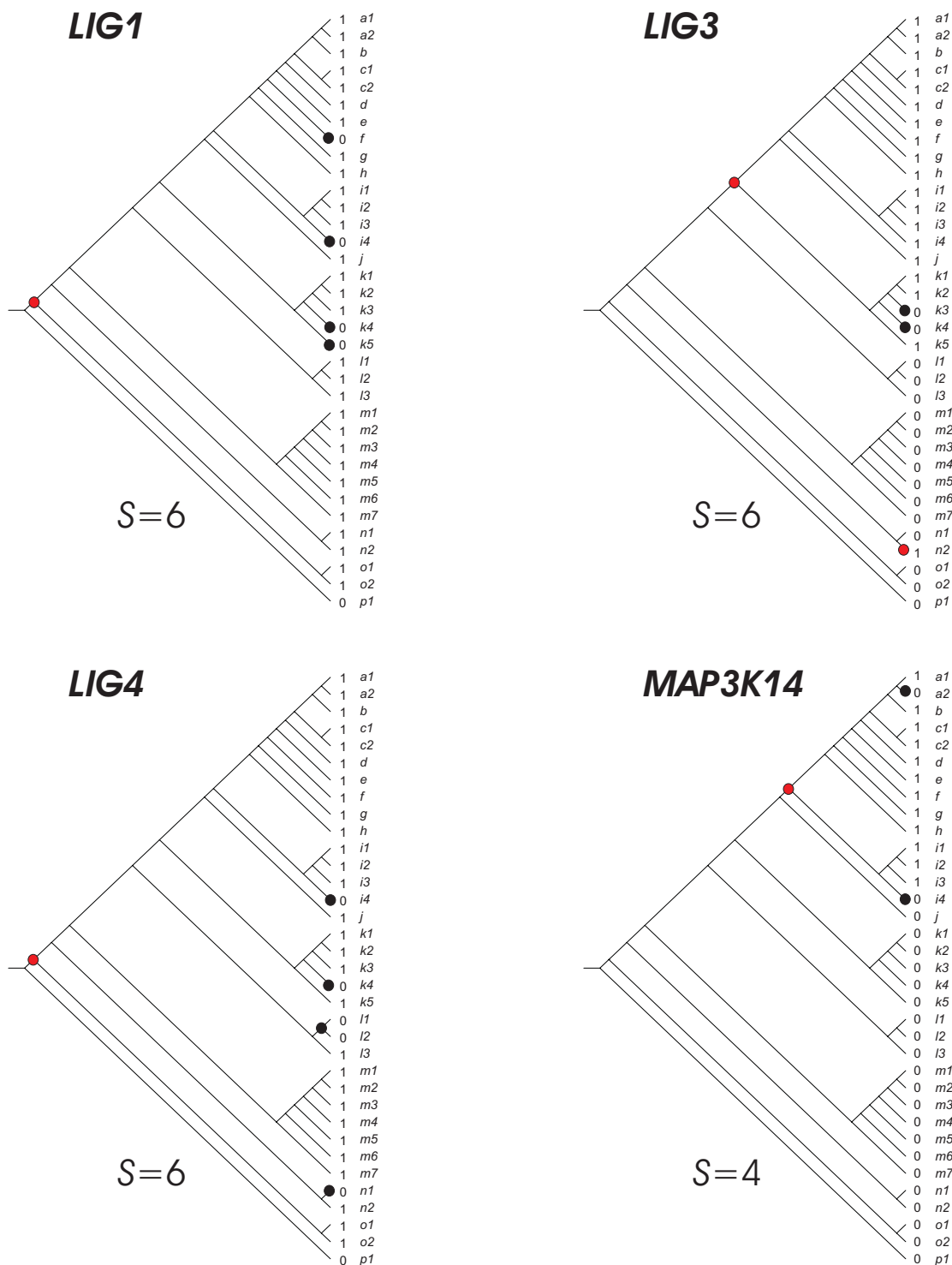
**Supplementary Figure S66.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



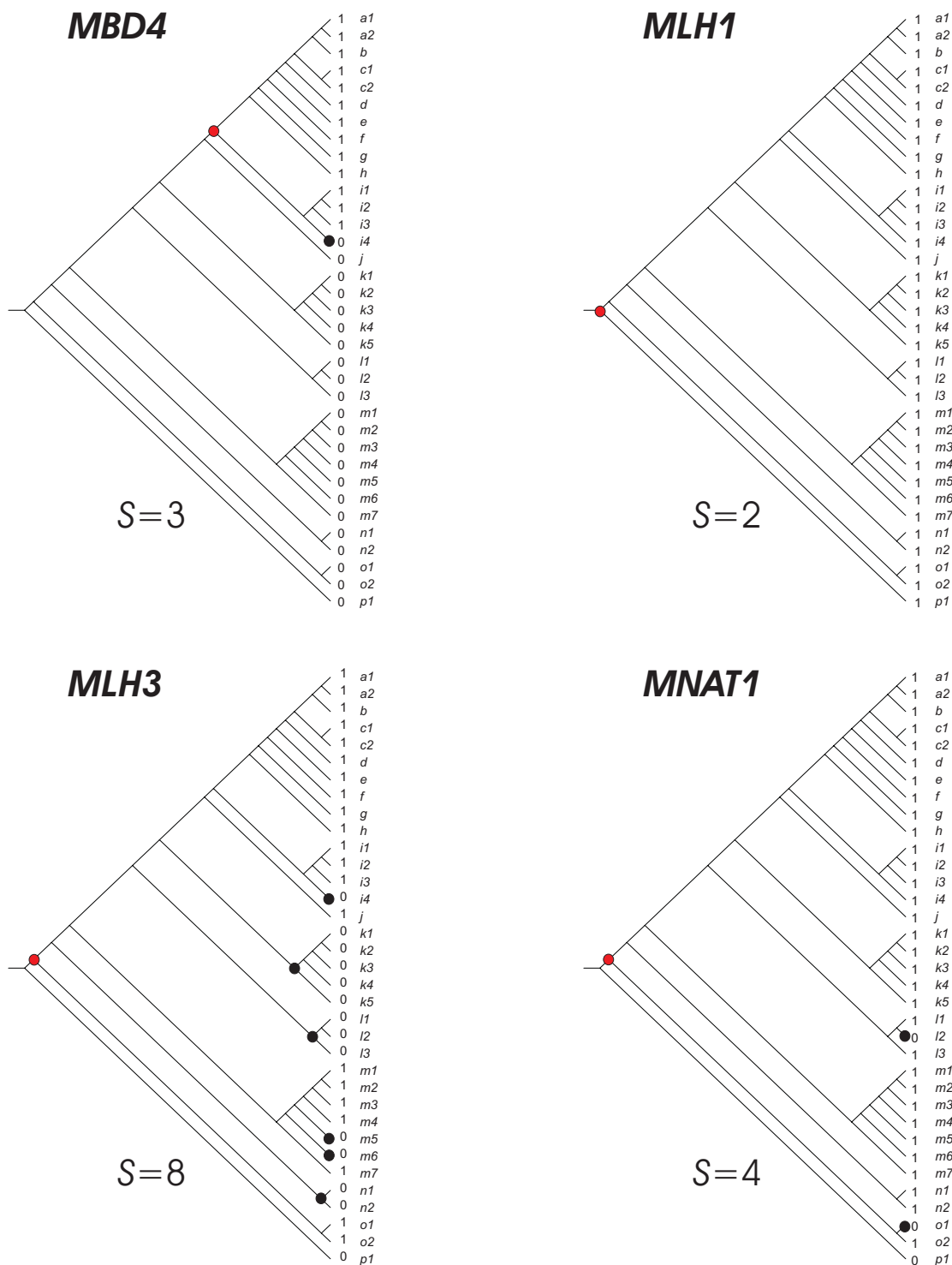
**Supplementary Figure S67.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



