

1       **The evolution, structure, and function of RubisCO and its**  
2                               **homolog the RubisCO-like protein**

3  
4  
5       **F. Robert Tabita<sup>\*,†,1</sup>, Thomas E. Hanson<sup>1</sup>, Sriram Satagopan<sup>†</sup>, Brian H.**  
6       **Witte<sup>†</sup>, and Nathan E. Kreel<sup>1</sup>**

7  
8       **Department of Microbiology<sup>†</sup>, the Plant Molecular Biology Biology/Biotechnology**  
9       **Program<sup>‡</sup>, and the OSU Biochemistry Program<sup>1</sup>**  
10       **The Ohio State University**  
11       **484 West 12th Avenue**  
12       **Columbus, Ohio 43210-1292;**  
13       **Graduate College of Marine and Earth Studies<sup>1</sup>, Delaware Biotechnology Institute,**  
14       **University of Delaware, 127 DBI, 15 Innovation Way, Newark, DE 19711**

15  
16  
17  
18       **\*To whom correspondence should be addressed at:**  
19       **Department of Microbiology**  
20       **The Ohio State University**  
21       **484 West 12<sup>th</sup> Avenue**  
22       **Columbus, Ohio 43210-1292**  
23       **Telephone: 614-292-4297**  
24       **Fax: 614-292-6337**  
25       **E-mail: [Tabita.1@osu.edu](mailto:Tabita.1@osu.edu)**

26  
27  
28  
29       **Running Title: RubisCO and RLP**

30  
31  
32       **Key words: RubisCO, structure/function, evolution, different forms, CBB cycle**

33

34 **Abstract**

35

36 Ribulose 1,5-bisphosphate (RuBP) carboxylase/oxygenase (RubisCO) catalyzes the key  
37 reaction by which inorganic carbon may be assimilated into organic carbon. Phylogenetic  
38 analyses indicate that there are three classes of bona fide RubisCO proteins, forms I, II, and III,  
39 that all catalyze the same reactions. In addition, there exists another form of RubisCO, form IV,  
40 which does not catalyze typical RubisCO reactions. Form IV is actually a homolog of RubisCO  
41 and is called the RubisCO-like protein (RLP). Both RubisCO and RLP appear to have evolved  
42 from a methanogenic archaeal ancestor protein and comprehensive analyses indicate that the  
43 various forms (I, II, III, and IV) contain various subgroups, with individual sequences derived  
44 from representatives of all three kingdoms of life. The diversity of RubisCO molecules, many of  
45 which function in distinct milieus, have provided convenient model systems to study the ways in  
46 which the active site of this protein has evolved to accommodate necessary molecular  
47 adaptations. Such studies have proven useful to help provide a framework for understanding the  
48 molecular basis for many important aspects of RubisCO catalysis, including the elucidation of  
49 factors or functional groups that impinge on RubisCO carbon dioxide/oxygen substrate  
50 discrimination.

## 51 Introduction

52 Ribulose 1,5-bisphosphate (RuBP) carboxylase/oxygenase (RubisCO) is a protein that  
53 arguably catalyzes the most important biochemical reaction in biology, as it is responsible for  
54 the vast majority of all the organic carbon found in the biosphere. The enzyme, found in  
55 phototrophic and chemoautotrophic organisms, basically functions to catalyze the removal and  
56 sequestration of carbon dioxide from the environment by reducing this oxidized gas to the level  
57 of organic carbon, in the process providing the organic building blocks needed to sustain life.  
58 However, RubisCO has an innate problem in that the enediolate intermediate of RuBP that acts  
59 as CO<sub>2</sub> acceptor may also be attacked by other ligands, including molecular oxygen, such that  
60 CO<sub>2</sub> and O<sub>2</sub> compete for the enediolate intermediate at the same active site. The promiscuity of  
61 the enediol intermediate thus enables RubisCO to function as an internal monooxygenase, as  
62 well as a carboxylase, with the unique product, 2-phosphoglycolate (2-PG), formed as a result  
63 of O<sub>2</sub> fixation (Fig. 1). Efficient RubisCO catalysis is thus dependent on the inherent ability of  
64 the enzyme to discriminate between CO<sub>2</sub> and O<sub>2</sub> (the substrate specificity or  $\Omega$  or  $\tau$  value) at the  
65 relative concentration of CO<sub>2</sub> and O<sub>2</sub> employed in a particular reaction. The rates of the two  
66 reactions ( $v_c$  and  $v_o$ ) may be defined by  $v_c/v_o = \Omega [CO_2]/[O_2]$ . Thus,  $\Omega = v_c [O_2]/v_o [CO_2]$  and  $\Omega =$   
67  $V_c K_o/V_o K_c$  with  $V_c$  and  $V_o$  representing maximum velocities for carboxylation and oxygenation,  
68 respectively, and  $K_c$  and  $K_o$  the relative Michaelis constants for CO<sub>2</sub> and O<sub>2</sub>, respectively.

69 Because of the great disparity in ambient O<sub>2</sub> concentrations relative to CO<sub>2</sub> levels,  
70 aerobic organisms produce considerable amounts of 2-PG. The resultant oxidative metabolism  
71 of 2-PG is inimical to maximal CO<sub>2</sub> fixation (and biomass productivity) because up to 50 % of  
72 the CO<sub>2</sub> that is fixed via RubisCO may be released as a result of 2-PG metabolism (Fig. 1).  
73 Clearly, a major key to solving the CO<sub>2</sub> sequestration issue and perhaps improving  
74 photosynthetic crop yields is to somehow improve RubisCO's ability to discriminate between  
75 CO<sub>2</sub> and O<sub>2</sub> (Long et al. 2006a,b). In this article we review likely scenarios as to how RubisCO

76 may have evolved from primordial ancestors, how the active site from different forms of  
77 RubisCO became adapted to diverse intracellular milieus, and how we might take advantage of  
78 these findings.

79

## 80 **Different molecular forms of RubisCO for the same and different functions**

81         Sixty years after the discovery of Fraction One Protein (i.e., RubisCO) from plant leaf  
82 extracts, studies with prokaryotic CO<sub>2</sub> assimilatory organisms have made it abundantly clear  
83 that nature has evolved different structural forms of this enzyme (Tabita 1999; 2007). In most  
84 instances RubisCO, as expected, catalyzes the typical carboxylation reaction so necessary for  
85 CO<sub>2</sub> reduction via the Calvin-Benson-Bassham (CBB) reductive pentose phosphate pathway  
86 (Fig. 1). However, as the different molecular forms were discovered and analyzed, it also soon  
87 became apparent that RubisCO may be used for purposes other than for primary carbon  
88 assimilation. Indeed, in certain archaea, RubisCO plays primarily an anaplerotic role in providing  
89 a way in which intermediates of purine/pyrimidine metabolism may be salvaged (reviewed in  
90 Tabita et al. 2007). Indeed, representative archaea have developed a novel means to  
91 synthesize RuBP (Finn & Tabita 2004; Sato et al. 2007) and require RubisCO to remove this  
92 metabolite, which has no known metabolic fate beyond serving as a substrate for RubisCO.  
93 Moreover, in some organisms, a type of RubisCO analog has evolved which does not even  
94 catalyze RuBP-dependent CO<sub>2</sub> or O<sub>2</sub> fixation.

95         What are these different forms of RubisCO and in what type organisms are they found?

96         What distinguishes these enzymes? On the basis of amino acid sequences, four known forms of  
97 RubisCO, forms I, II, III, and IV, are found in nature. A summary of the different forms of  
98 RubisCO is described in Table 1. X-ray structures are available for representatives of each of  
99 these proteins, so there is considerable baseline information relative to how their function is  
100 related to their structure (Tabita et al. 2007). The most abundant form of RubisCO is the form I  
101 protein. This is the classic high molecular weight protein originally found in plants and once

102 called Fraction One Protein. Form I RubisCO is comprised of large, approximately 50 KDa  
103 (catalytic) subunits encoded by either the *rbcL* or *cbbL* genes. In addition, the form I  
104 holoenzyme is comprised of small, approximately 15 Kda polypeptides encoded by the *rbcS* or  
105 *cbbS* genes. [Note, in proteobacteria the *cbbL* and *cbbS* genes are often found associated with  
106 other structural genes of the CBB pathway (*cbb* genes) in discrete operons; thus these genes all  
107 have the same three letter (*cbb*) prefix]. Eight large subunits are arranged as four dimeric pairs  
108 to form a catalytic core that is decorated on the top and bottom by four small subunits to form  
109 basically an  $L_8S_8$  structure. Form I is widespread amongst higher plants, eukaryotic algae, and  
110 bacteria. There appear to be four subclasses of form I RubisCO, termed IA, IB, IC, and ID,  
111 found in different organisms (Tabita 1999). In prokaryotes and in nongreen algae the form I  
112 genes are cotranscribed behind a single promoter, while in plants and green algae, the *rbcL*  
113 gene is chloroplast-encoded and the *rbcS* gene is nuclear-encoded.

114 Form II RubisCO is found in different types of proteobacteria, and in one group of  
115 eukaryotes, the dinoflagellates. It is comprised of a single large subunit, which shares only  
116 about 30 percent sequence identity with form I large subunits (Tabita 1999). Depending on the  
117 organism, different numbers of dimeric pairs comprise the quaternary structure of form II  
118 RubisCO. Indeed, the fundamental unit common to all forms of RubisCO is the large subunit  
119 dimer in which the active site is formed from the interface between the N-terminal domain of one  
120 subunit and the  $\alpha/\beta$  barrel of the C-terminal domain of the second subunit. Form II proteins  
121 discriminate  $CO_2$  from  $O_2$  less well than form I RubisCO. Form II RubisCO is often found in  
122 organisms that also contain form I. In such instances, form II RubisCO appears not to be the  
123 major means to acquire carbon, but rather this RubisCO, along with other enzymes of the CBB  
124 pathway, functions to allow  $CO_2$  to be used as an electron acceptor to balance the intracellular  
125 redox potential when organic carbon is oxidized (Dubbs & Tabita 2004). In such organisms,  
126 form I RubisCO appears to be selectively synthesized when carbon or  $CO_2$  is limiting, consistent  
127 with the form I enzymes having a higher affinity for  $CO_2$  than form II proteins (Tabita 1999).

128 Form III RubisCO is found (thus far) only in archaea. In these organisms, as mentioned  
129 previously, the enzyme basically serves as a means to remove RuBP, which is produced by the  
130 isomerization of ribose 1,5-bisphosphate during purine/pyrimidine metabolism (Finn & Tabita  
131 2004, Sato et al. 2007). Many of the form III proteins thus far studied come from anaerobic  
132 extremophiles. Indeed, we have shown that several of these enzymes are highly oxygen  
133 sensitive (Finn & Tabita 2003), due largely to the extremely high affinity of these enzymes for O<sub>2</sub>  
134 (Kreel & Tabita 2007). Structurally, the proteins from methanogens and *Archaeoglobus fulgidus*  
135 (Finn & Tabita 2003; Watson et al. 1999) are found as dimers, although there are interesting  
136 exceptions such as the pentamer of dimers found in *Thermococcus kodakaraensis* (Ezaki et al.  
137 1999).

138 Form IV is also called the RubisCO-like protein (RLP) because this protein does not  
139 catalyze RubisCO activity (Hanson & Tabita 2001). Yet, there are similarities in the primary  
140 sequence (Hanson & Tabita 2001) and tertiary structure (Imker et al. 2007; Li et al. 2005) of  
141 these proteins that clearly indicate that RLPs are homologs of RubisCO and are derived from  
142 some common ancestor (Tabita et al. 2007). RLPs cannot catalyze RubisCO activity because  
143 there are alterations in key active site residues. The *Bacillus subtilis* (Ashida et al. 2003) and  
144 *Geobacillus kaustophilus* (Imker et al. 2007) RLP (or YkrW/MtnW), the cyanobacterial RLP  
145 from *Microcystis aeruginosa* (Carre-Mlouka et al. 2006), and RLPs from the photosynthetic  
146 bacteria *Rhodospirillum rubrum* and *Rhodopseudomonas palustris* (Singh and Tabita  
147 unpublished observations) participate in a methionine salvage pathway and catalyze the  
148 enolization of the RuBP analog, 2,3-diketo-5-methylthiopentyl-1-P in a reaction much akin to the  
149 enolase reaction catalyzed by RubisCO (Ashida et al. 2003; Imker et al. 2007). Physiological  
150 results indicate that RLP from the green sulfur bacterium *Chlorobium tepidum* is involved with  
151 some aspect of thiosulfate oxidation (Hanson & Tabita 2001; 2003), however the precise  
152 reaction catalyzed has not been identified as yet. Other clades of RLP molecules from different  
153 organisms have not yet been assigned a function; interestingly many of these latter RLP genes

154 do not complement RLP knockout strains from organisms with defined functions, indicative of  
155 different physiological roles (Singh & Tabita unpublished observations). Only one archaeon, *A.*  
156 *fulgidus*, and one eukaryote, the alga *Ostreococcus taurii*, have thus far been shown to contain  
157 an RLP gene (Tabita et al. 2007).

158

## 159 **Phylogenetic Relationships and the Evolution of RubisCO and the RubisCO-Like Protein:**

160

### 161 **Re-evaluating the Global Ocean Survey data set's contribution to RubisCO/RLP diversity**

162       Recent phylogenetic and bioinformatic analyses of RubisCO and RLP amino acid  
163 sequences provide a useful framework to understand the relationship of the different forms and  
164 how they may have evolved from a common ancestor (Tabita et al. 2007). Moreover, structural  
165 and functional studies impinge on these analyses as a coherent picture begins to emerge as to  
166 how the active site of this protein might have evolved and became adapted to different  
167 intracellular milieus. The overall conclusion from these studies was that a form III RubisCO from  
168 a methanogenic archaeon ancestor was the most probable source of all RubisCO and RLP  
169 lineages. Clearly, this conclusion is directly tied to the data currently available; thus it is certainly  
170 possible that additional sequences, when and if they are discovered, might cause this  
171 hypothesis to be reexamined.

172       In the report describing the construction and analysis of protein families encoded by  
173 DNA sequenced via the Global Ocean Survey (GOS) program (Yooseph et al. 2007), the claim  
174 was made that the GOS data contained thousands of previously unrecognized protein families  
175 and remarkably expanded our understanding of the sequence diversity present in well known  
176 protein families, including the RubisCO/RLP super family. However, the phylogeny presented  
177 for RubisCO in that work was strikingly different than those previously published by ourselves  
178 and others (Ashida et al. 2005; Carre-Mlouka et al. 2006; Hanson & Tabita 2001; Hanson &  
179 Tabita 2003; Tabita et al. 2007). To rigorously examine these claims with respect to the

180 RubisCO superfamily in the context of our previous work, here we present our re-analysis of the  
181 GOS-derived RubisCO sequences.

182 GOS sequences for analysis were collected in two batches, which were treated  
183 independently. First, all sequences within cluster 3734 used by Yooseph et al. in constructing  
184 their RubisCO/RLP phylogeny (Yooseph et al. 2007) were retrieved (Huiying Li, personal  
185 communication). Second, sequences representing each major lineage of the RubisCO  
186 superfamily were used as queries in BLASTP searches of the NCBI environmental non-  
187 redundant protein sequence database (env-nr), retaining all sequences that matched the  
188 example RubisCO/RLP sequences with E-values of less than  $1 \times 10^{-10}$ . Each batch of GOS-  
189 derived sequences was added independently to a FASTA file containing the recently described  
190 non-redundant collection of RubisCO and RLP amino acid sequences (Tabita et al. 2007). Both  
191 collections were then clustered with CD-Hit software (Li & Godzik 2006) at a cut off value of  
192 83% identity, which was chosen based on the ability of this value to discriminate and preserve  
193 sub-groups observed within major RubisCO/RLP lineages (i.e. form IA vs. IB). The longest  
194 sequence seeding each cluster was retained and all sequences representing clusters of less  
195 than 200 amino acid residues in length (<50% of the average  $448 \pm 40$  residue length of well  
196 characterized RubisCO/RLP's) were discarded. Remaining sequences were aligned by  
197 ClustalW (Thompson et al. 1994) and phylogenetic trees constructed in MEGA 4.0 (Tamura et  
198 al. 2007) using the Minimum Evolution method with a p-distance model of amino acid  
199 substitution and a gamma parameter of 1.554 for rate distributions between lineages.

200 Together, the approaches above created two sequence sets, one of 44 sequences from  
201 the cluster 3734 data and one of 35 sequences from the BLASTP search of env-nr. The sets  
202 had 31 sequences in common, leaving 13 unique sequences in the env-nr data (Table 2).  
203 When each set was used in concert with known RubisCO/RLP's in phylogenetic  
204 reconstructions, the trees that we observed for each set collected were nearly identical to one  
205 another and to that which we recently reported (Fig. 2). Indeed, the trees for each set reveal



206 that identical sequences were recovered and placed in extremely similar positions within each  
207 tree. The minor discrepancies between the trees are likely derived from unique sequences  
208 recovered within each set causing slight differences during the CD-Hit clustering of the  
209 independently collected sequence groups. Both trees indicate that there is a single,  
210 monophyletic group of GOS-derived RubisCO/RLP sequences that do not have any  
211 correspondent in currently analyzed genome sequences. As noted in our recent review (Tabita  
212 et al. 2007), the active site of this clade (IV-GOS) does not appear to be functional, thus making  
213 it an RLP lineage. The branching position of the IV-GOS clade within our trees is slightly  
214 different, though both positions are consistent with the group being an RLP lineage. With data  
215 from GOS cluster 3734 (Fig. 2A), the IV-GOS group emerges from between the IV-NonPhoto  
216 and IV-YkrW/MtnW clades. However, with GOS sequences retrieved by BLASTP (Fig. 2B), the  
217 IV-GOS clade is an early branch of the RLP's relative to the bona fide RubisCO lineage. The  
218 bona fide RubisCO's still appear to be a monphyletic clade in either analysis and are likely to  
219 have originated from an ancestral methanogen RubisCO as per our favored model for the  
220 evolution of RubisCO/RLP clades (see Fig. 9 of Tabita et al. 2007).

221 The GOS data set contained sequences that tree within every major clade of RubisCO  
222 and RLP with two exceptions. No GOS sequences were placed within the IV-YkrW/MtnW or IV-  
223 EnvOnly clades. The IV-YkrW/MtnW clade is restricted to Gram-positive organisms related to  
224 *Bacillus subtilis* while the IV-EnvOnly clade is derived exclusively from an acid mine drainage  
225 community metagenome (Tyson et al. 2004). This observation likely reflects a limited  
226 distribution of *Bacillus* spp. and acidophilic microbes in the surface ocean waters that provided  
227 the bulk of the GOS sequence data released to date (Yooseph et al. 2007).

228 The phylogenies presented here do not agree well with that presented by Yooseph et al  
229 (Yooseph et al. 2007). In particular, the distances between major clades within the  
230 RubisCO/RLP superfamily appear to be significantly shorter in the Yooseph et al. tree with a  
231 larger number of intervening branches, which were interpreted to mean that a large number of

232 additional RubisCO/RLP lineages were represented in the GOS data set (Yooseph et al. 2007).  
233 In addition, branch lengths between relatively closely related form I subgroups appear to be  
234 unusually long in the Yooseph *et al.* RubisCO tree. These incongruencies are likely due to three  
235 major differences in the approaches used for phylogenetic reconstruction.