**Supplementary Figure S68.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

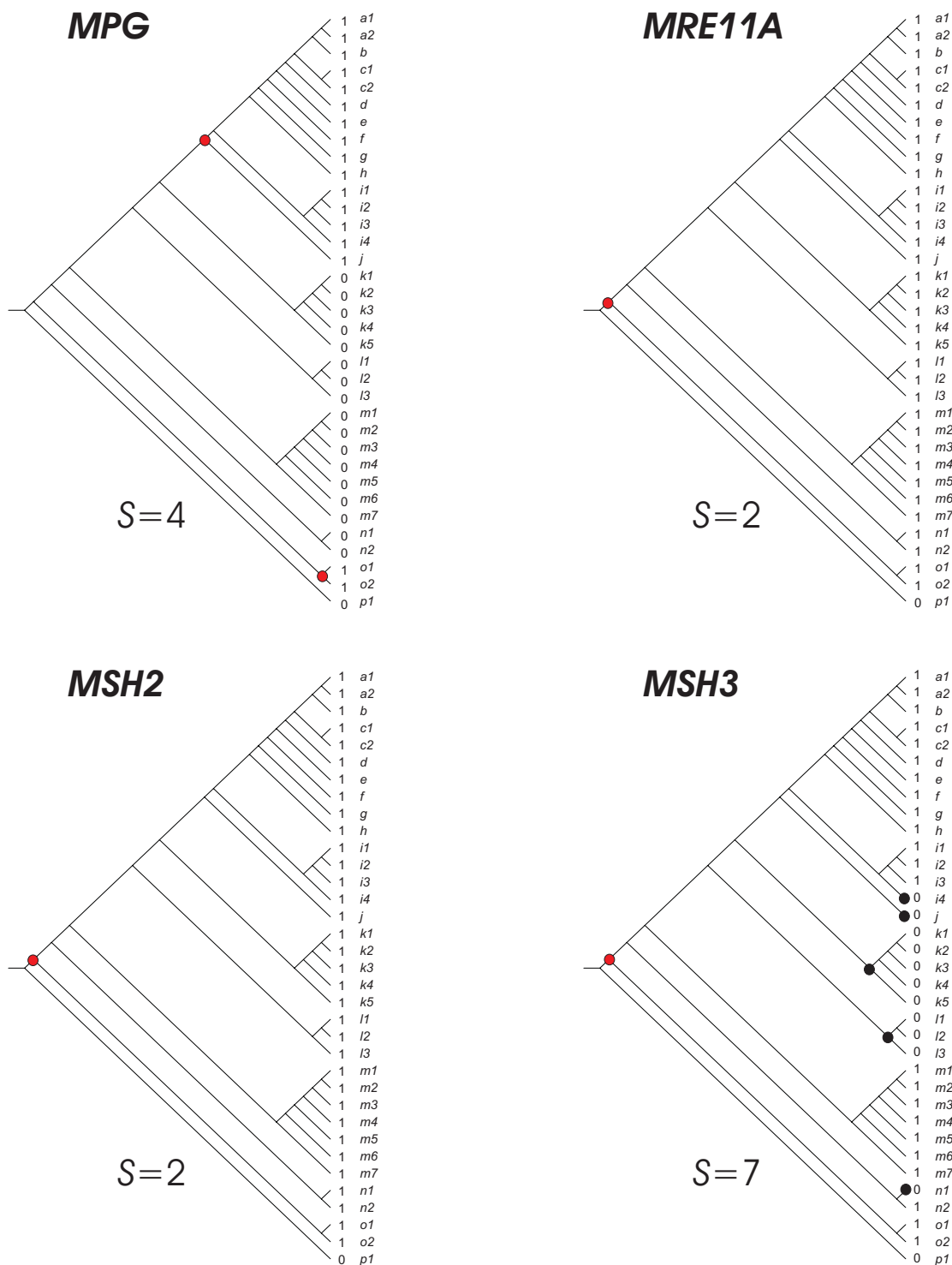


**Supplementary Figure S69.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

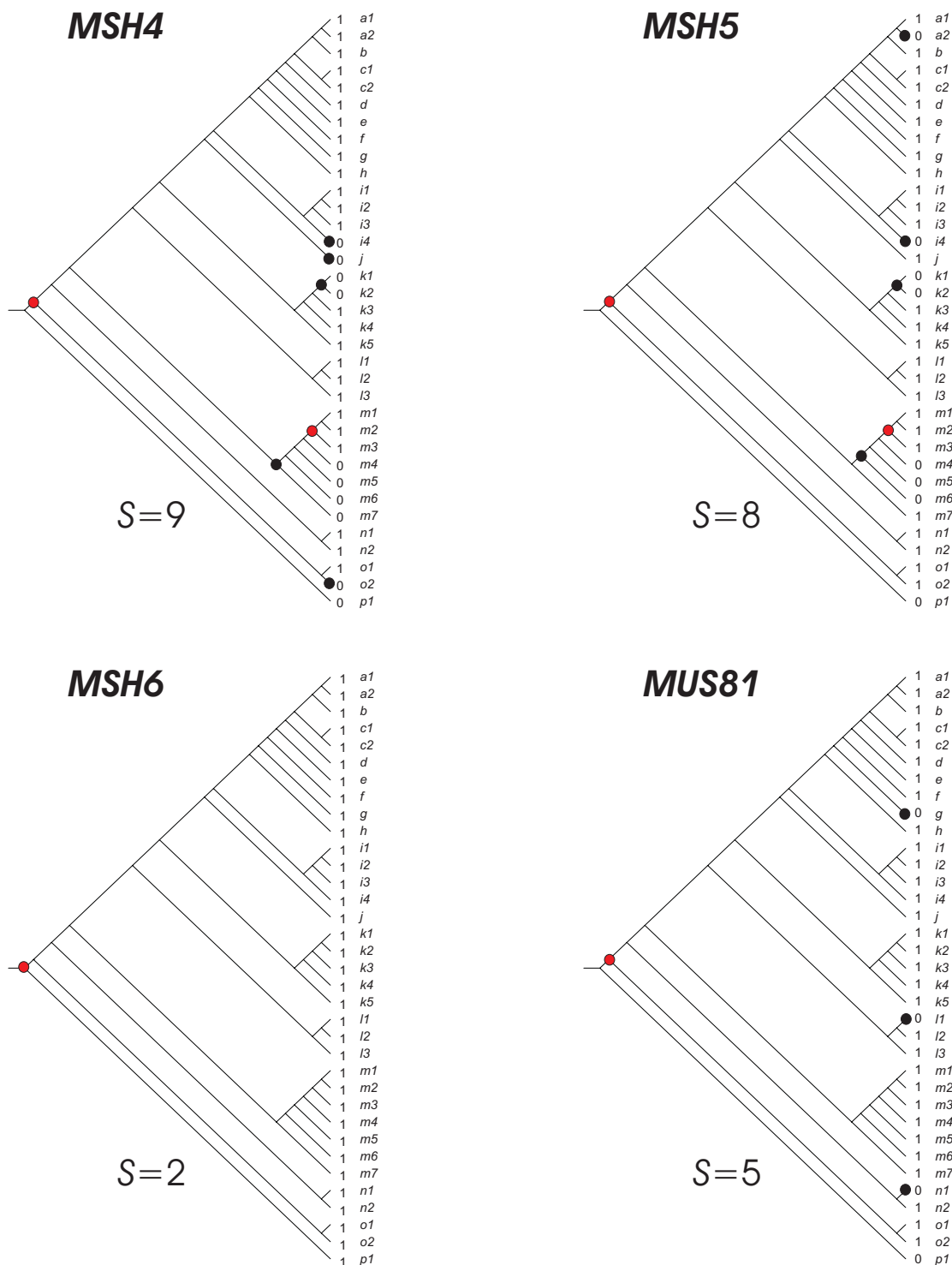


**Supplementary Figure S70.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

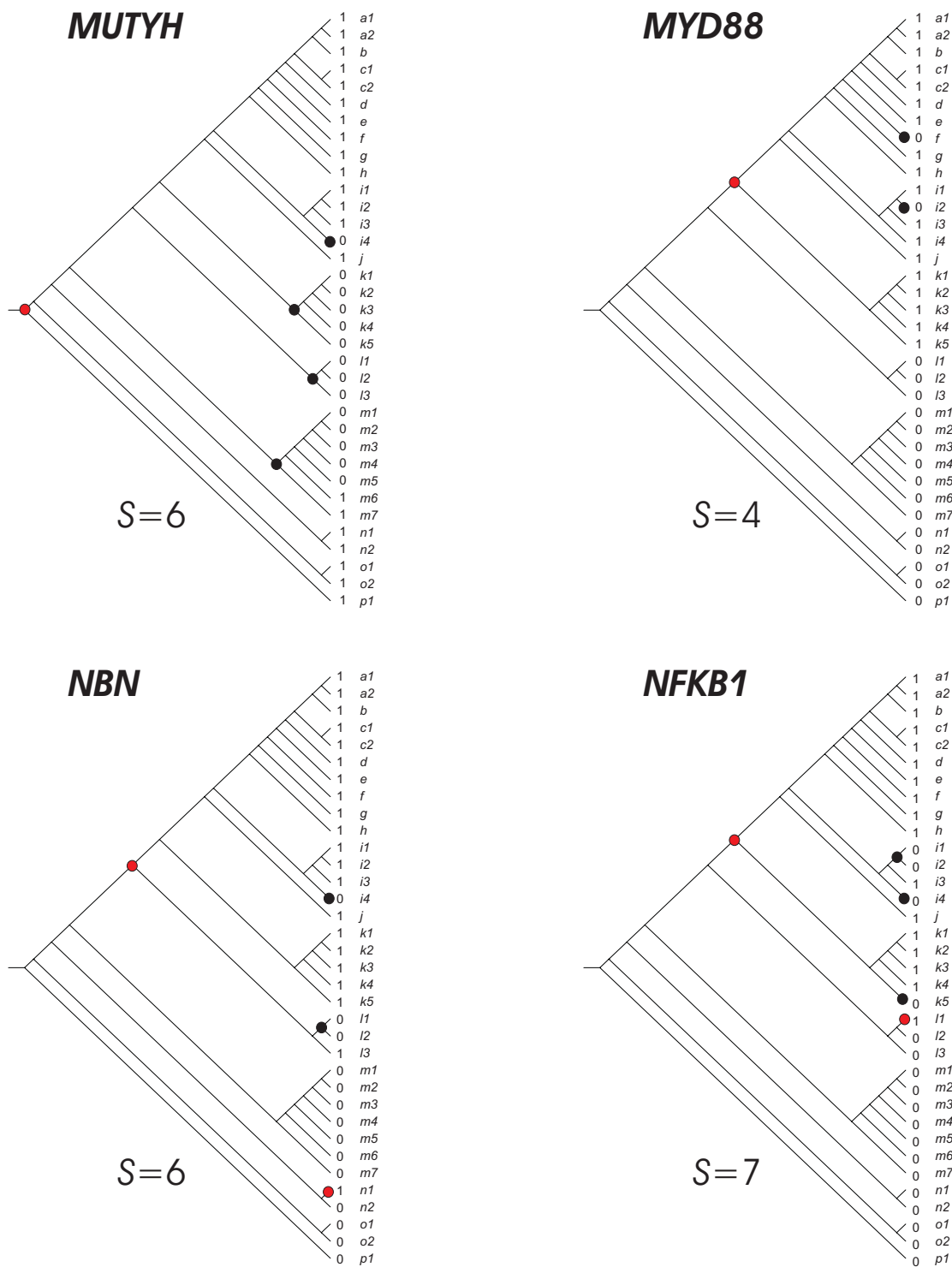




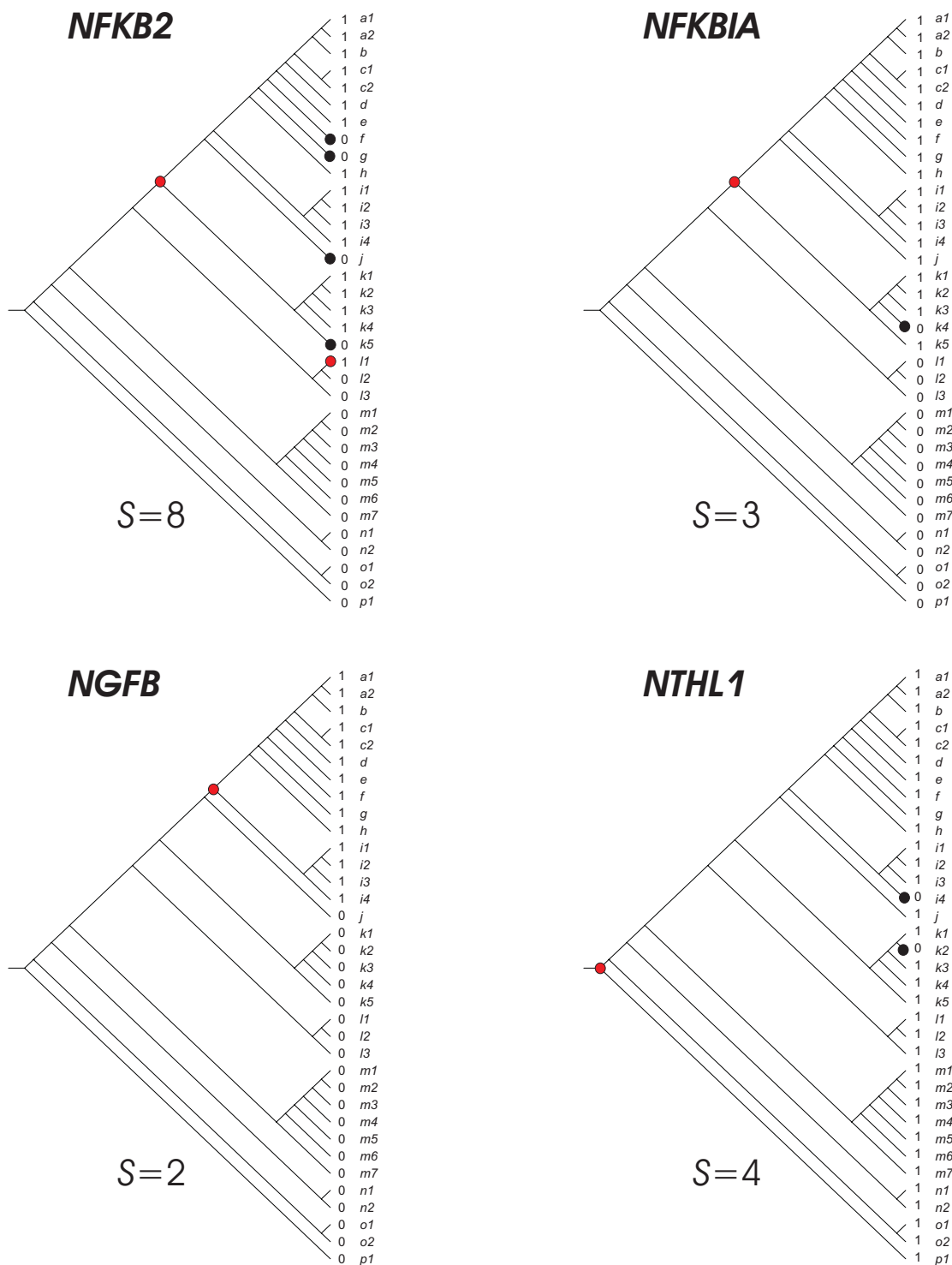
**Supplementary Figure S71.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



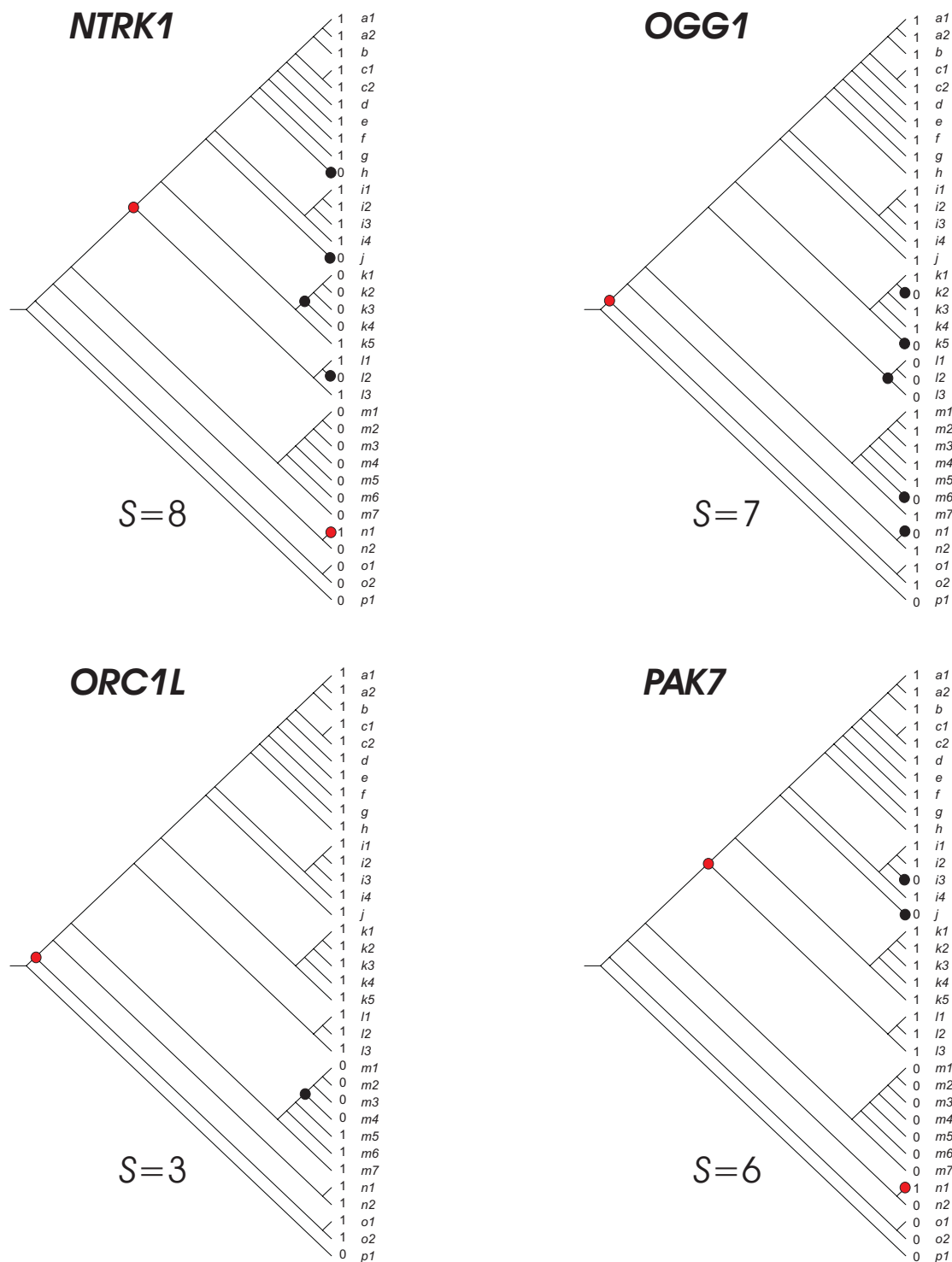
**Supplementary Figure S72.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



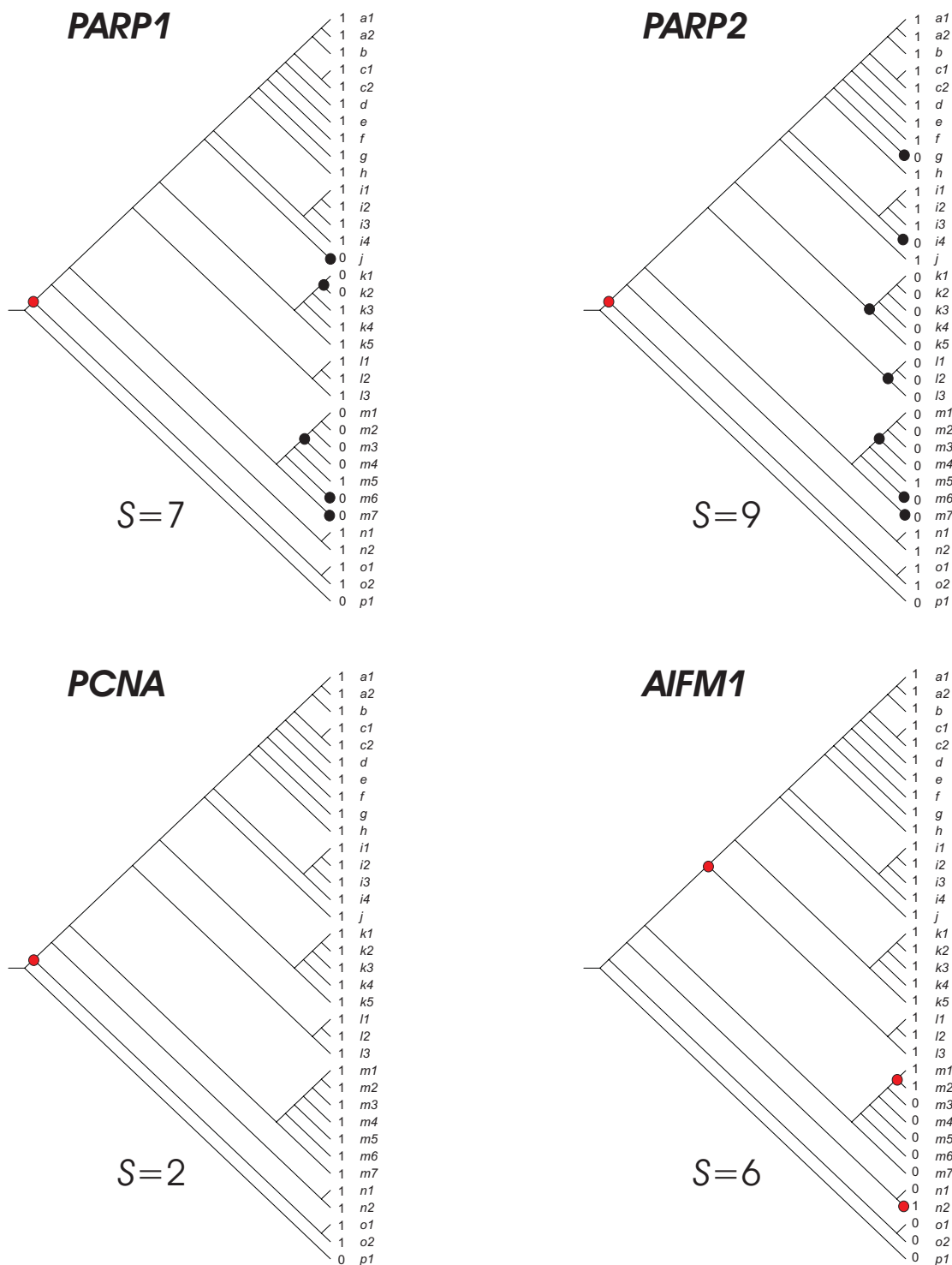
**Supplementary Figure S73.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



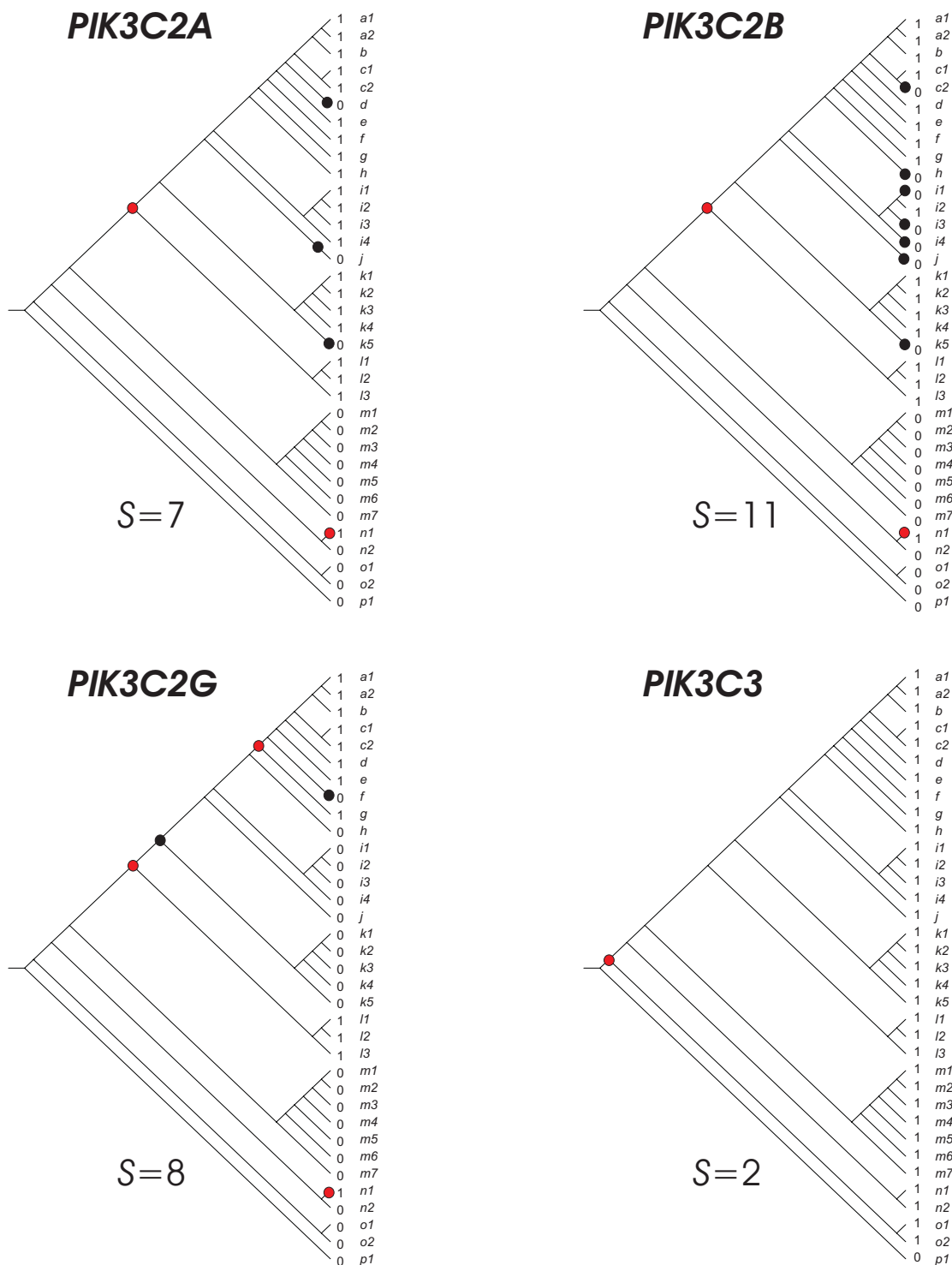
**Supplementary Figure S74.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



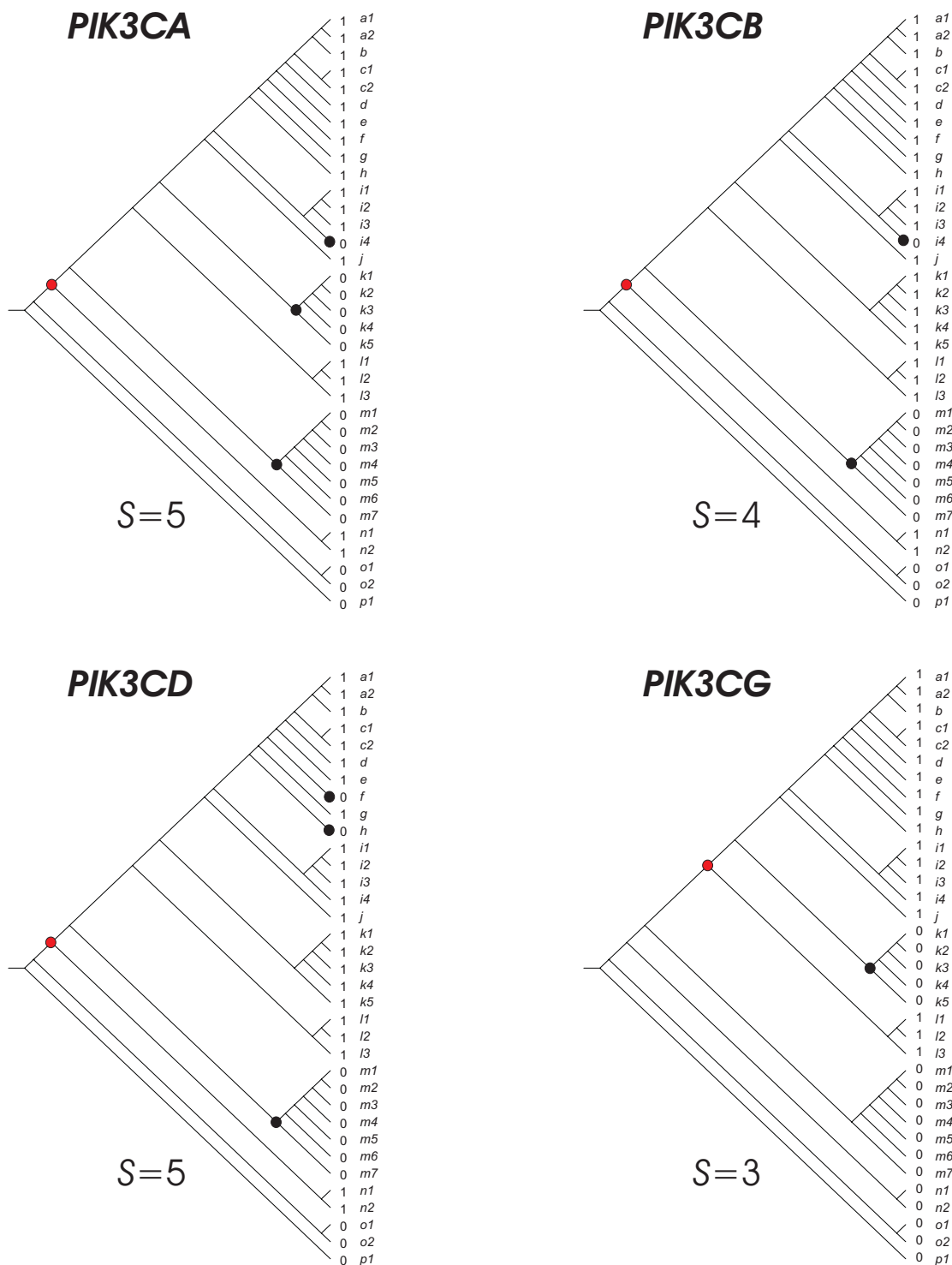
**Supplementary Figure S75.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