236         The first difference relates to ORF prediction and length. Our current, as well as  
237 previous analyses (Tabita et al. 2007), utilized sequences that constitute a significant fraction of  
238 the length of currently known RubisCO/RLP sequences obtained from genomic sequences. In  
239 many of these instances, biochemical evidence existed for the start site and protein size. In  
240 contrast, the GOS data set, retrieved by either BLASTP or from cluster 3734, contains  
241 sequences as short as 61 amino acids with 35-45% of the data set having less than 200 amino  
242 acids (Fig. 3A). All protein start sites in the GOS sequences are predicted computationally, with  
243 no supporting biochemical evidence. The inclusion of short sequences presents a major  
244 difficulty in the reconstruction of phylogenetic relationships as they limit the number of  
245 informative positions in an alignment and increase the weight of each informative position  
246 towards determining both branching order and branch length. This problem is particularly acute  
247 if complete deletion of gaps is used in calculating distances from alignments for use in trees.  
248 This was not specified in details of tree construction methods used by these authors (Yooseph  
249 et al. 2007).

250         Given the uncertainties in start site prediction, it is not clear whether or not short  
251 sequences within the GOS data set represent authentic short variants of RubisCO/RLP or  
252 fragmentary sequences of longer genes where the remainder of the sequence had not been  
253 determined. One indication of fragmentation in the GOS data sets is that the initiation residue  
254 for the amino acid sequences was frequently not the canonical methionine (Fig. 3B). While  
255 alternative translation initiation sites do occur in prokaryotes (Makita et al. 2007; Tech &  
256 Meinicke 2006), this is usually a rare occurrence where leucine or valine specified by the  
257 codons T/CTG and GTG replace methionine. Over half (59%) of the ORFs in cluster 3734 data

258 set and just under half (47%) of GOS sequences retrieved by BLASTP did not initiate with M, L,  
259 or V while over 97% of the RubisCO/RLP sequences in our prior data set did (Fig. 3B).  
260 Therefore, it seems likely that the GOS ORF set may contain a large proportion of partial and/or  
261 corrupted ORF sequences that carry irrelevant sequence attached to true ORFs which may  
262 have driven spurious localization in phylogenetic tree reconstruction.

263         The observed difference in percentages of non-canonical initiation between the two GOS  
264 derived sets was found to depend on the database from which the sequences were retrieved. A  
265 large percentage of sequences in the cluster 3734 data set that initiated with a non-canonical  
266 amino acid were found in the BLASTP data set (retrieved from GenBank) to initiate with the  
267 canonical methionine. This is most easily seen by comparing the relative percentages of  
268 sequences initiating with I and M. The percentage of I-initiating sequences in the cluster 3734  
269 data is high relative to the BLASTP data while these two groups are reversed in M-initiating  
270 sequences (i.e. BLASTP percentage is greater than cluster 3734). Within the twenty nine  
271 sequences common between the final sets used to construct the trees, this I for M substitution  
272 was observed in ten sequences, or 34%. It is not clear how this substitution came about, but  
273 when sequences were aligned, this substitution was the only difference observed. Thus, users  
274 of the GOS data should be skeptical about the identity of the initiating residue when these  
275 sequences are retrieved from Genbank.

276         The second major difference between our analyses and those of Yooseph et al. (2007)  
277 has to do with the assumption of constant rates of sequence change across multiple lineages  
278 within protein superfamilies. As others have noted (DeBry 1992), the failure to account for rate  
279 variation can have profound effects on branch lengths observed in phylogenetic reconstructions  
280 and also has the tendency to confound certain types of phylogenetic reconstruction methods,  
281 notably UPGMA and Maximum Parsimony. In the specific case of the RubisCO/RLP lineages,  
282 analysis of the sequence data indicates that a moderate amount of rate variation between

283 lineages must be included to accurately reconstruct sequence relationships (Tabita et al. 2007).  
284 Rate variation was apparently not taken into account in the neighbor-joining phylogeny  
285 presented by Yooseph et al. nor was the data set evaluated for appropriate substitution model  
286 via a tool like ProtTest (Abascal et al. 2005).

287         The third and final difference has to do with clustering of sequences to reduce  
288 complexity. The GOS protein sequence data were clustered at values of 100% identity and  
289 98% similarity to produce non-redundant sets of protein sequences for further analysis. While  
290 this level will clearly delineate unique individual sequences, it is likely much too inclusive for the  
291 identification of clearly defined families or sub-families. Within all the RubisCO clades, average  
292 in-group pair-wise amino acid sequence identity range from a low of 39% within the IV-DeepYkr  
293 clade to a high of 85% within the form IB sub-family (Tabita et al. 2007). This is far below the  
294 100% identity criterion used to cluster the GOS data, likely leading to the inclusion of  
295 unnecessary sequences from a diversity analysis standpoint.

296         While the value of the GOS data is clear in terms of improving our understanding of the  
297 overall genetic diversity within marine microbial communities, the re-analysis we present here of  
298 the GOS RubisCO/RLP sequences suggests that a more rigorous analysis of the claim that  
299 "...the predicted GOS proteins also add a great deal of diversity to known protein families and  
300 shed light on their evolution" must be performed on a case by case basis. In the case of  
301 RubisCO/RLP, we conclude that only one truly novel lineage of RubisCO/RLP sequences was  
302 found within the GOS data and that the major contribution was to provide evidence for the  
303 existence of genes encoding previously defined major clades of RubisCO and RLP in marine  
304 environments.

305         In summary, currently available curated sequence data, as well as proper phylogenetic  
306 analyses, are compatible with there being three distinct lineages of bona fide RubisCO forms  
307 (forms I, II, and III); each group contains varying numbers of sub-groups. Within the form IV

308 (RLP) group of the RubisCO super family, six different clades of RLP molecules have thus far  
309 been identified (Tabita et al. 2007) and confirmed in the current study.

310

### 311 **Exploiting RubisCO Diversity to Learn More About Function**

312 Phylogenetic, bioinformatic, and evolutionary considerations provide thought-provoking  
313 discussions, but how can these analyses help us with solving some of the “mysteries” of  
314 RubisCO catalysis, such as the molecular basis for CO<sub>2</sub>/O<sub>2</sub> discrimination or why the affinities  
315 for CO<sub>2</sub> and O<sub>2</sub> vary? Our approach has always been to use what nature provides. Thus, the  
316 different forms and structural adaptations of RubisCO available, from organisms that assimilate  
317 and metabolize CO<sub>2</sub> under diverse and even extreme environments, may provide useful insights  
318 as to how all RubisCO molecules function.

319

### 320 **The interesting case of *Methanococcoides burtonii***

321 The recently sequenced genome of *Methanococcoides burtonii* revealed the presence of  
322 an ORF apparently coding for a RubisCO uniquely intermediate between form II and III (Fig. 4).  
323 The predicted protein sequence places the *M. burtonii* RubisCO (MBR) slightly closer to form II  
324 (40% identity to *Rhodobacter capsulatus* CbbM via BLASTP) than form III (35% identity to  
325 *Methanocaldococcus janaschii* RbcL). A more detailed look at the sequence, however, reveals  
326 a number of catalytic-site residues that are more similar to form II enzymes than to form III.

327 Biochemically, the classification of MBR is ambiguous as the enzyme has not been  
328 purified or characterized. *M. burtonii* is an obligately anaerobic methanogenic archaeon which  
329 would suggest that MBR functions in a role closer to that of previously characterized form III  
330 RubisCOs (Finn & Tabita 2004; Sato et al. 2007), namely, in pathways involved in recycling key  
331 metabolites rather than primary carbon fixation. Further, *M. burtonii* is apparently unable to  
332 grow with CO<sub>2</sub> as the sole carbon source (Franzmann et al. 1992). Several lines of evidence  
333 support this supposition. First, although MBR has been detected in whole-cell proteomic

334 studies, expression appears to be at levels consistent with pathways of secondary importance.  
335 Secondly, there is no evidence of other CBB cycle enzymes, especially phosphoribulokinase,  
336 the other enzyme unique to the CBB cycle. Finally, homologs of DeoA and E2B2 are present in  
337 the genome and are expressed under normal growth conditions (Goodchild et al. 2004a;  
338 Goodchild et al. 2005; Goodchild et al. 2004b; Saunders et al. 2005) indicative of the potential  
339 presence of an AMP-recycling pathway in which form III RubisCO is hypothesized to participate  
340 (Sato et al. 2006)

341 Complicating the classification of MBR, however, is the presence of a novel structural  
342 motif. Alignments of MBR amino acid sequence with form II and form III sequences indicates a  
343 26 residue sequence (EQTWSKIMDTDKDVINLVNEDLAHHVI) near the C-terminus with no  
344 homology to extant RubisCO sequences or, indeed, to any sequence currently deposited in  
345 GenBank. SwissProt models of the MBR structure, whether using form II (9rub) or form III  
346 (1geh) as a template, suggest the presence of a loop (Fig. 5), possibly forming an anti-parallel  
347  $\beta$ -sheet motif. This predicted motif is distant from the active site and may play a role in  
348 maintaining holoenzyme structure without directly impacting catalysis; certainly this remains to  
349 be determined. Given the huge array of sequences already known, however, this novel motif  
350 indicates that, at the very least, the RubisCO family still possesses surprises.

351

### 352 **Interactions of archaeal RubisCO with oxygen**

353 As discussed above, the profusion of genomic sequencing projects has played a  
354 dramatic role in altering perceptions as to how RubisCO might have evolved. With respect to  
355 structure/function studies, the discovery of form III enzymes from anaerobic archaea is  
356 particularly cogent as such organisms obviously evolved in the complete absence of oxygen. In  
357 particular, the finding that several representative archaeal RubisCO enzymes exhibit unusually  
358 high sensitivity to molecular oxygen, even in the presence of high levels of CO<sub>2</sub> (Watson et al.  
359 1999; Finn & Tabita 2003; Kreel & Tabita 2007), was deemed to be particularly interesting.