**Supplementary Figure S76.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

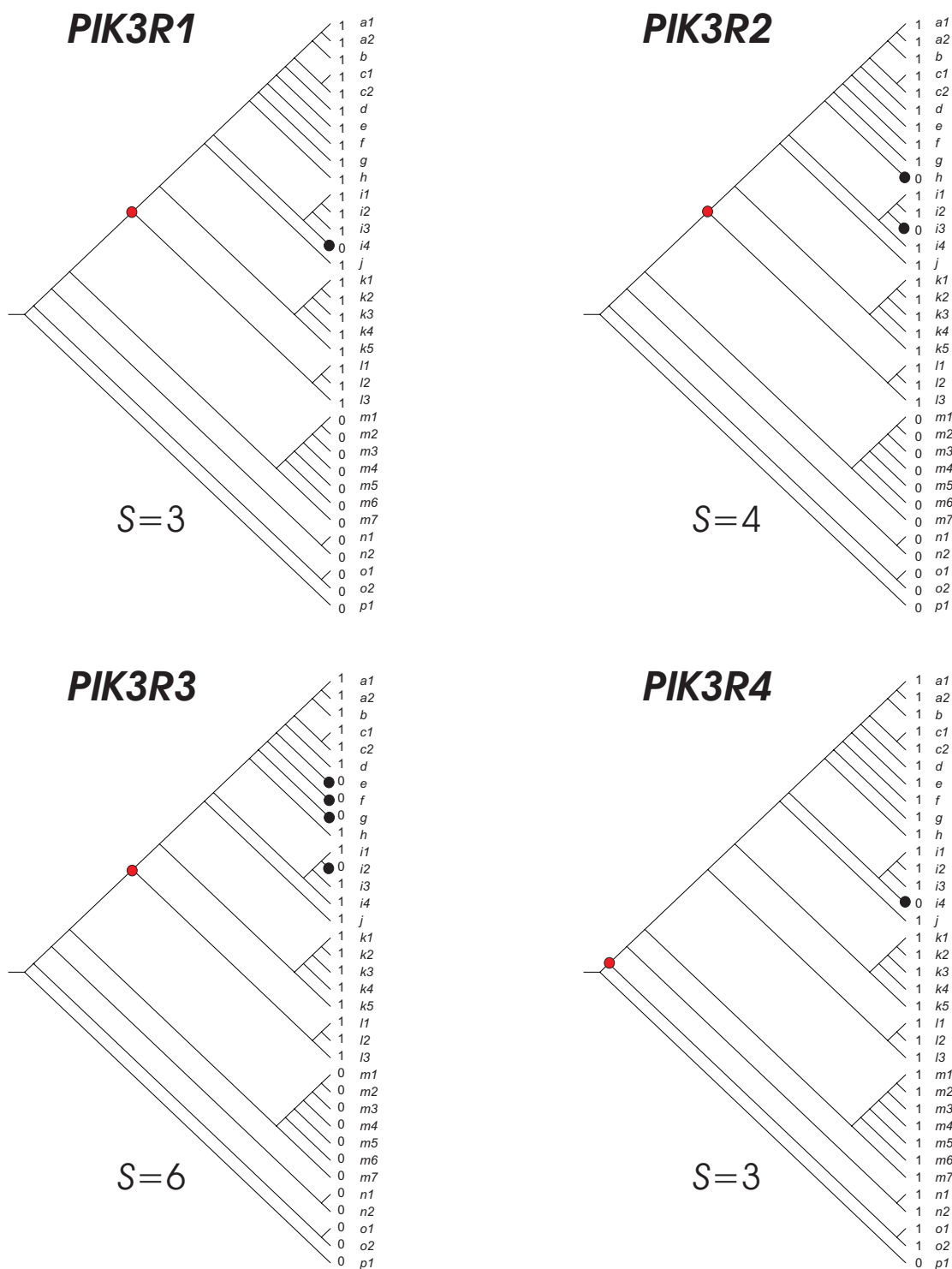


**Supplementary Figure S77.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

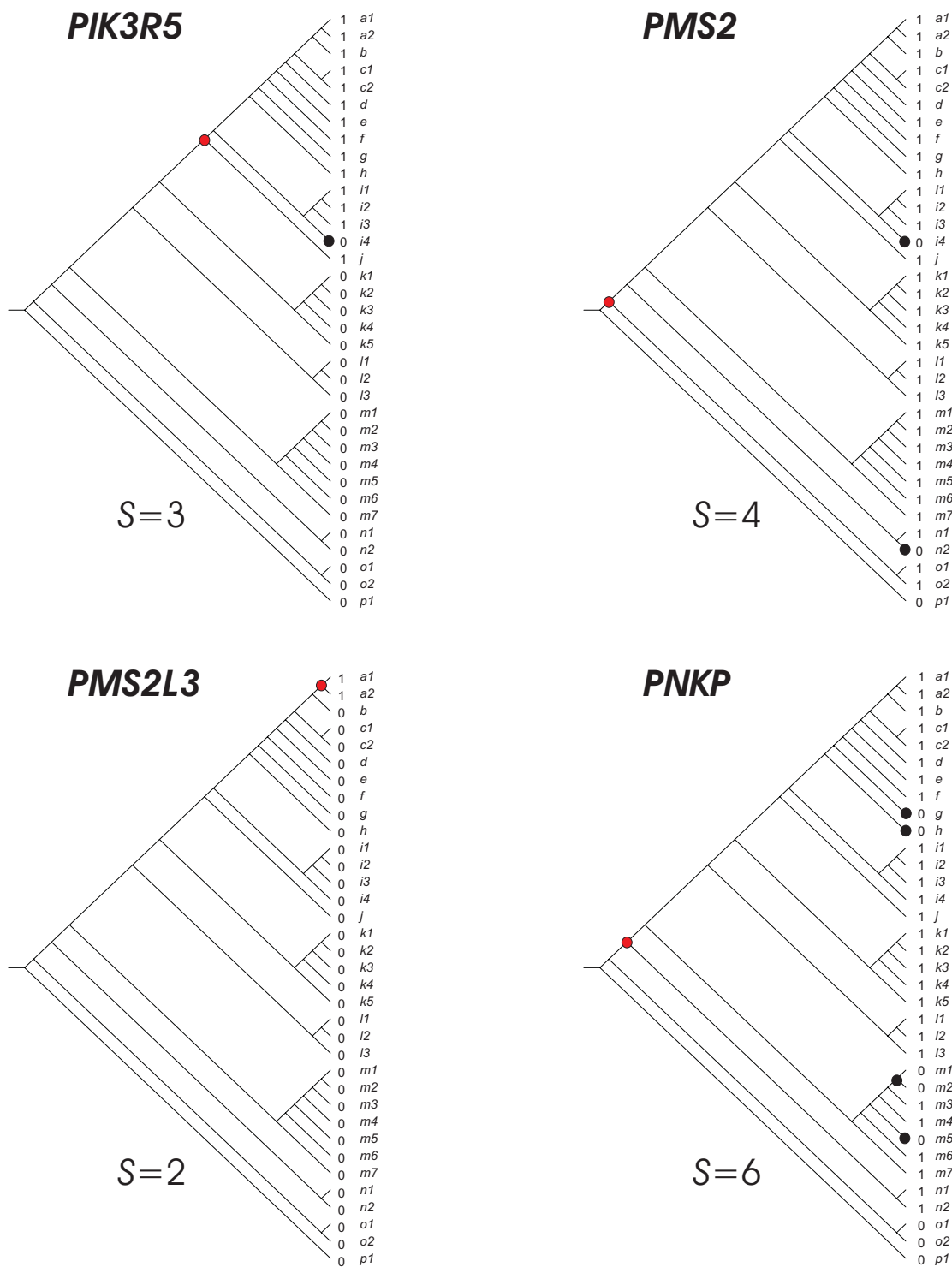


**Supplementary Figure S78.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

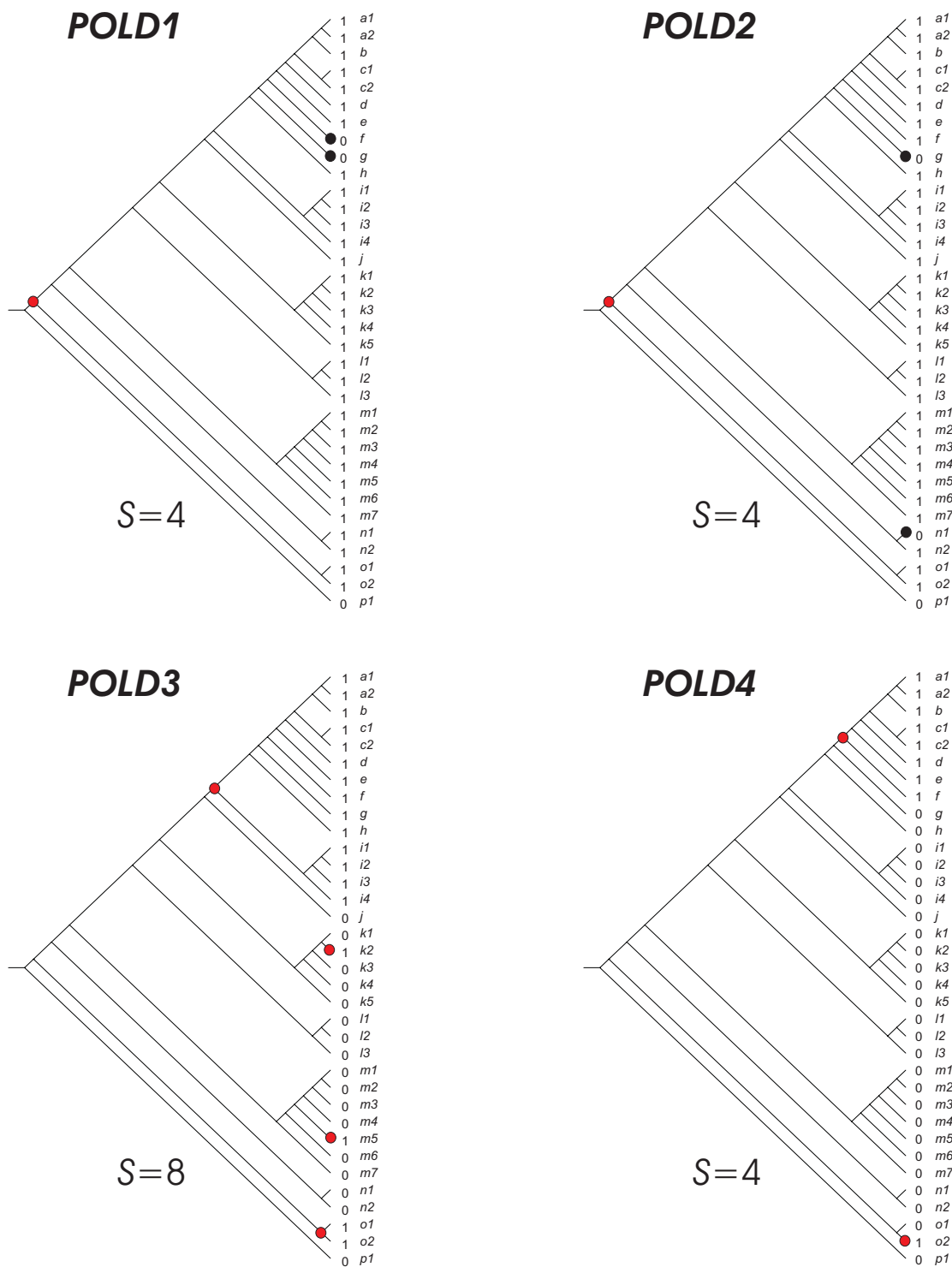




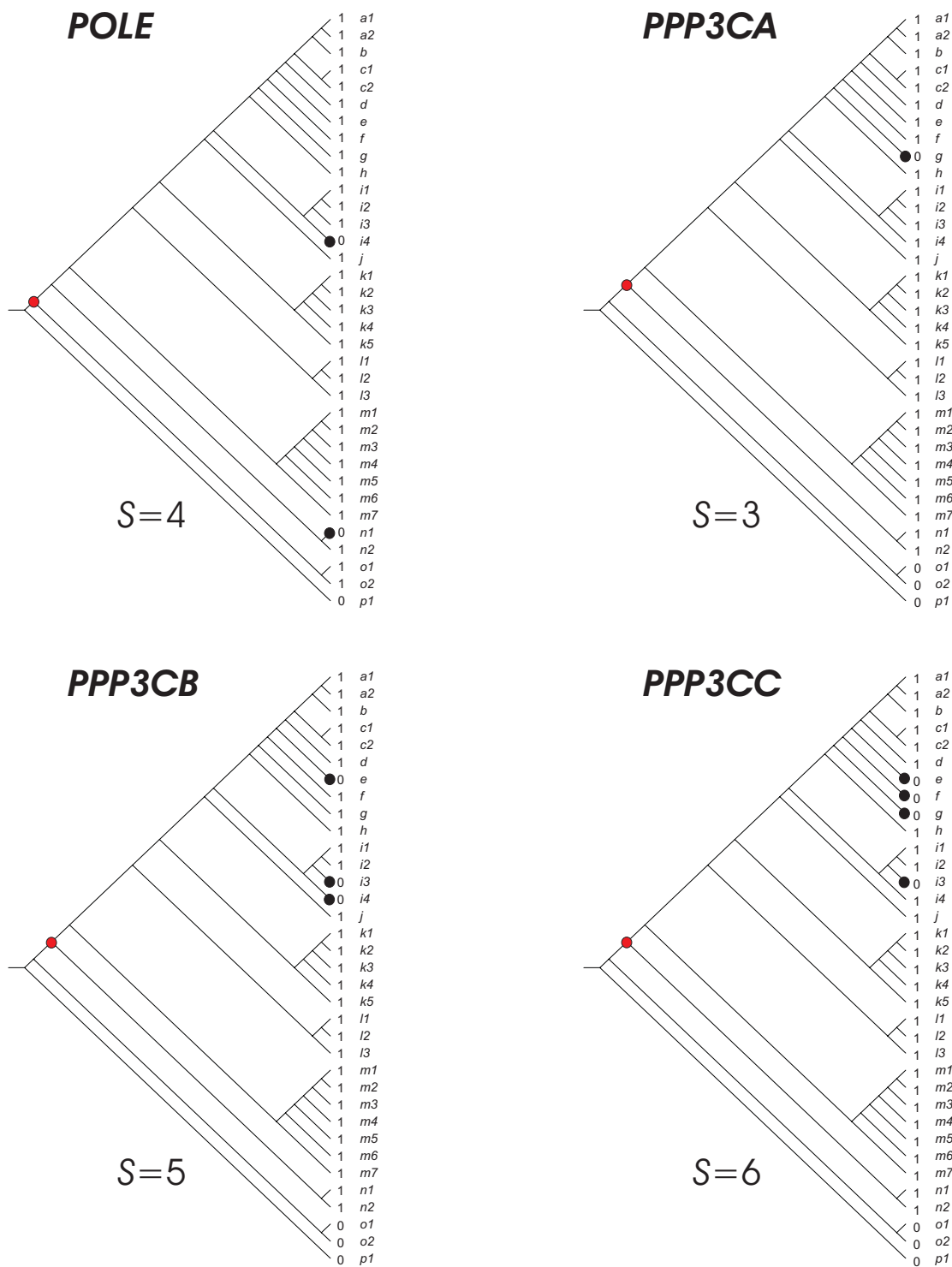
**Supplementary Figure S79.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