360 Thus, experiments were initiated to determine the basis for this sensitivity to O<sub>2</sub>, with the thought  
361 that these findings would relate or provide clues as to how all forms of RubisCO interact with O<sub>2</sub>.  
362 Using the RubisCO from *Archaeoglobus fulgidus* as a model, it was shown that O<sub>2</sub> sensitivity  
363 correlated with the fact that this enzyme showed an extremely high affinity for O<sub>2</sub>, with a K<sub>o</sub> of 5  
364 μM (Kreel & Tabita), some two orders of magnitude lower than the values obtained for form I or  
365 form II enzymes K<sub>o</sub> (which vary from 500 to 1000 μM). Two residues, Met-295 and Ser-363, of  
366 the *A. fulgidus* enzyme appeared to influence the K<sub>o</sub> value. Changing Met-295 to an aspartate  
367 (M295D) enhanced the ability of the enzyme to recover from O<sub>2</sub> inactivation with a consequent  
368 increase in the K<sub>o</sub> to 24 μM. This was accompanied by a three-fold increase in the substrate  
369 specificity factor. Similar results were obtained with Ser-363, a residue which sits in a  
370 hydrophobic pocket near the active site. Moreover, double mutants (i.e., M295D/S363I) showed  
371 an additive effect in the ability of such an enzyme to recover from oxygen inactivation (Kreel and  
372 Tabita 2007). The M295D/S363I enzyme has a K<sub>o</sub> of about 400-500 μM (Kreel & Tabita,  
373 manuscript in preparation), a K<sub>o</sub> that approximates the value for form I (aerobic and plant) and  
374 form II RubisCO (i.e., nearly 100-fold greater than the wild-type *A. fulgidus* enzyme). Despite a  
375 trade off in which the k<sub>cat</sub> has been reduced, this double mutant might serve as the starting point  
376 to select for a high activity, high K<sub>o</sub> enzyme.

377 Structurally, Met-295 of the *A. fulgidus* enzyme was found to be in close proximity to a  
378 highly conserved residue, Arg-279, found in all other forms of RubisCO and shown to be  
379 necessary for substrate (RuBP) binding (Zhang et al. 1994). The model structure suggests that  
380 a mutation to an aspartate residue at the Met-295 position would allow for an interaction  
381 between one of the hydroxyl side chains of the aspartate residue with one of the side chain  
382 nitrogen atoms of Arg-279. In wild-type *A. fulgidus* RubisCO, this hydrogen bond is absent.  
383 However, there is definite hydrogen-bonding to the equivalent Arg residue in all other form I and  
384 form II RubisCO structures. In addition, the model structure of *A. fulgidus* RubisCO shows an  
385 interaction of the side chain of Ser-363 with highly conserved and catalytically important

386 residues Gly-313 and Thr-314 of *A. fulgidus* RubisCO (Fig. 6). Gly-313 and Thr-314, found in  
387 all forms of RubisCO, show no ionic interactions with the amino acid residue equivalent to Ser-  
388 363 of *A. fulgidus* RubisCO in form I and form II enzymes. Thus, not only does a mutation in  
389 Ser-363 influence O<sub>2</sub>-sensitivity of the *A. fulgidus* enzyme but the degree of conservation of  
390 residues in this region may also provide clues as to the importance of this region in all RubisCO  
391 enzymes. In the model structure of the S363I and S363V mutants, it appears as though the  
392 introduction of a bulky hydrophobic amino acid in the hydrophobic pocket not only eliminates the  
393 hydrogen bonding interaction with highly conserved Gly-313 and Thr-314 (Fig. 6), but these  
394 substitutions may also cause conformational changes that likely affect the folding of the enzyme  
395 either before or during catalysis.

396

#### 397 **Interactions of form I RubisCO with oxygen**

398 Previous work showed that directed enzyme evolution procedures might be applied to  
399 bioselect useful mutant forms of prokaryotic RubisCO after random mutagenesis using a  
400 *Rhodobacter capsulatus* host with its host RubisCO genes deleted (Smith & Tabita 2003; 2004).  
401 This system (using *R. capsulatus* RubisCO deletion strain SBI/II<sup>-</sup>) offers many advantages,  
402 including the ability to select mutant enzymes that may or may not allow growth under aerobic  
403 conditions due to specific changes in the enzyme. During the course of isolating suppressor  
404 mutations that overcame the negative effects of a previously isolated D103V cyanobacterial  
405 (*Synechococcus* sp. strain PCC6301) RubisCO mutant, a residue was identified, Ala-375, that  
406 proved to be extremely interesting in another context (Satagopan et al. submitted for  
407 publication). The original suppressor, a D103V/A375V double mutant, was able to support  
408 growth of *R. capsulatus* under anaerobic conditions where the D103V enzyme could not.  
409 Moreover, upon kinetic analysis of the D103V/A375V, as well as an A375V mutant, it was found  
410 that the K<sub>o</sub> of both enzymes was considerably increased. Interestingly, structural analysis  
411 indicated that Ala-375 was situated in a hydrophobic pocket similar to that previously discussed



412 for the archaeal (*A. fulgidus*) enzyme. Indeed, Ala-375 of the form I cyanobacterial enzyme is  
413 equivalent to residue Ser-363 of the form III archaeal enzyme and this hydrophobic pocket is  
414 found in all forms of bona fide RubisCO (Fig. 7).

415         Since mutations of both Ser-363 of the *A. fulgidus* form III (archaeal) enzyme and Ala-  
416 375 of the *Synechococcus* form I enzyme affected the  $K_o$ , it was surmised that the A375V  
417 mutant enzyme might show a phenotypic response, especially under aerobic growth conditions.  
418 Thus, the affect of the A375V mutation was determined using the *R. capsulatus* system; i.e. *R.*  
419 *capsulatus* strain SBI/II<sup>-</sup> was complemented with the gene containing the A375V change (Fig. 8).  
420 In this experiment it is clear that only enzymes containing the A375V mutation are able to  
421 support luxuriant growth in strain SBI/II<sup>-</sup> (equivalent to wild-type *R. capsulatus* strain SB1003) in  
422 a complex medium under aerobic conditions. Neither the wild-type nor the D103V mutant  
423 *Synechococcus rbcL* gene supported luxuriant growth in strain SBI/II<sup>-</sup> using complex media (Fig.  
424 8). Furthermore, under aerobic chemoautotrophic conditions, with CO<sub>2</sub> as the sole source of  
425 carbon, only the A375V mutant supported substantial growth (Satagopan et al, submitted for  
426 publication). These studies indicate that residue Ala-375 influences the ability of RubisCO to  
427 interact with oxygen and its effect is manifested in a physiologically significant way.

428

## 429 **Conclusions**

430         From the analyses presented here it is clear that recent claims for new RubisCO families  
431 that fall outside the three classes of bona fide RubisCO proteins (forms I, II, and III) and the six  
432 clades of form IV or RLP molecules are not tenable. RubisCO from extremophile organisms that  
433 live in the absence of oxygen also provide convenient systems to learn more about how the  
434 active site of the enzyme has evolved to accommodate this special milieu. These studies also  
435 provide clues as to how RubisCO may be modified to enhance oxygen insensitivity in form I  
436 type enzymes.

437

438

439

440 **Acknowledgments**

441 FRT was supported by grant GM24497 from the National Institutes of Health and grants

442 DE-FG02-01ER63241 and DE-FG02-91ER20033 from the offices of Biological & Environmental

443 Research (Genomics: GTL Program) and Energy Biosciences, respectively, of the U. S.

444 Department of Energy. TEH was supported by by NSF Career Award MCB-0447649.

445

446

447

448

449

450

451

452

453 **Literature Cited**

454

455

456 Abascal, F., Zardoya, R. & Posada, D. 2005 ProtTest: selection of best-fit models of protein  
457 evolution. *Bioinformatics* **21**, 2104-2105.

458

459 Ashida, H., Danchin, A. & Yokota, A. 2005 Was photosynthetic RubisCO recruited by acquisitive  
460 evolution from RubisCO-like proteins involved in sulfur metabolism? *Res. Microbiol.* **156**, 611-  
461 618.

462

463 Ashida, H., Saito, Y., Kojima, C., Kobayashi, K., Ogasawara, N. & Yokota, A. 2003 A functional  
464 link between RubisCO-like protein of *Bacillus* and photosynthetic RubisCO. *Science* **302**, 286-  
465 290.

466

467 Carre-Mlouka, A., Mejean, A., Quillardet, P., Ashida, H., Saito, Y., Yokota, A., Callebaut, I.,  
468 Sekowska, A., Dittmann, E., Bouchier, C. & de Marsac, N. T. 2006 A new RubisCO-like protein  
469 coexists with a photosynthetic RubisCO in the planktonic cyanobacteria *Microcystis*. *J. Biol.*  
470 *Chem.* **281**, 24462-24471.

471

472 DeBry, R. W. 1992 The consistency of several phylogeny-inference methods  
473 under varying evolutionary rates. *Mol. Biol. Evol.* **9**, 537-551.

474

475 Dubbs, J. M. & Tabita, F. R. 2004 Regulators of nonsulfur purple phototrophic bacteria and the  
476 interactive control of CO<sub>2</sub> assimilation, nitrogen fixation, hydrogen metabolism and energy  
477 generation. *FEMS Microbiol. Rev.* **28**, 353-376.

478

479 Ezaki, S., Maeda, N., Kishimoto, T., Atomi, H. & Imanaka, T. 1999 Presence of a structurally  
480 novel type ribulose-bisphosphate carboxylase/oxygenase in the hyperthermophilic archaeon,  
481 *Pyrococcus kodakaraensis* KOD1. *J. Biol. Chem.* **274**, 5078-5082.

482

483 Finn, M. W. & Tabita, F. R. 2003 Synthesis of catalytically active form III ribulose 1,5-  
484 bisphosphate carboxylase/oxygenase in archaea. *J. Bacteriol.* **185**, 3049-3059.

485

486 Finn, M. W. & Tabita, F. R. 2004 Modified pathway to synthesize ribulose 1,5-bisphosphate in  
487 methanogenic archaea. *J Bacteriol* **186**, 6360-6366.

488

489 Franzmann, P. D., Springer, N., Ludwig, W., Demacario, E. C. & Rohde, M. 1992 A  
490 Methanogenic Archaeon from Ace Lake, Antarctica - *Methanococcoides-Burtonii* Sp-Nov.  
491 *System. & Applied Microbiol.* **15**, 573-581.

492

493 Goodchild, A., Raftery, M., Saunders, N. F. W., Guilhaus, M. & Cavicchioli, R. 2004a Biology of  
494 the cold adapted archaeon, *Methanococcoides burtonii* determined by proteomics using liquid  
495 chromatography-tandem mass spectrometry. *J.Proteome Res.* **3**, 1164-1176.

496

497 Goodchild, A., Raftery, M., Saunders, N. F. W., Guilhaus, M. & Cavicchioli, R. 2005 Cold  
498 adaptation of the Antarctic archaeon, *Methanococcoides burtonii* assessed by proteomics using  
499 ICAT. *J.Proteome Res.* **4**, 473-480.

500

501 Goodchild, A., Saunders, N. F. W., Ertan, H., Raftery, M., Guilhaus, M., Curmi, P. M. G. &  
502 Cavicchioli, R. 2004b A proteomic determination of cold adaptation in the Antarctic archaeon,  
503 *Methanococcoides burtonii*. *Mol.Microbiol.* **53**, 309-321.

504

505 Hanson, T. E. & Tabita, F. R. 2001 A ribulose-1,5-bisphosphate carboxylase/oxygenase  
506 (RubisCO)-like protein from *Chlorobium tepidum* that is involved with sulfur metabolism and the  
507 response to oxidative stress. *Proc. Natl. Acad. Sci. U. S. A.* **98**, 4397-402.  
508

509 Hanson, T. E. & Tabita, F. R. 2003 Insights into the stress response and sulfur metabolism  
510 revealed by proteome analysis of a *Chlorobium tepidum* mutant lacking the RubisCO-like  
511 protein. *Photosynth. Res.* **78**, 231-248.  
512

513 Imker, H. J., Fedorov, A. A., Fedorov, E. V., Almo, S. C. & Gerlt, J. A. 2007 Mechanistic  
514 diversity in the RubisCO superfamily: the "enolase" in the methionine salvage pathway in  
515 *Geobacillus kaustophilus*. *Biochemistry* **46**, 4077-4089.  
516

517 Kreeel, N. E. & Tabita, F. R. 2007 Substitutions at methionine 295 of *Archaeoglobus fulgidus*  
518 ribulose-1,5-bisphosphate carboxylase/oxygenase affect oxygen binding and CO<sub>2</sub>/O<sub>2</sub>  
519 specificity. *J. Biol. Chem.* **282**, 1341-1351.  
520

521 Li, H., Sawaya, M. R., Tabita, F. R. & Eisenberg, D. 2005 Crystal structure of a RubisCO-like  
522 protein from the green sulfur bacterium *Chlorobium tepidum*. *Structure* **13**, 779-789.  
523

524 Li, W. & Godzik, A. 2006 Cd-hit: a fast program for clustering and comparing large sets of  
525 protein or nucleotide sequences. *Bioinformatics* **22**, 1658-1659.  
526

527 Long, S. P., Ainsworth, E. A., Leakey, A. D. B., Nosberger, J. & Ort, D. R. 2006a. Food for  
528 thought: lower-than-expected crop yield stimulation with rising CO<sub>2</sub> concentrations. *Science* **312**,  
529 1918-1921.  
530

531 Long, S. P., Zhu, X.-G., Naidu, S. L. & Ort, D. R. 2006b. Can improvement in photosynthesis  
532 increase crop yields? *Plant Cell & Environ.* **29**, 315-330.

533

534 Makita, Y., de Hoon, M. J. & Danchin, A. 2007 Hon-yaku: a biology-driven Bayesian  
535 methodology for identifying translation initiation sites in prokaryotes. *BMC Bioinformatics* **8**, 47.

536

537 Sato, T., Atomi, H. & Imanaka, T. 2007 Archaeal type III RubisCOs function in a pathway for  
538 AMP metabolism. *Science* **315**, 1003-1006.

539

540 Saunders, N. F. W., Goodchild, A., Raftery, M., Guilhaus, M., Curmi, P. M. G. & Cavicchioli, R.  
541 2005 Predicted roles for hypothetical proteins in the low-temperature expressed proteome of the  
542 Antarctic archaeon *Methanococcoides burtonii*. *J. Proteome Res.* **4**, 464-472.

543

544 Smith S. A. & Tabita F. R. 2003 Positive and negative bioselection of mutant forms of  
545 prokaryotic (cyanobacterial) ribulose-1, 5-bisphosphate carboxylase/oxygenase. *J.Mol.Biol.* **331**,  
546 557-569.

547

548 Smith, S. A. & Tabita, F. R. 2004 Glycine 176 affects catalytic properties and stability of the  
549 *Synechococcus* sp. strain PCC 6301 ribulose 1,5-bisphosphate carboxylase/oxygenase. *J. Biol.*  
550 *Chem.* **279**, 25632-25637.

551

552

553 Tabita, F. R. 1999 Microbial ribulose 1,5-bisphosphate carboxylase/oxygenase: A different  
554 perspective. *Photosynth. Res.* **60**, 1-28.

555

556 Tabita, F. R., Hanson, T. E., Li, H., Satagopan, S., Singh, J. & Chan, S. 2007 Function,  
557 structure, and evolution of the RubisCO-like proteins and their RubisCO homologs. *Microbiol.*  
558 *Mol. Biol. Rev.* **71** In Press.

559

560 Tamura, K., Dudley, J., Nei, M. & Kumar, S. 2007 MEGA4: Molecular Evolutionary Genetics  
561 Analysis (MEGA) software version 4.0. *Mol. Biol. Evol.* **24**, 1596-1599.

562

563 Tech, M. & Meinicke, P. 2006 An unsupervised classification scheme for improving predictions  
564 of prokaryotic TIS. *BMC Bioinformatics* **7**, 121.

565

566 Thompson, J. D., Higgins, D. G. & Gibson, T. J. 1994 CLUSTAL W: improving the sensitivity of  
567 progressive multiple sequence alignment through sequence weighting, position-specific gap  
568 penalties and weight matrix choice. *Nucleic Acids Res.* **22**, 4673-4680.

569

570 Tyson, G. W., Chapman, J., Hugenholtz, P., Allen, E. E., Ram, R. J., Richardson, P. M.,  
571 Solovyev, V. V., Rubin, E. M., Rokhsar, D. S. & Banfield, J. F. 2004 Community structure and  
572 metabolism through reconstruction of microbial genomes from the environment. *Nature* **428**, 37-  
573 43.

574

575 Watson, G. M., Yu, J. P. & Tabita, F. R. 1999 Unusual ribulose 1,5-bisphosphate  
576 carboxylase/oxygenase of anoxic Archaea. *J. Bacteriol.* **181**, 1569-1575.

577

578 Yooseph, S., Sutton, G., Rusch, D. B., Halpern, A. L., Williamson, S. J., Remington, K., Eisen,  
579 J. A., Heidelberg, K. B., Manning, G., Li, W., Jaroszewski, L., Cieplak, P., Miller, C. S., Li, H.,  
580 Mashiyama, S. T., Joachimiak, M. P., van Belle, C., Chandonia, J. M., Soergel, D. A., Zhai, Y.,  
581 Natarajan, K., Lee, S., Raphael, B. J., Bafna, V., Friedman, R., Brenner, S. E., Godzik, A.,

582 Eisenberg, D., Dixon, J. E., Taylor, S. S., Strausberg, R. L., Frazier, M. & Venter, J. C. 2007 The  
583 Sorcerer II Global Ocean Sampling Expedition: Expanding the Universe of Protein Families.  
584 *PLoS Biol.* **5**, e16.

585

586 Zhang, K. Y., Cascio, D., & Eisenberg, D. 1994. Crystal structure of the unactivated ribulose  
587 1,5-bisphosphate carboxylase/oxygenase complexed with a transition state analog, 2-carboxy-  
588 D-arabinitol 1,5-bisphosphate. *Protein Sci.* **3**, 64-69.

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609



610 **Table I. Summary of the properties of different forms of RubisCO**

RubisCO form	Quaternary structure	Types of organisms
I	$L_8S_8$	Plants, algae, proteobacteria, cyanobacteria
II	$(L_2)_n$	Proteobacteria, dinoflagellates
III	$(L_2)_n$	Archaea
IV	$L_2$	Proteobacteria, cyanobacteria, archaea, algae