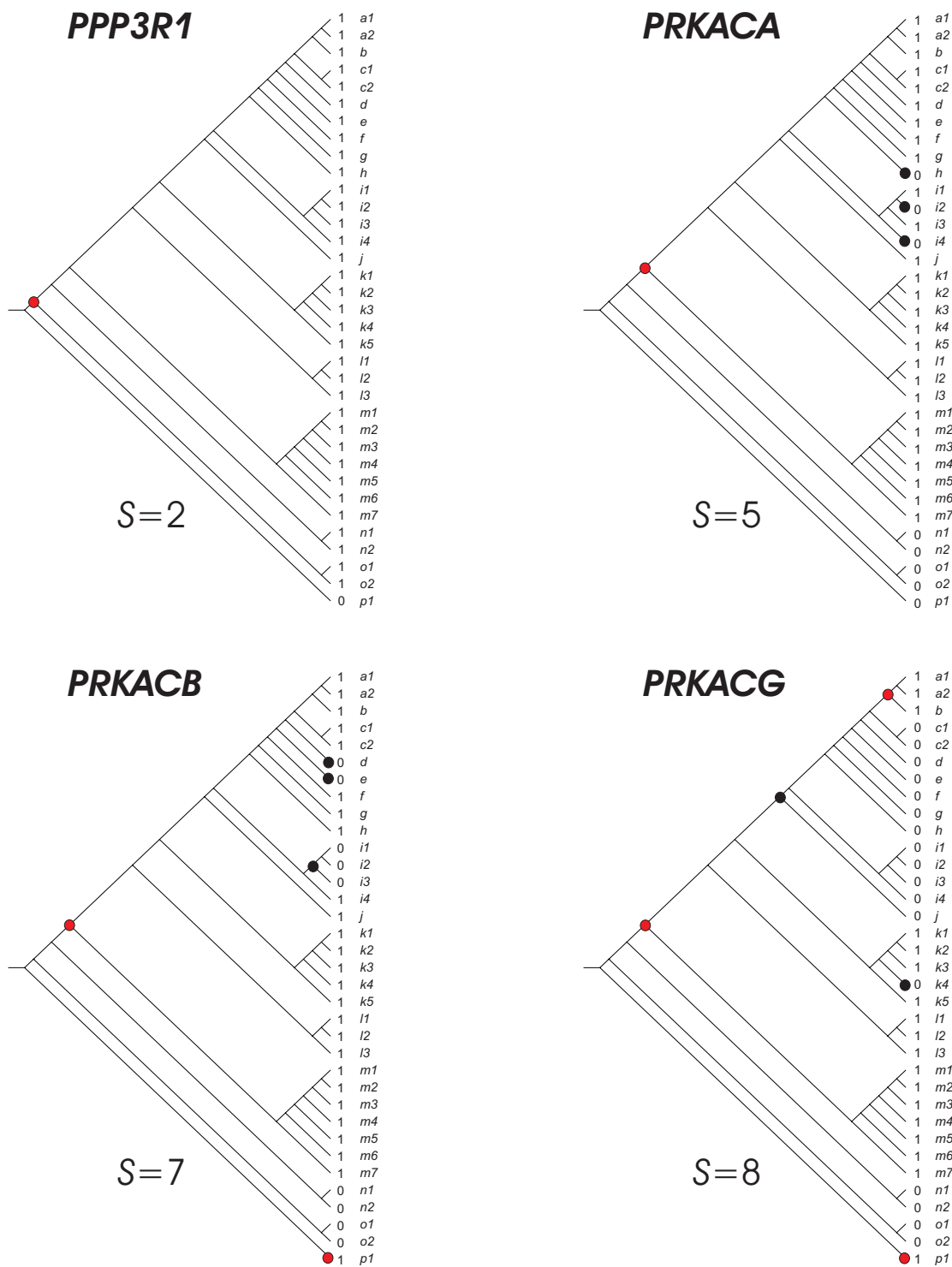
**Supplementary Figure S80.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



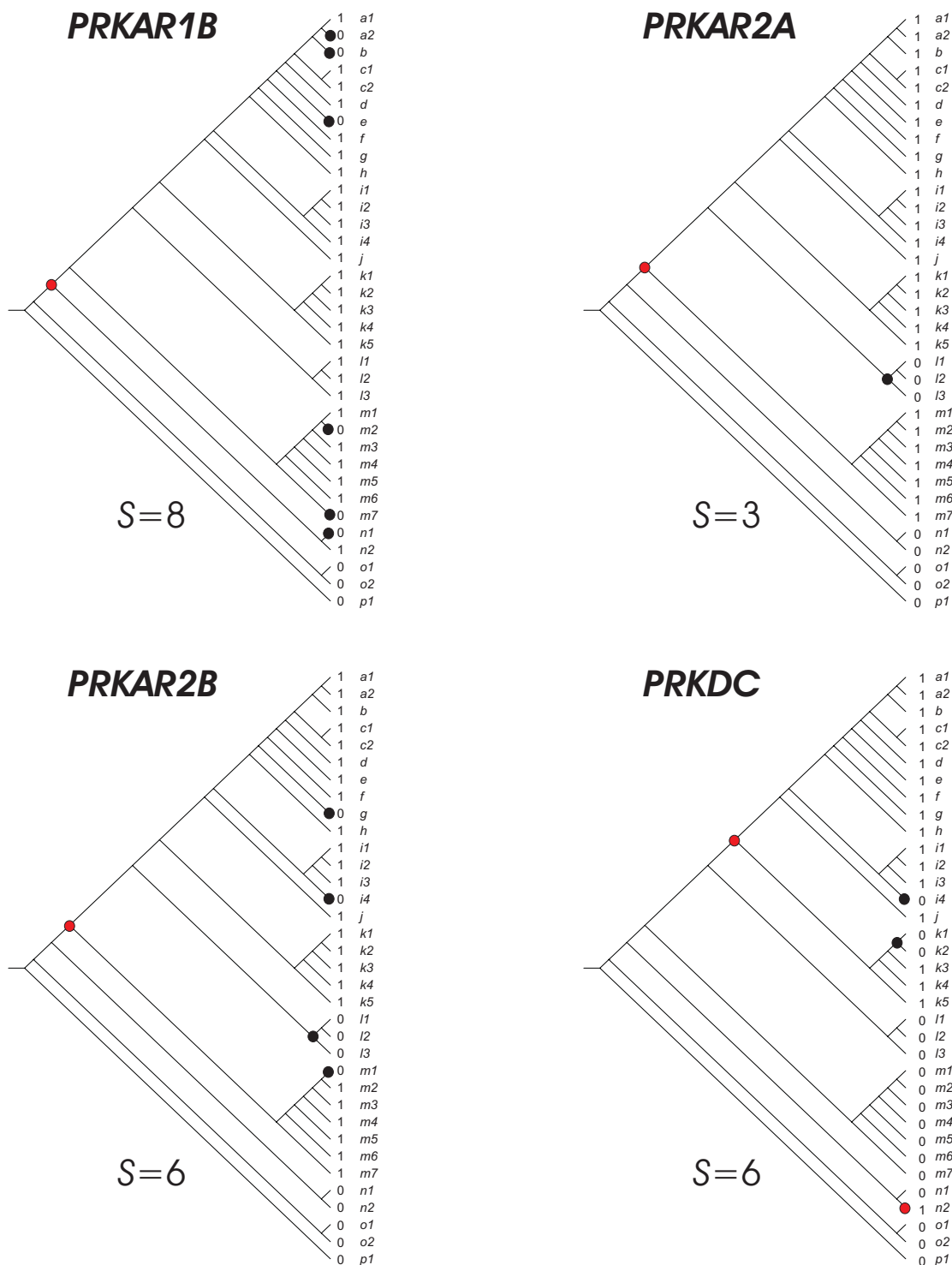
**Supplementary Figure S81.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



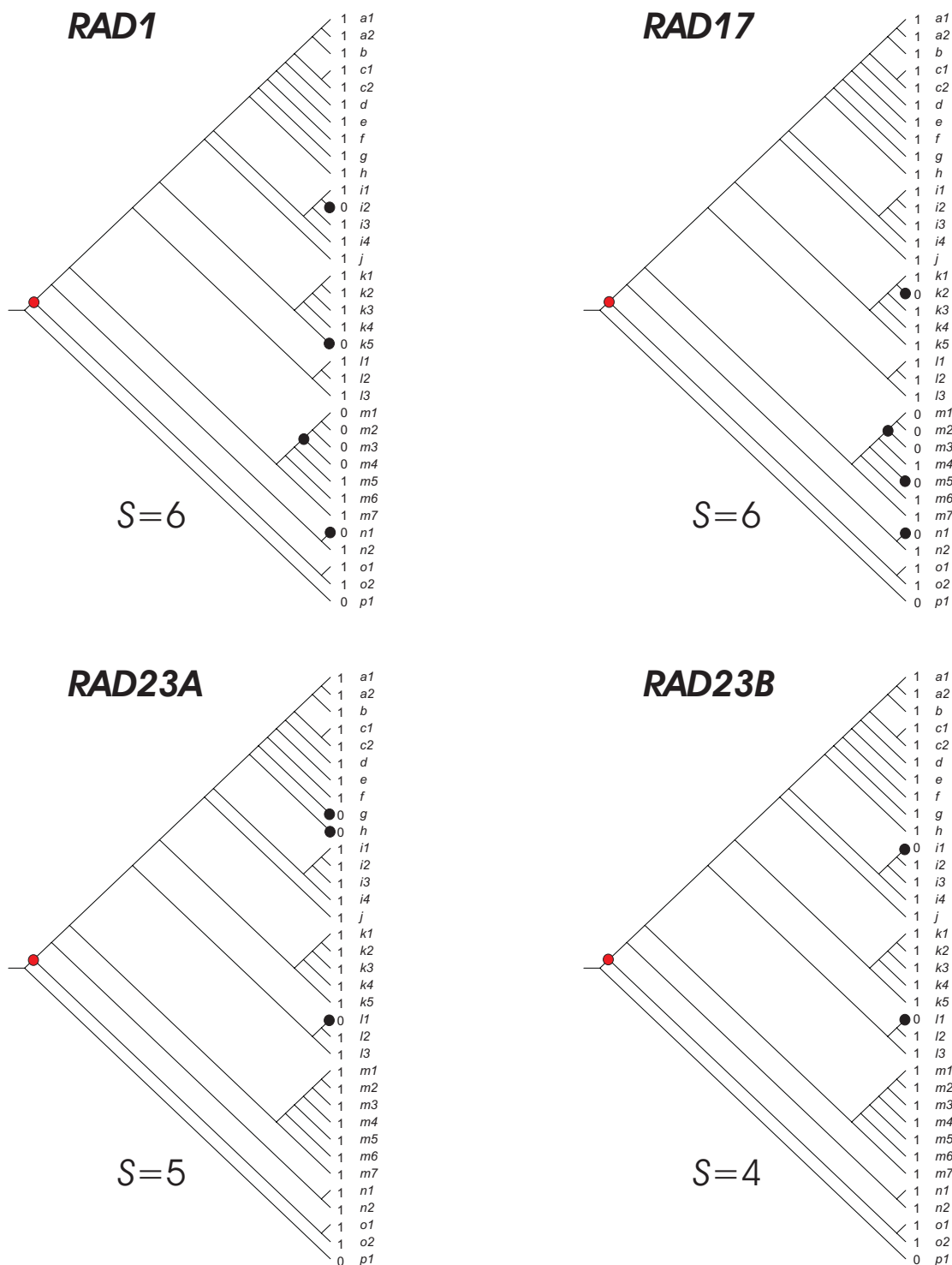
**Supplementary Figure S82.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