611 **Table 2. All sequences used in this study.**

Family	Organism or Environment	Accession	Label
IA	Bradyrhizobium sp. BTAi1	ZP_00862357	Bra.sp-BTAi1-1
	Hydrogenophilus thermoluteolus	Q51856	Hyd.the
	Nitrobacter vulgaris	Q59613	Nit.vul
IB	Anabaena variabilis ATCC 29413	ABA23512	Ana.var-ATCC 29413
	Nicotiana tabacum	P00876	Nic.tab
IC	Rubrivivax gelatinosus PM1	ZP_00243775	IC-1-Rub.gel
	Xanthobacter autotrophicus Py2	ZP_01200598	Xan.aut-Py2-1
	Chlorothrix halobium	D. Bryant Pers Comm	ChlhalRbc1
ID	Acidiphilium cryptum JF-5	ZP_01146719	Aci.cry-JF-5-1
	Aurantimonas sp. SI85-9A1	AAB41464	Mn-Ox-SI85-9A1
	Burkholderia xenovorans LB400	ZP_00283421	ID-Bur.xen-LB400
	Nitrosococcus oceani ATCC 19707	ABA56859	Nit.oce-ATCC 19707
	Nitrosospora multififormis ATCC 251966	YP_411385	Nit.mul-ATCC251966
	Odontella sinensis	NP_043654	Odo.sin
	Pleurochrysis carterae	Q08051	Ple.car
II	Hydrogenovibrio marinus	Q59462	II-Hyd.mar
	Lingulodinium polyedrum	AAA98748	Lin.pol
	Magnetospirillum magnetotacticum AMB-1	YP_422059	Mag.mag-AMB-1
	Rhodoferrax ferrireducens T118	YP_522655	Rho.fer-T118
		CAA25080	II-Rho.rub
	Symbiodinium sp.	AAG37859	Sym.sp
	Thiomicrospira crunigena XCL-2	ABB41020	II-Thi.cru-XCL-2
Methanococcoides burtonii DSM 6242	ZP_00563653	Met.bur-DSM6242	
Unaffiliated archaea	Methanosaeta thermophila PT	ZP_01153096	Met.the-PT
	Methanospirillum hungatei JF-1	YP_503739	Met.hun-JF-1
III	Archaeoglobus fulgidus DSM 4304	NP_070466	III-Arc.ful-DSM 4304
	Haloferax volcanii	contig 3020 <sup>1</sup>	Hal.vol
	Hyperthermus butylicus DSM 5456	YP_001012710	Hyp.the-DSM5456
	Methanocaldococcus jannaschii	AAB99239	Met.jan
	Methanoculleus marisnigri JR1	ZP_01391373	Met.cul-JR1
	Methanosarcina acetivorans C2A	AAM07894	Met.ace-C2A
	Natronomonas pharaonis DSM 2160	CAI49476	Nat.pha-DSM 2160
	Pyrococcus horikoshii OT3	BAA30036	Pyr.hor-OT3
	Thermococcus kodakaraensis KOD1	BAD86479	The.kod-KOD1
	Thermofilum pendens Hrk 5	YP_920628	The.pen-Hrk5
	IV-Non-photo	Acidiphilium cryptum JF-5	ZP_01146529
Bordetella bronchiseptica RB50		CAE31534	Bor.bro-RB50
Burkholderia xenovorans LB400		ZP_00284840	IV-Bur.xen-LB400
Chromohalobacter salexigens DSM 3043		ZP_00471249	Chr.sal-DSM 3043
Delftia acidovorans SPH1		ZP_01577127	Del.aci-SPH1
Fulvimarina pelagi HTCC2506		ZP_01438569	Ful.mar-HTCC2506
Jannaschia sp. CCS1		YP_511005	Jan.sp-CCS1
Mesorhizobium loti		BAB53192	Mes.lot
Polaromonas sp. JS666		ZP_00502320	IV-1-Pol.sp-JS666
Polaromonas sp. JS666		ZP_00502381	IV-2-Pol.sp-JS666

	<i>Pseudomonas putida</i> F1	ZP_00900417	Pse.put-F1
	<i>Rhizobium leguminosarum</i>	pRL120396 <sup>2</sup>	Rhi.leg
	<i>Roseobacter</i> sp. MED193	ZP_01056409	Ros.sp-MED193
	<i>Sinorhizobium meliloti</i> 1021	CAC48779	IV-Sin.mel
	<i>Thermomicrobium roseum</i>	contig:2220 <sup>3</sup>	The.ros
	<i>Xanthobacter autotrophicus</i> Py2	ZP_01199940	Xan.aut-Py2-2
IV-DeepYkr	<i>Alkalilimnicola ehrlichei</i> MLHE-1	YP_742007	Alk.ehr-MLHE-1-2
	<i>Halorhodospira halophila</i> SL1	YP_001002057	Hal.hal-SL1-2
	<i>Heliobacillus mobilis</i>	ABH04879	Hel.mob
	<i>Ostreococcus tauri</i>	Ot07g01830 <sup>4</sup>	IV-1-Ost.tau
	<i>Ostreococcus tauri</i>	Ot08g02600 <sup>4</sup>	IV-2-Ost.tau
	<i>Rhodopseudomonas palustris</i> BisA53	YP_782588	IV-1-Rho.pal-BisA53
	<i>Rhodopseudomonas palustris</i> BisB18	YP_532057	IV-1-Rho.pal-BisB18
	<i>Rhodopseudomonas palustris</i> BisB5	YP_569369	IV-1-Rho.pal-BisB5
	<i>Rhodopseudomonas palustris</i> CGA009	CAE27610	IV-1-Rho.pal-CGA009
	<i>Rhodospirillum rubrum</i>	ABC22798	IV-Rho.rub
IV-EnvOnly	Acid Mine Consortium	AADL01000066	AMC-066
	Acid Mine Consortium	AADL01000179	AMC-179
	Acid Mine Consortium	AADL01000602	AMC602
IV-Photo	<i>Allochromatium vinosum</i>	BAB44150	IV-All.vin
	<i>Chlorobium chlorochromatii</i> CaD3	ABB28892	Chl.chl-CaD3
	<i>Chlorobium phaeobacteroides</i> DSM 266	ZP_00527577	Chl.pha-DSM 266
	<i>Chlorobium tepidum</i> TLS1	AAM72993	Chl.tep-TLS1
	<i>Prosthecochloris aestuarii</i> DSM 271	ZP_00590874	Pro.aes-DSM 271
	<i>Rhodopseudomonas palustris</i> BisB18	YP_530146	IV-2-Rho.pal-BisB18
IV-YkrW	Acid Mine Consortium	AADL01000541	AMC-541
	<i>Bacillus cereus</i> E33L	AAU16474	Bac.cer-E33L
	<i>Bacillus clausii</i> KSM-K16	BAD64310	Bac.cla-KSM-K16
	<i>Bacillus licheniformis</i> ATCC 14580	AAU23062	Bac.lic-ATCC-14580
	<i>Bacillus subtilis</i> subsp. <i>subtilis</i> str. 168	CAB13232	Bac.sub-168
	<i>Exiguobacterium sibiricum</i> 255-15	ZP_00539172	Exi.sib-255-15
	<i>Geobacillus stearothermophilus</i>	Contig321 <sup>5</sup>	Geo.ste
Unique GOS Sequences	GOS Cluster 3734	EB060441	EB060441
		EBH72800	EBH72800
		EBK38093	EBK38093
		ECC63254	ECC63254
		ECH76386	ECH76386
		ECI09606	ECI09606
		ECK47660	ECK47660
		ECM15273	ECM15273
		ECN61796	ECN61796
		ECV58360	ECV58360
		ECV93956	ECV93956
		ECW52658	ECW52658
		EDE24211	EDE24211
	BLASTP vs. env_nr	EBH72800	EBH72800
		ECE36765	ECE36765
		EDE26668	EDE26668

Common to 3734 and env\_nr

EBE68134	EBE68134
EBG47321	EBG47321
EBH57905	EBH57905
EBJ85584	EBJ85584
EBK48460	EBK48460
EBO94763	EBO94763
EBO96920	EBO96920
EBS14603	EBS14603
EBV75076	EBV75076
ECA82964	ECA82964
ECD28674	ECD28674
ECE21553	ECE21553
ECO29831	ECO29831
ECQ06045	ECQ06045
ECQ25263	ECQ25263
ECQ91781	ECQ91781
ECR75886	ECR75886
ECU18502	ECU18502
ECU62179	ECU62179
ECV95656	ECV95656
ECW18074	ECW18074
ECW49692	ECW49692
EDD64574	EDD64574
EDE27295	EDE27295
EDG90708	EDG90708
EDH70269	EDH70269
EDH87249	EDH87249
EDJ07695	EDJ07695
EDJ25372	EDJ25372
EDJ35923	EDJ35923
EDJ36865	EDJ36865

---

1-Source: <http://zdna2.umbi.umd.edu/Public/hvo/>

2-Source: [http://www.sanger.ac.uk/Projects/R\\_leguminosarum/](http://www.sanger.ac.uk/Projects/R_leguminosarum/)

3-Source: [http://tigrblast.tigr.org/ufmg/index.cgi?database=t\\_roseum%7Cseq](http://tigrblast.tigr.org/ufmg/index.cgi?database=t_roseum%7Cseq)

4-Source: <http://bioinformatics.psb.ugent.be/blast/public/?project=ostreococcus>

5-Source: [http://www.genome.ou.edu/bstearo\\_blast.html](http://www.genome.ou.edu/bstearo_blast.html)

613 **Figure Legends**

614

615 **Fig. 1.** Reactions catalyzed by RubisCO showing the inherent problem that aerobic organisms  
616 face as a result of the O<sub>2</sub> fixation reaction catalyzed by the enzyme. CO<sub>2</sub>/O<sub>2</sub> substrate  
617 discrimination may be determined experimentally as described in the text.

618

619 **Fig. 2.** Minimum evolution phylogenetic trees of the RubisCO and RLP protein superfamily. (A)  
620 Tree built with sequences derived from GOS cluster 3734. (B) Tree built with GOS sequences  
621 derived from BLASTP searches of the GenBank env\_nr database. Major lineages are labeled  
622 and color coded identically in each panel. Bona fide RubisCO lineages were given individual  
623 colors (I-green and red, II-purple, III-gold), while all RLP's are colored Cyan except for the  
624 unique GOS RLP lineage, IV-GOS, which is colored red-brown. Abbreviations for non-GOS  
625 sequences are listed in Table X. GOS sequences are identified by their NCBI accession  
626 number.

627

628 **Fig. 3.** Comparison of starting data sets used to construct phylogenetic trees. (A) The  
629 percentage of total sequences falling into a given length category are shown for sequences in  
630 GOS cluster 3734 (dark grey bars, n = 195), GOS RubisCO/RLP sequences retrieved by  
631 BLASTP from Genbank (light grey bars, n = 377), and the set previously used for phylogenetic  
632 reconstruction (white bars, n = 191). B) Distribution of initiating residues for RubisCO and RLP  
633 sequences within data sets used to construct phylogenetic trees. The initiating residue is  
634 indicated on the horizontal axis. The bars for panel (B) denote the same data sets as in panel  
635 (A).

636

637

638 **Fig. 4.** Phylogenetic position of *Methanococcoides burtonii* (*M. bur* in figure)) RubisCO  
639 amongst known form II and form III RubisCO sequences. This is a subset of the full RubisCO  
640 phylogenetic tree (Fig. 2) that highlights the *M. burtonii* enzyme's unique position.

641  
642 **Fig. 5.** Modeled structure of the *M. burtonii* RubisCO with the loop structure shown in gold. The  
643 red bar represents the approximate position of RuBP at the active site.

644  
645 **Fig. 6.** The hydrophobic pocket of the *Archaeoglobus fulgidus* form III RubisCO showing  
646 interactions of Ser-363 with conserved residues Gly-313 and Thr-314.

647  
648 **Fig. 7.** Hydrophobic pocket regions surrounding equivalent residues (in green) implicated in  
649 oxygen interactions. (A) Ala-375 of the form I *Synechococcus* PCC6301 RubisCO, (B) Ile-367  
650 of the form II *Rhodospirillum rubrum* enzyme, and (C) Ser-366 of the archaeal *Pyrococcus*  
651 *kodakaraensis* form III enzyme. The adjacent active-site Ser residue is colored blue in all three  
652 structures. Ligands at the active-site, i.e. CABP in (A) and sulfate in (C) are colored black.  
653 Residue side-chains colored in yellow are within 4 Å of the Ala-375/Ile-367/Ser-366 of the form  
654 I/II/III proteins, respectively.

655  
656 **Fig. 8.** Phenotypic response of mutant form I *Synechococcus* PCC6301 RubisCO constructs in  
657 *R. capsulatus* strain SBI/II<sup>-</sup>. All strains were grown under aerobic conditions in a complex  
658 peptone-yeast extract media. Sectors contained (A) *R. capsulatus* wild-type strain SB1003, (B)  
659 *R. capsulatus* strain SBI/II<sup>-</sup> with no complementing RubisCO gene; (C) *R. capsulatus* strain  
660 SBI/II<sup>-</sup> complemented with the wild-type *Synechococcus* PCC 6301 *rbcLS* genes; (D) *R.*  
661 *capsulatus* strain SBI/II<sup>-</sup> complemented with the D103V mutant *Synechococcus* PCC 6301 *rbcL*  
662 gene and wild-type *rbcS*; (E) *R. capsulatus* strain SBI/II<sup>-</sup> complemented with the A375V mutant  
663 *Synechococcus* PCC 6301 *rbcL* gene and wild-type *rbcS*; (F) *R. capsulatus* strain SBI/II<sup>-</sup>

664 complemented with the D103V/A375V mutant *Synechococcus* PCC 6301 *rbcL* gene and wild-  
665 type *rbcS*.