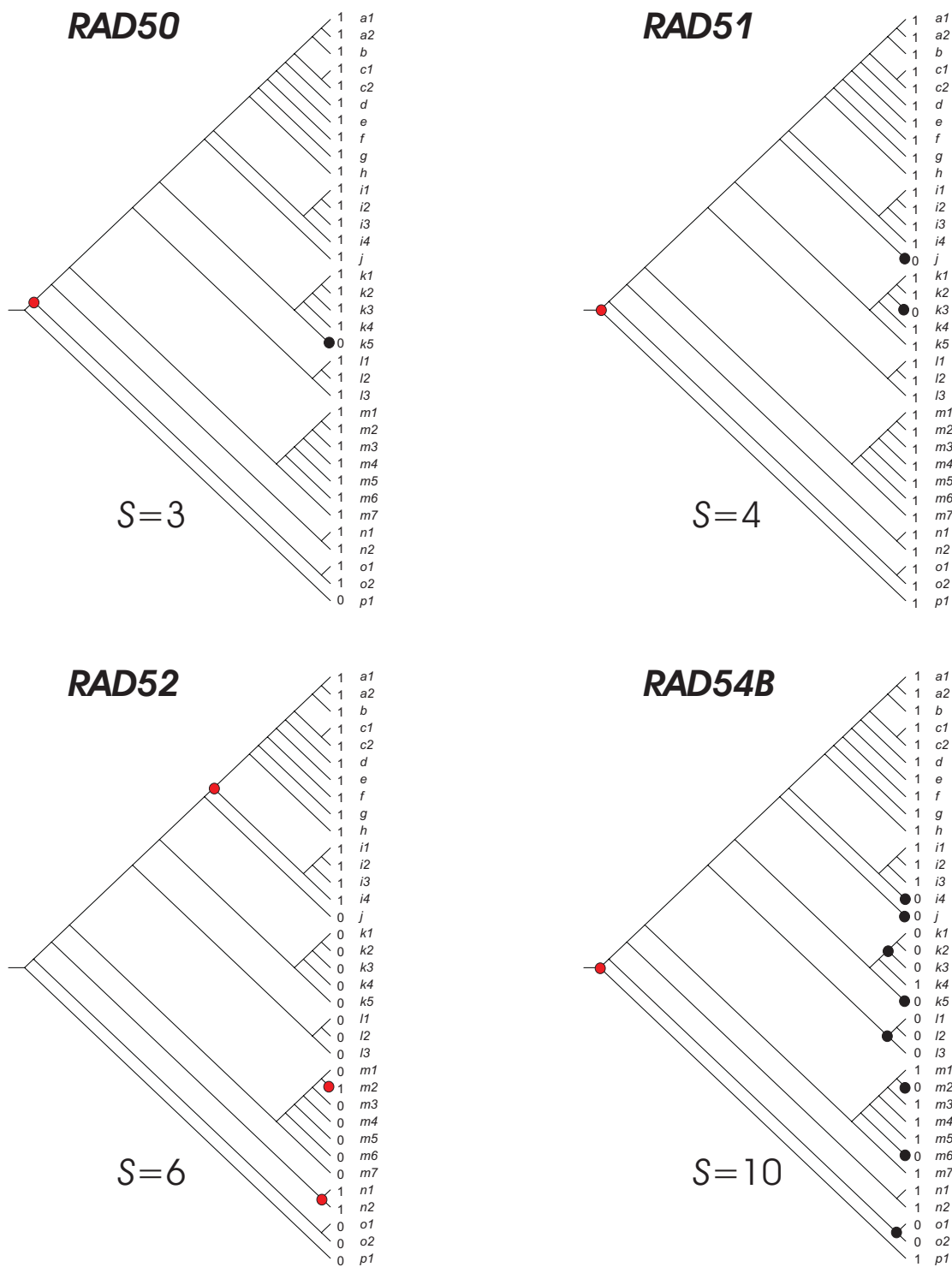
**Supplementary Figure S83.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



**Supplementary Figure S84.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

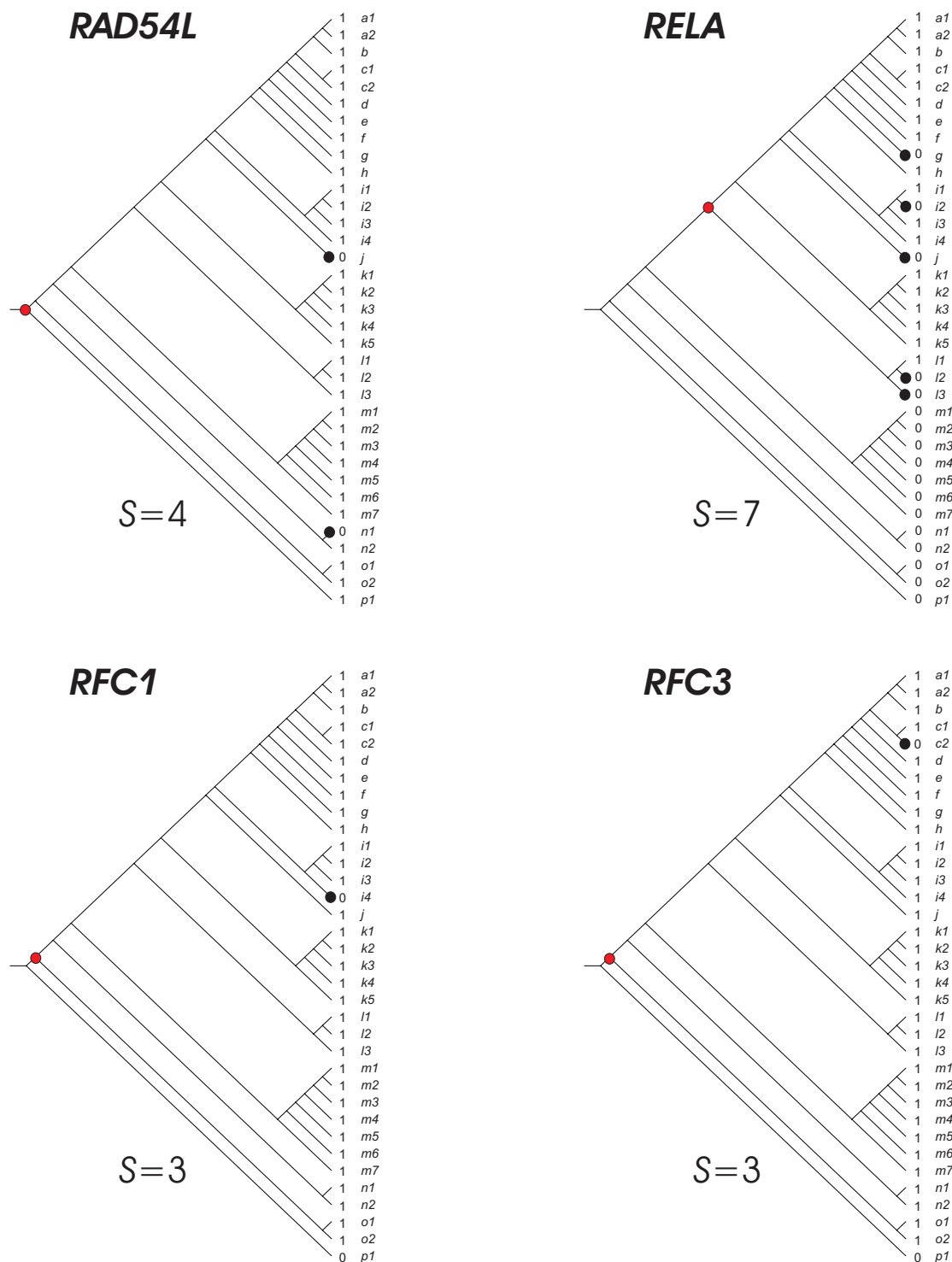


**Supplementary Figure S85.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

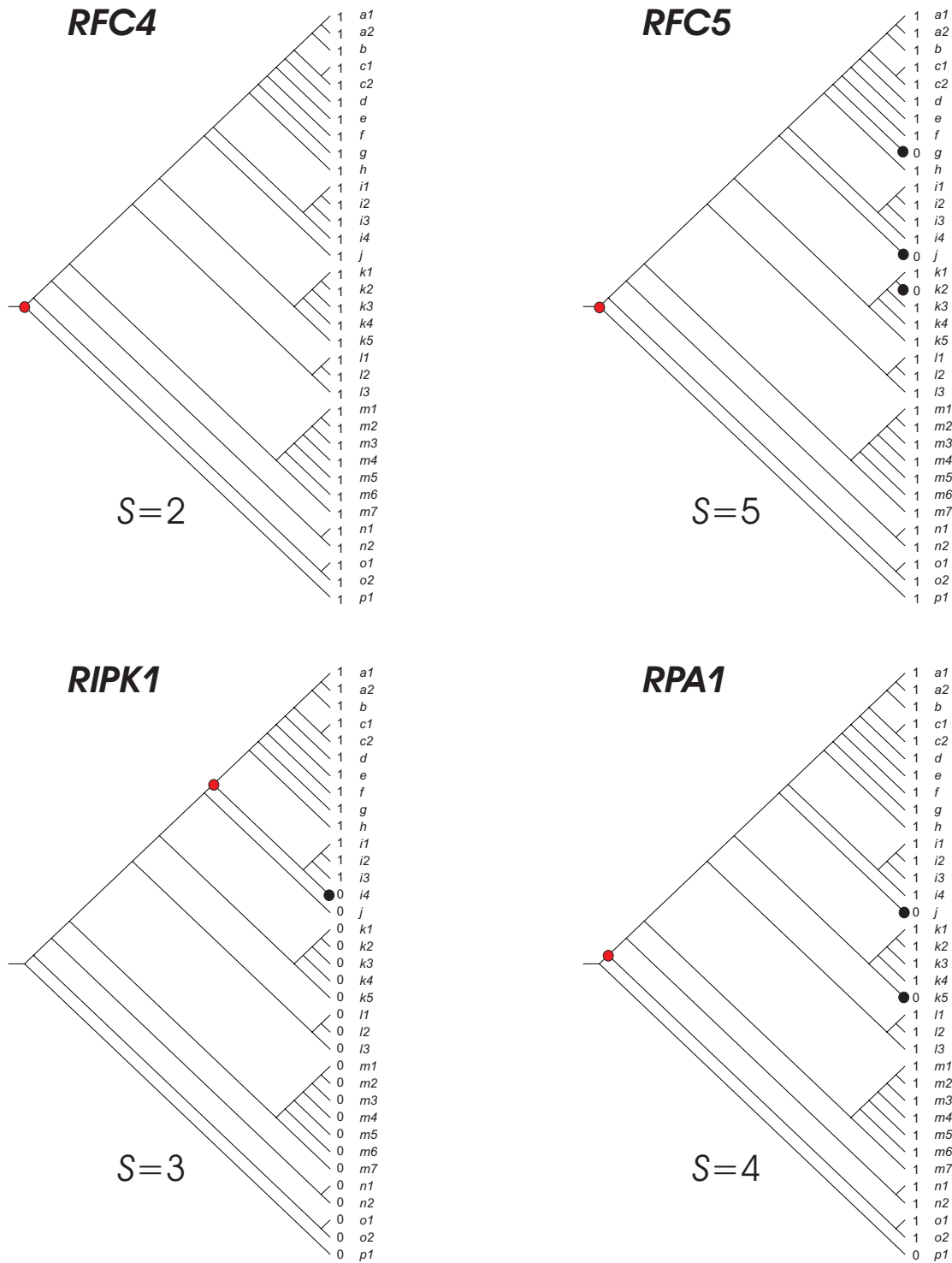


**Supplementary Figure S86.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.

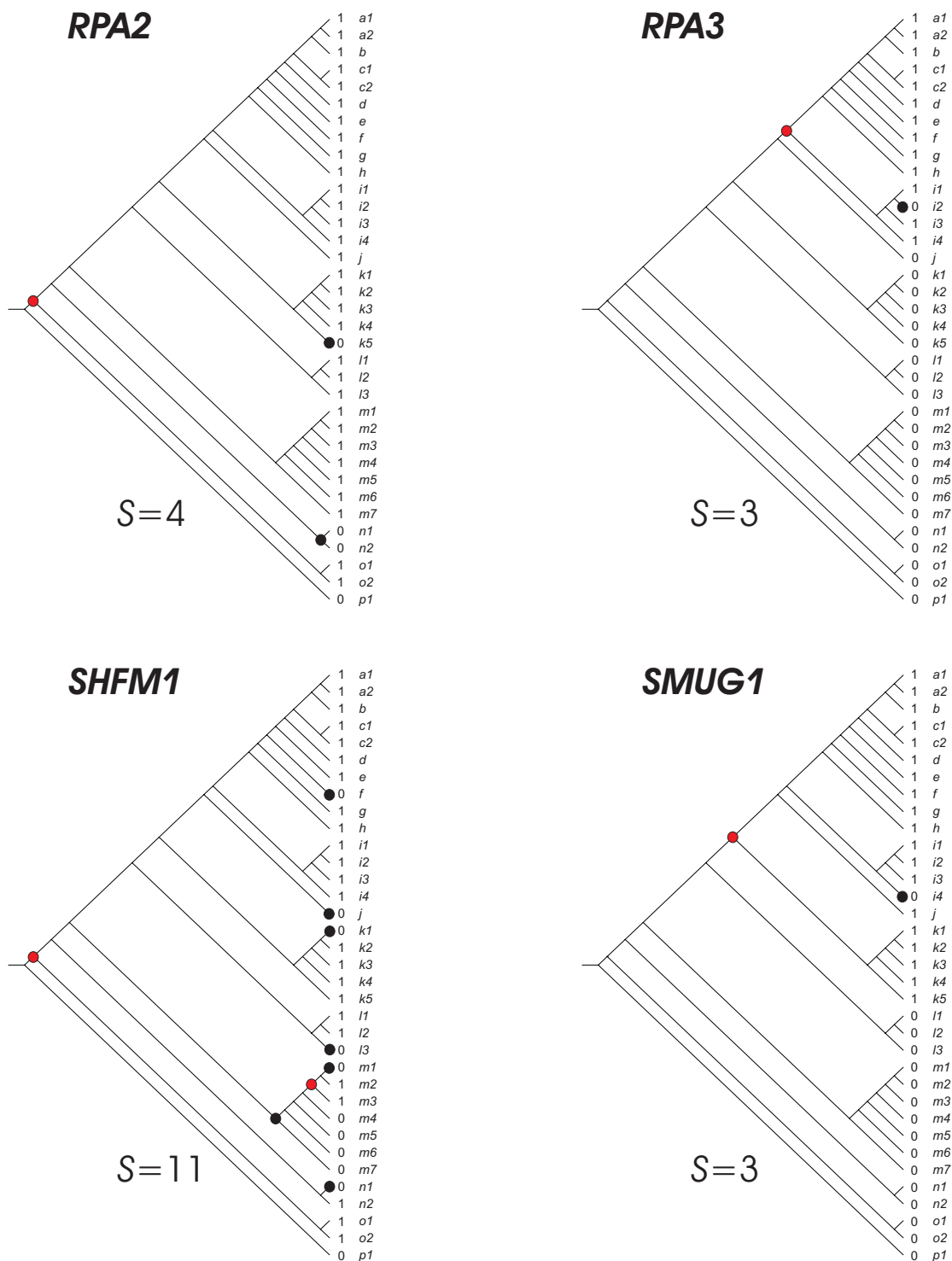




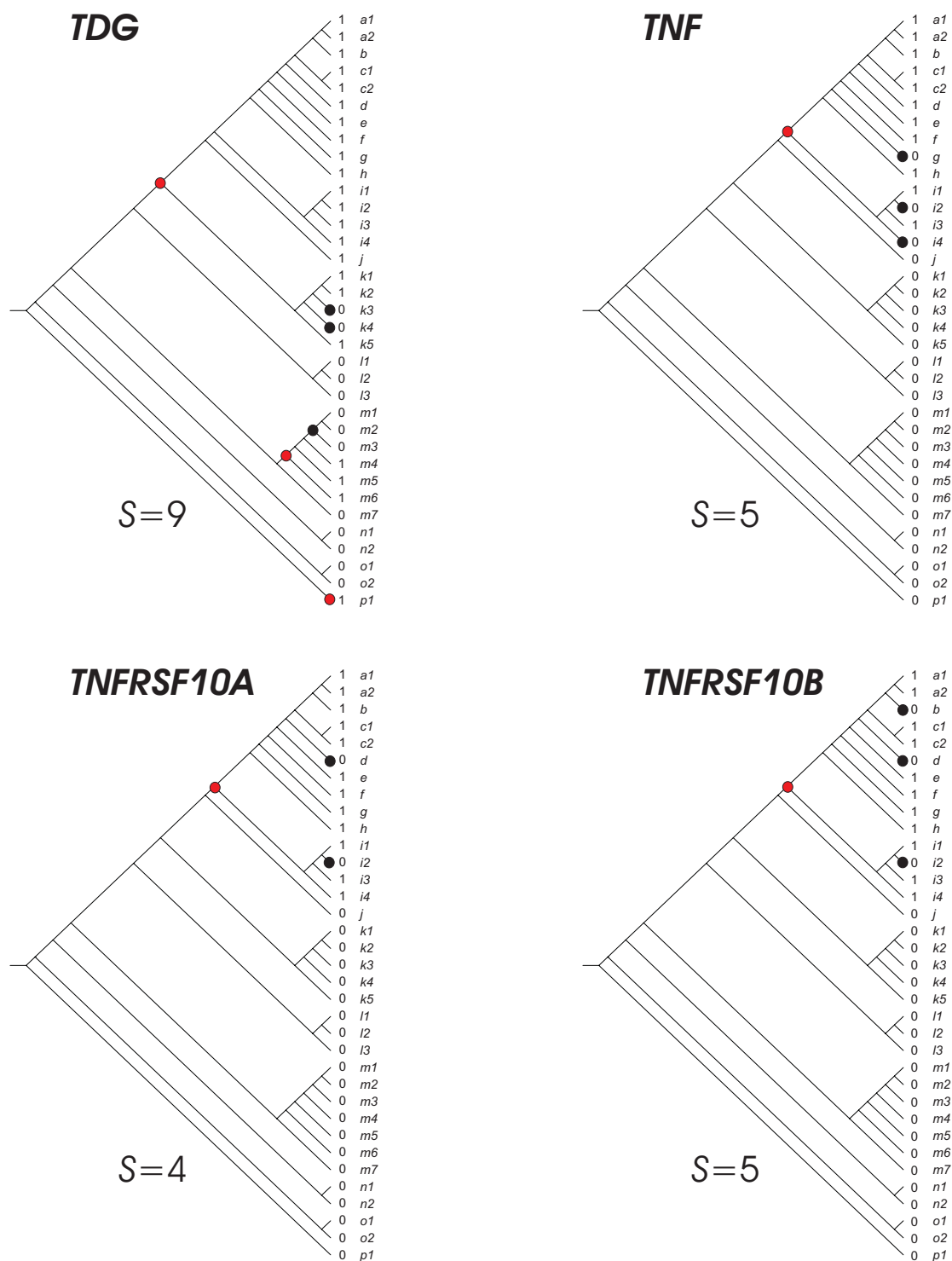
**Supplementary Figure S87.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



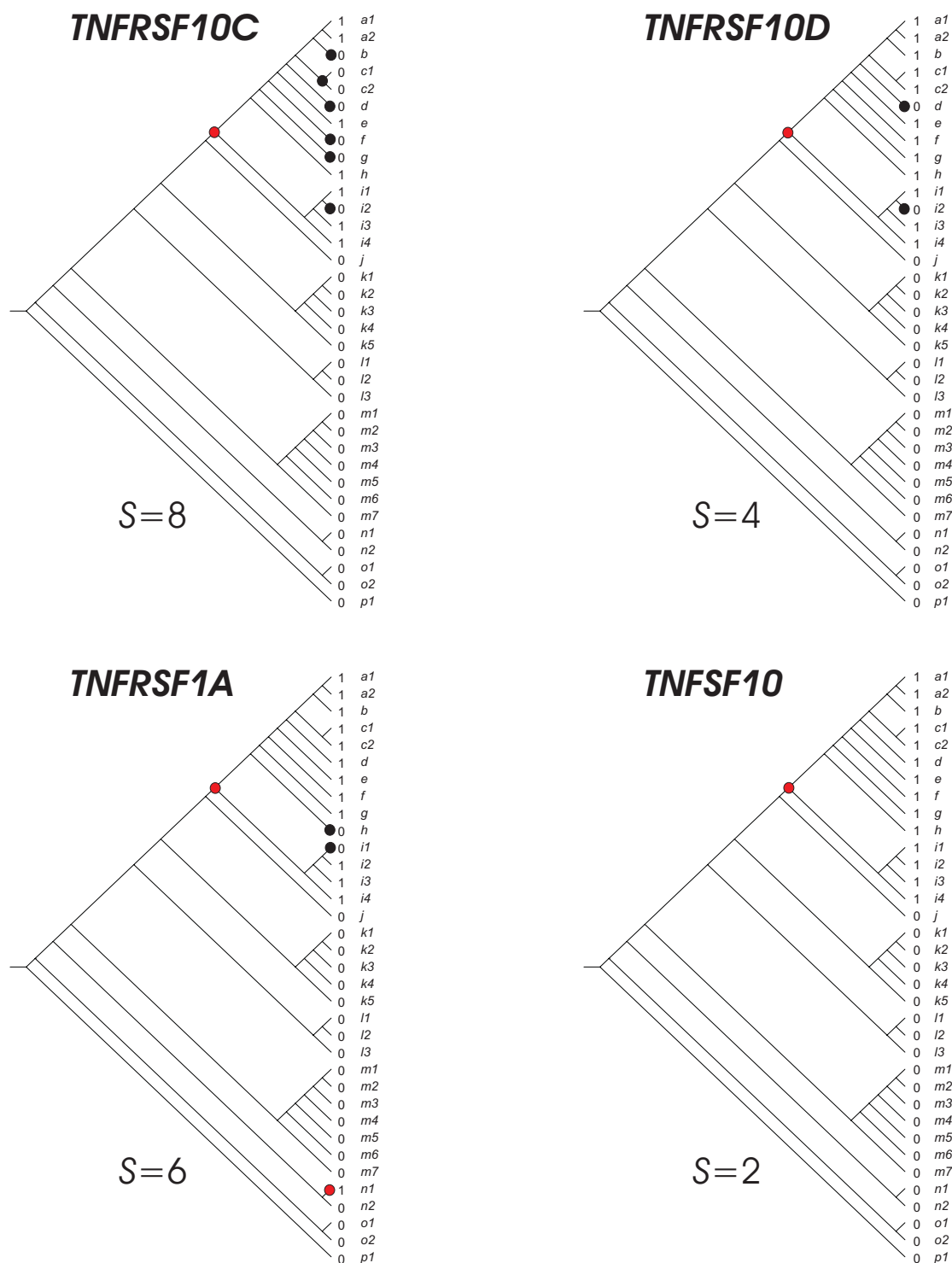
**Supplementary Figure S88.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



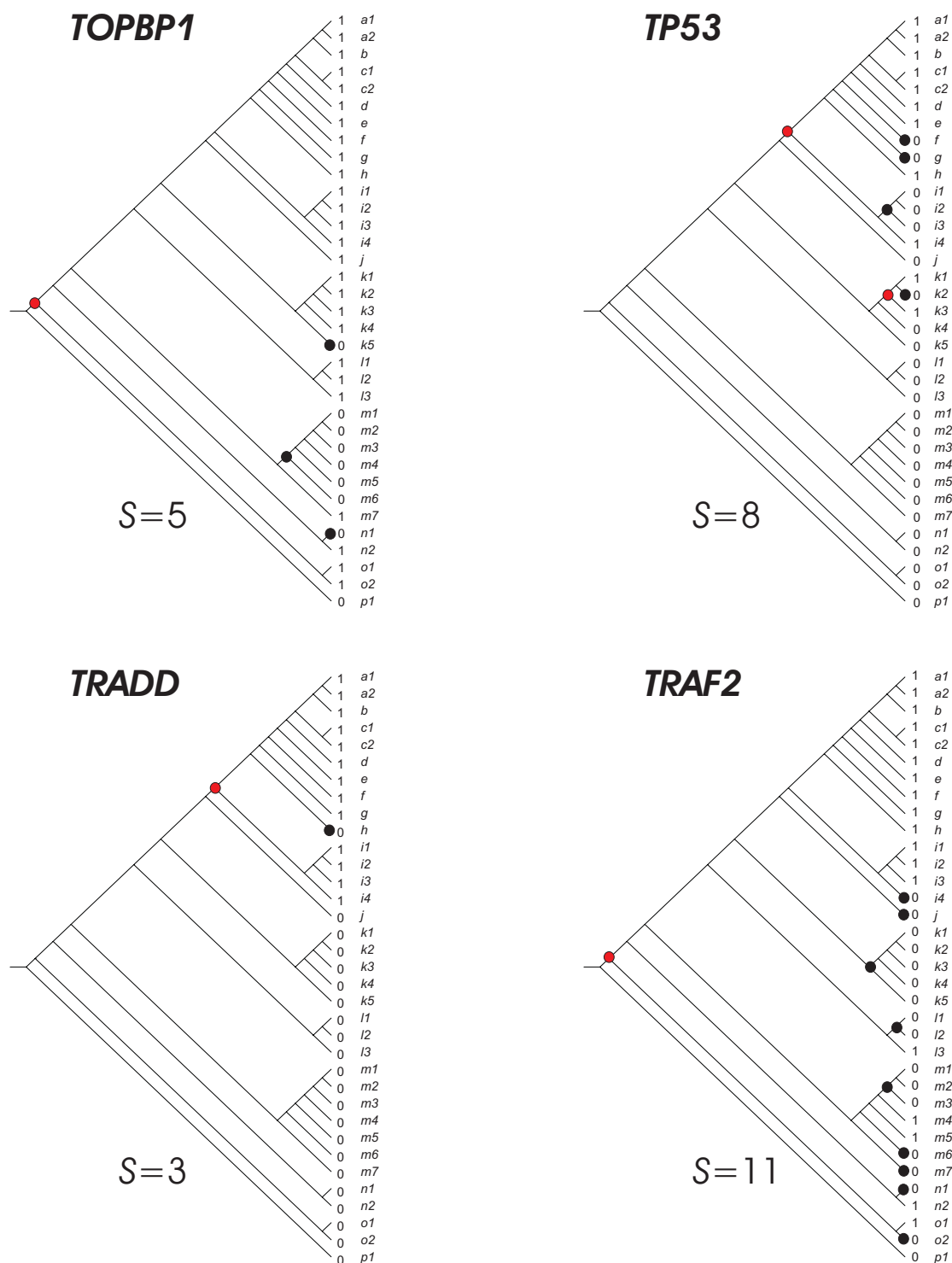
**Supplementary Figure S89.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



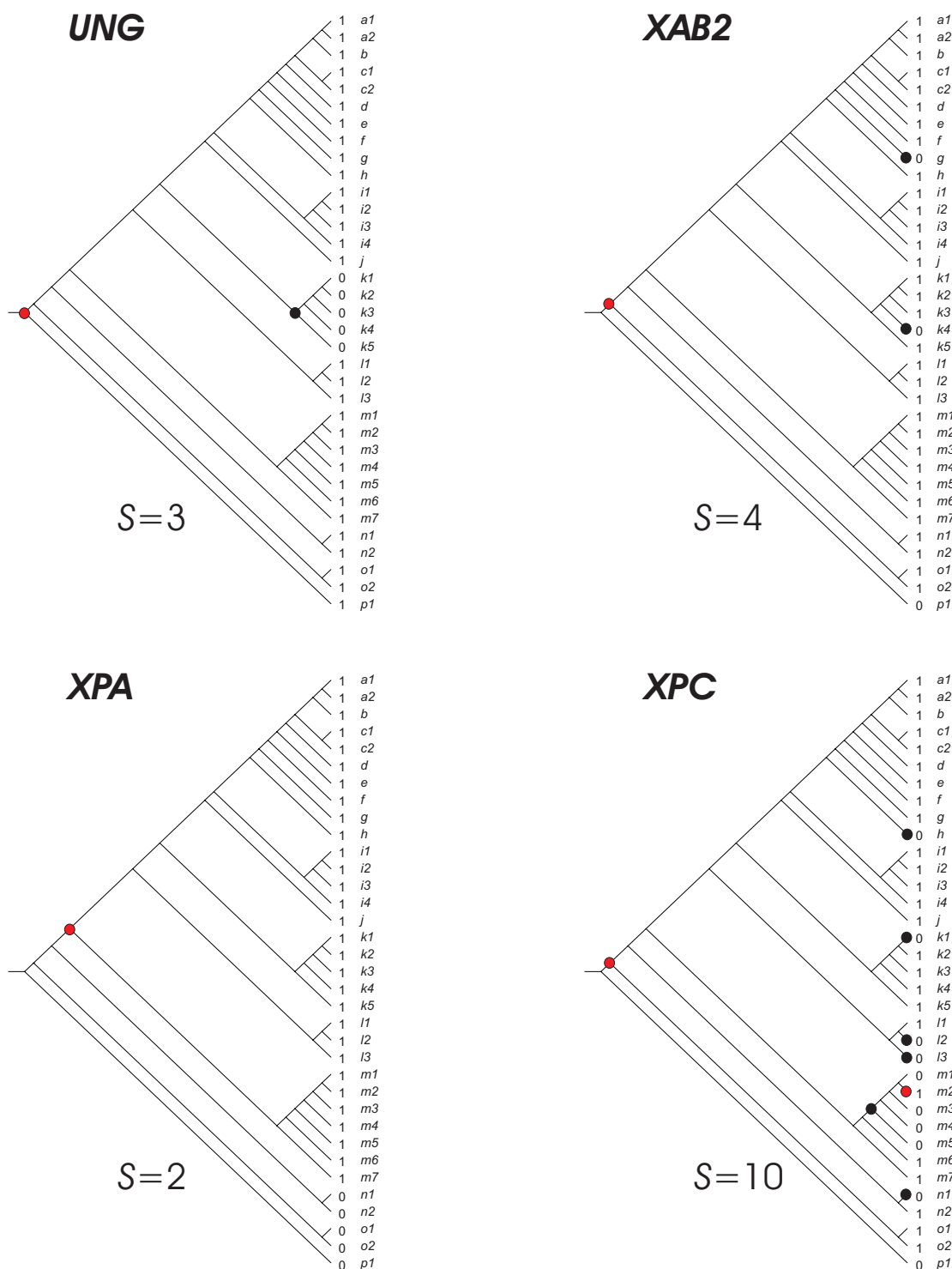
**Supplementary Figure S90.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



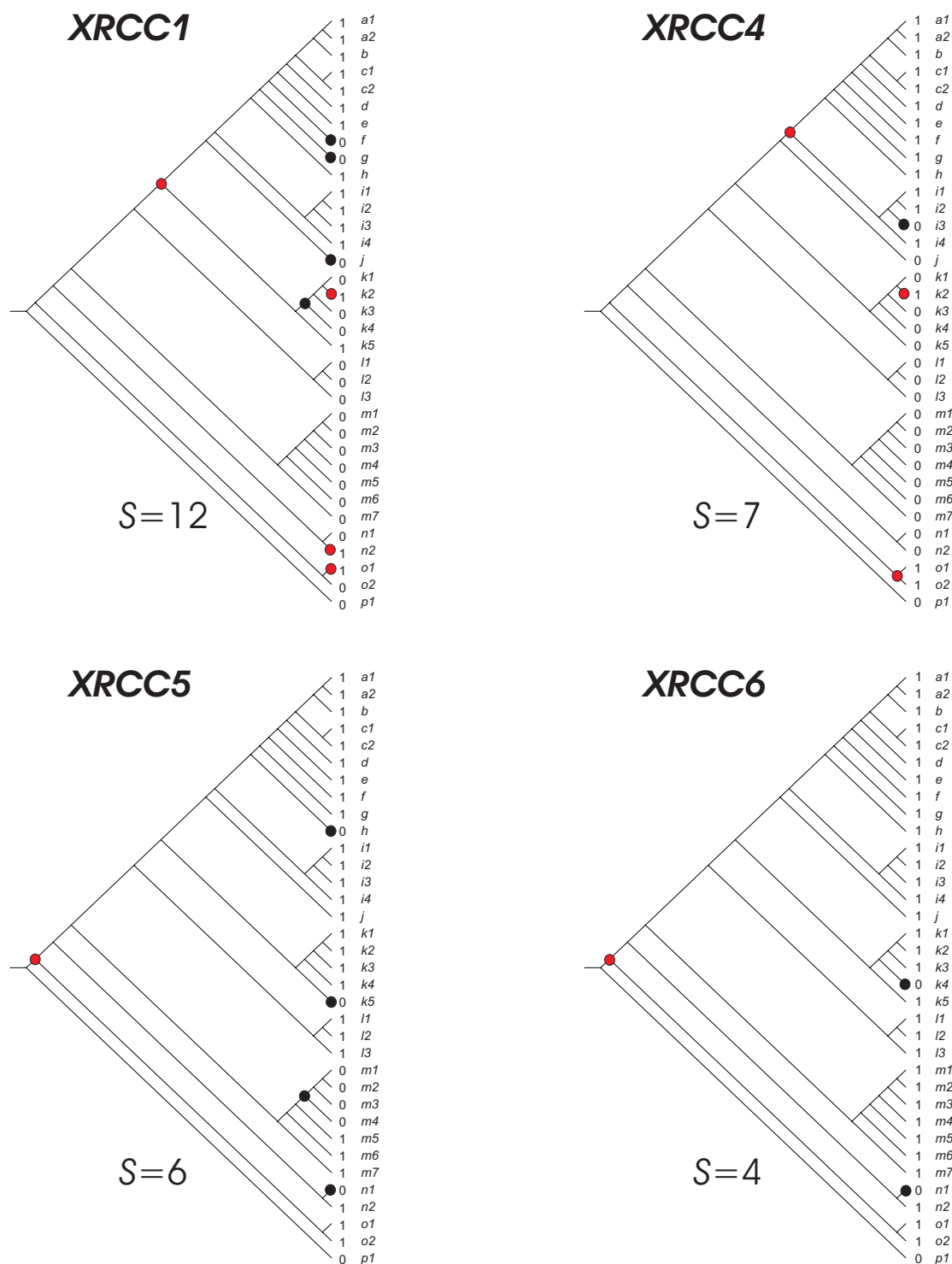
**Supplementary Figure S91.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



**Supplementary Figure S92.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



**Supplementary Figure S93.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



**Supplementary Figure S94.** Parsimony analysis of orthologous groups according to Inparanoid database. The inconsistency value  $S$  is indicated for each evolutionary scenario. The relative costs of the evolutionary events consider two cost units for gene birth or gene acquisition (red nodes), and one cost unit for gene loss (black nodes) as described in *material and methods*. Branch codes as in Supplementary Figure S50.



## 2. References

- Baldauf, S.L. 2003. The deep roots of eukaryotes. *Science* **300**:1703-1706.
- Best, A.A., Morrison, H.G., McArthur, A.G., Sogin, M.L., and Olsen, G.J. 2004. Evolution of eukaryotic transcription: Insights from the genome of *Giardia lamblia*. *Genome Res.* **14**:1537-1547.
- Castro, M.A.A., Mombach, J.C.M., de Almeida, R.M.C., and Moreira, J.C.F. 2007. Impaired expression of NER gene network in sporadic solid tumors. *Nucleic Acids Res.* **35**:1859-1867.
- Ciccarelli, F.D., Doerks, T., von Mering, C., Creevey, C.J., Snel, B., and Bork, P. 2006. Toward automatic reconstruction of a highly resolved tree of life. *Science* **311**:1283-1287.
- Delsuc, F., Brinkmann, H., and Philippe, H. 2005. Phylogenomics and the reconstruction of the tree of life. *Nat. Rev. Genet.* **6**:361-375
- Doolittle, W.F. 1999. Phylogenetic classification and the universal tree. *Science* **284**:2124-2128.
- Eppig, J.T., Blake, J.A., Bult, C.J., Kadin, J.A., Richardson, J.E., and the Mouse Genome Database Group. 2007. The mouse genome database (MGD): new features facilitating a model system. *Nucleic Acids Res.* **35**:D630-D637.
- Futreal, P.A., Coin, L., Marshall, M., Down, T., Hubbard, T., Wooster, R., Rahman, N., and Stratton, M.R. 2004. A census of human cancer genes. *Nat. Rev. Cancer* **4**:177-183.
- Harris, J.K., Kelley, S.T., Spiegelman, G.B., and Pace, N.R. 2003. The Genetic Core of the Universal Ancestor. *Genome Res.* GR-6528.
- Hirschman, J.E., Balakrishnan, R., Christie, K.R., Costanzo, M.C., Dwight, S.S., Engel, S.R., Fisk, D.G., Hong, E.L., Livstone, M.S., Nash, R., et al. 2006. Genome Snapshot: a new resource at the *Saccharomyces* Genome Database (SGD) presenting an overview of the *Saccharomyces cerevisiae* genome. *Nucleic Acids Res.* **34**:D442-D445.
- Katinka, M.D., Duprat, S., Cornillot, E., Metenier, G., Thomarat, F., Prensier, G., Barbe, V., Peyretailade, E., Brottier, P., Wincker, P., et al. 2001. Genome sequence and gene compaction of the eukaryote parasite *Encephalitozoon cuniculi*. *Nature* **414**:450-453
- Kunin, V. and Ouzounis, C.A. 2003. The balance of driving forces during genome evolution in prokaryotes. *Genome Res.* **13**:1589-1594.

- Letunic, I. and Bork, P. 2007. Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and annotation. *Bioinformatics* **23**:127-128
- Makarova, K.S., Wolf, Y.I., Mekhedov, S.L., Mirkin, B.G., and Koonin, E.V. 2005. Ancestral paralogs and pseudoparalogs and their role in the emergence of the eukaryotic cell. *Nucleic Acids Res.* **33**:4626-4638.
- Mirkin, B.G., Fenner, T.I., Galperin, M.Y., and Koonin, E.V. 2003. Algorithms for computing parsimonious evolutionary scenarios for genome evolution, the last universal common ancestor and dominance of horizontal gene transfer in the evolution of prokaryotes. *BMC Evol. Biol.* **3**.
- Pennisi, E. 2003. Drafting a Tree. *Science* **300**:1694
- Remm, M., Storm, C.E.V., and Sonnhammer, E.L.L. 2001. Automatic clustering of orthologs and in-paralogs from pairwise species comparisons. *J. Mol. Biol.* **314**:1041-1052.
- Snel, B., Bork, P., and Huynen, M.A. 2002. Genomes in Flux: The evolution of archaeal and proteobacterial gene content. *Genome Res.* **12**:17-25.
- Tatusov, R., Fedorova, N., Jackson, J., Jacobs, A., Kiryutin, B., Koonin, E., Krylov, D., Mazumder, R., Mekhedov, S., Nikolskaya, A., et al. 2003. The COG database: an updated version includes eukaryotes. *BMC Bioinformatics* **4**:41.