666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677

Fig. 1

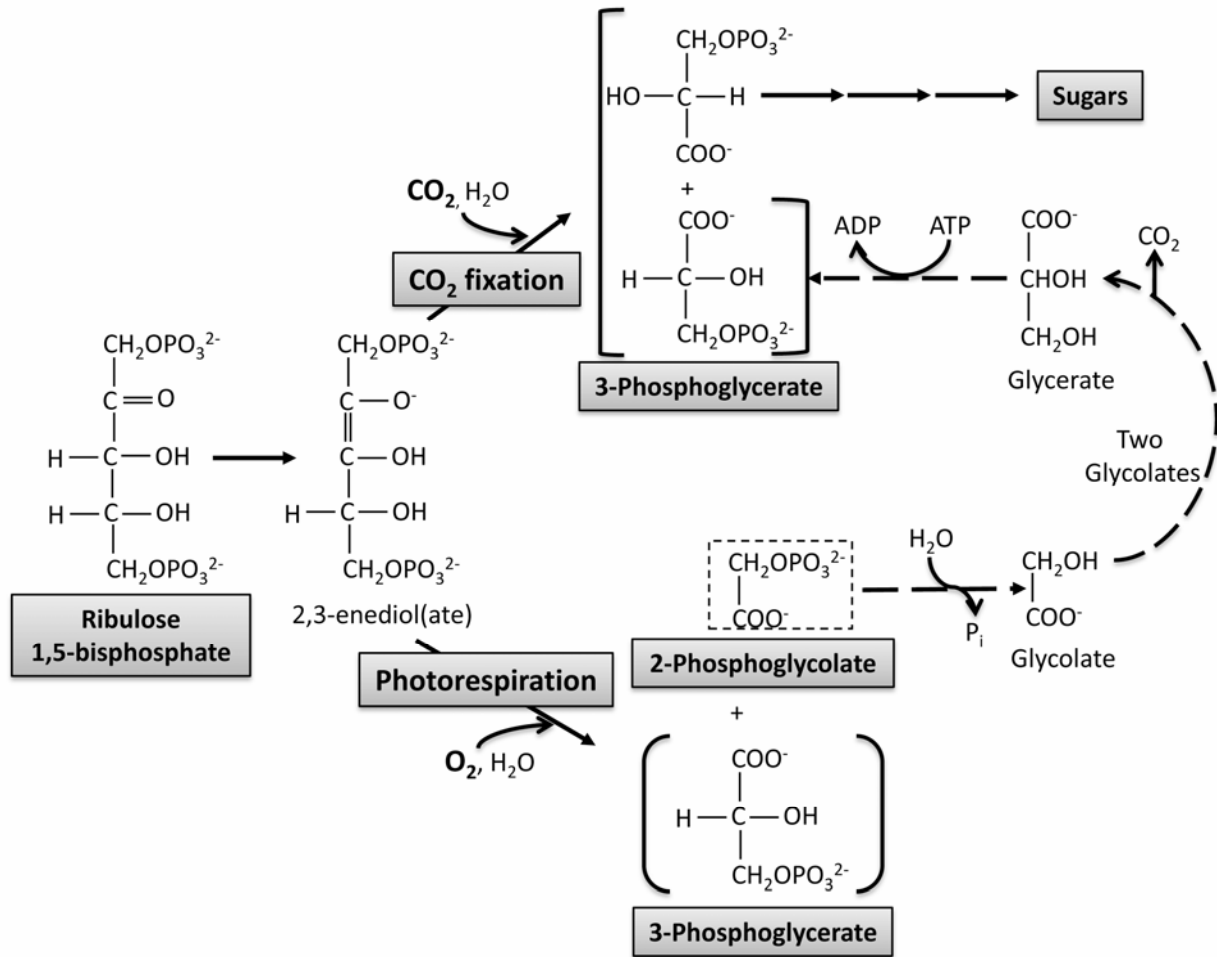




Fig. 2A

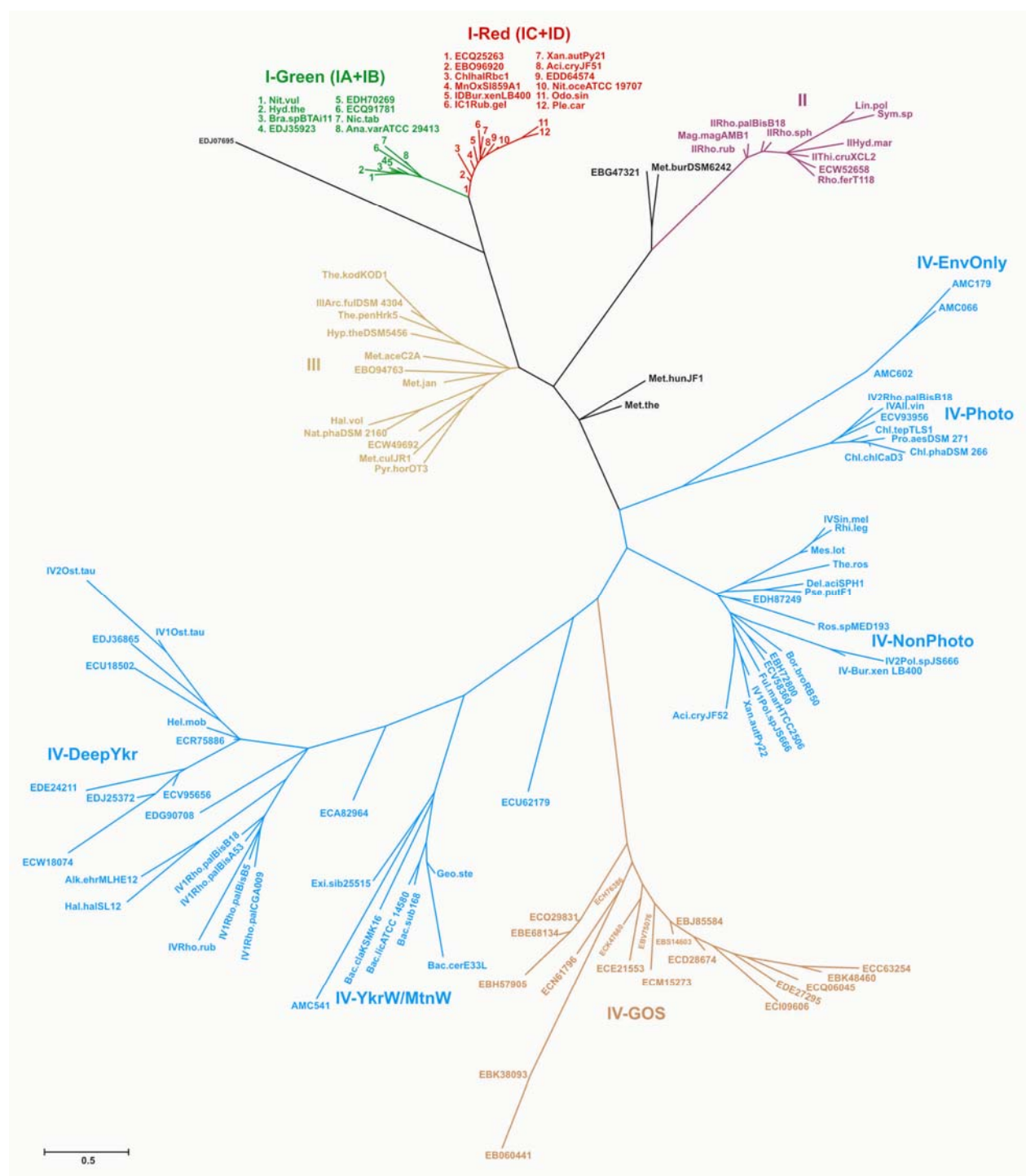


Fig. 2B

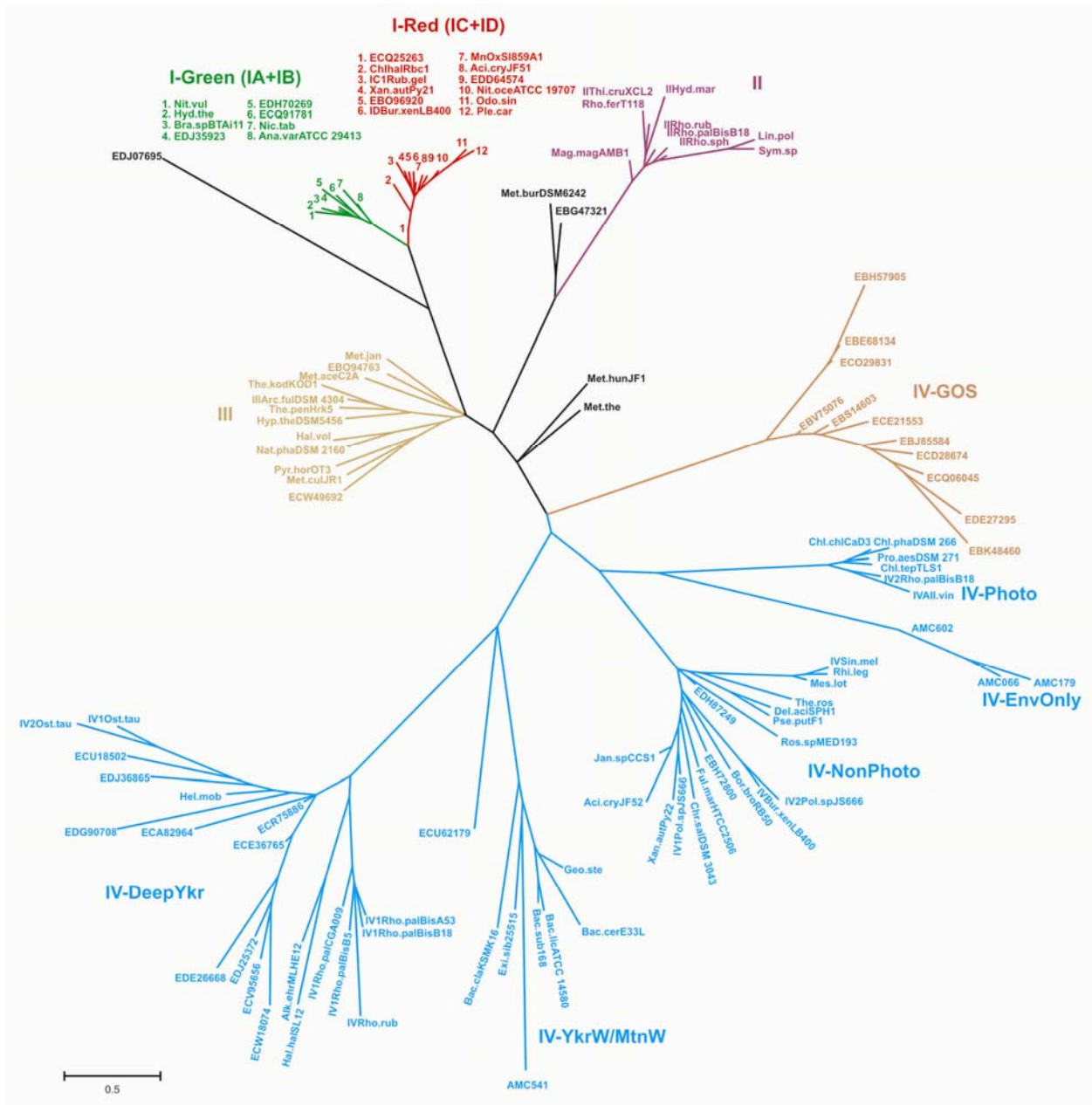


Fig. 3A

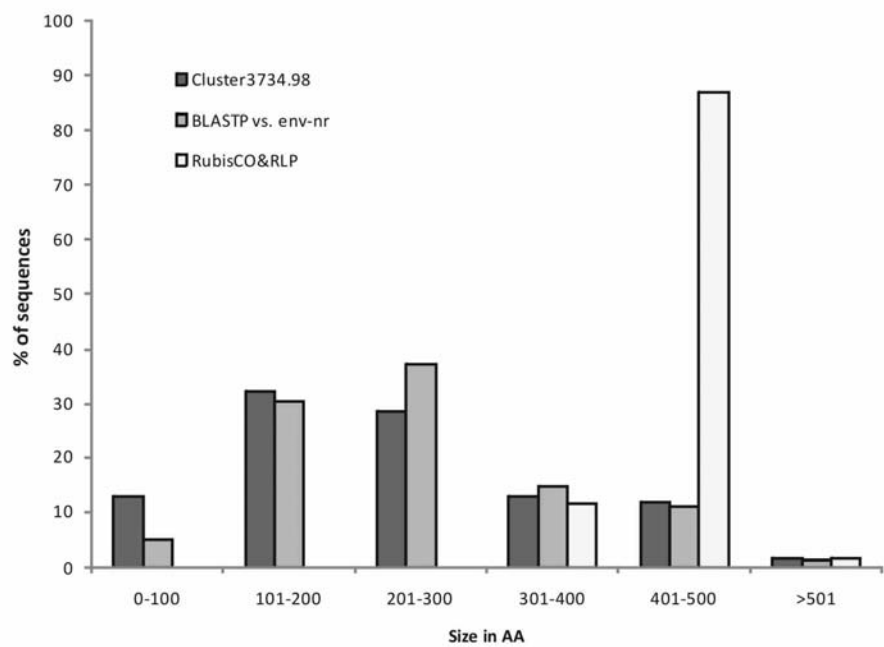


Fig. 3B

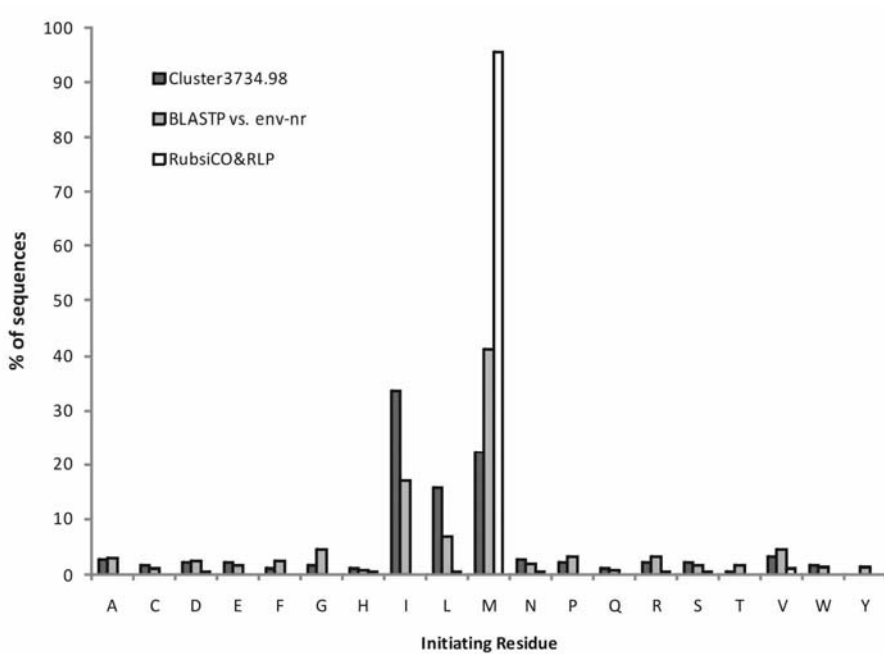


Fig. 4

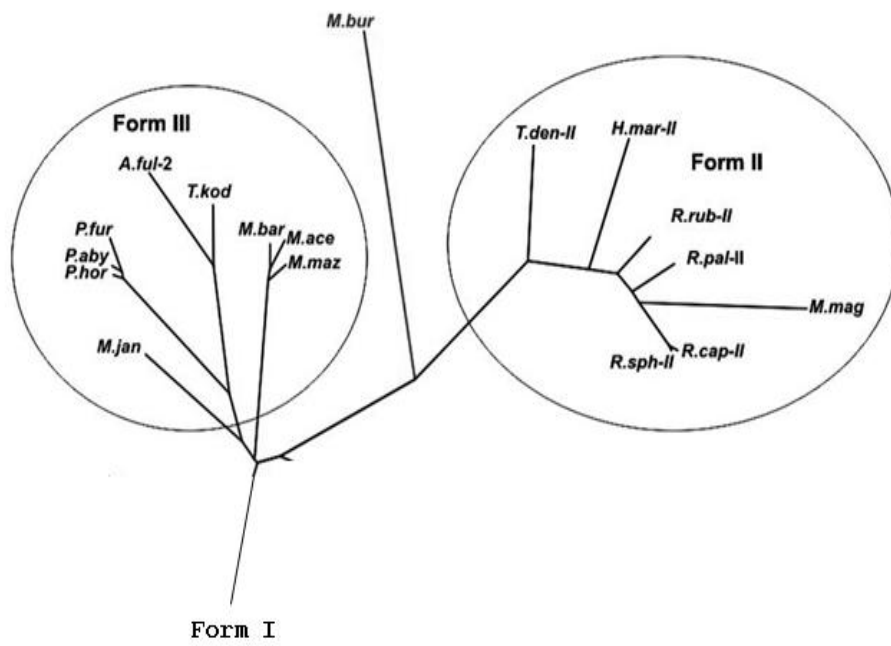


Fig. 5



Fig. 6

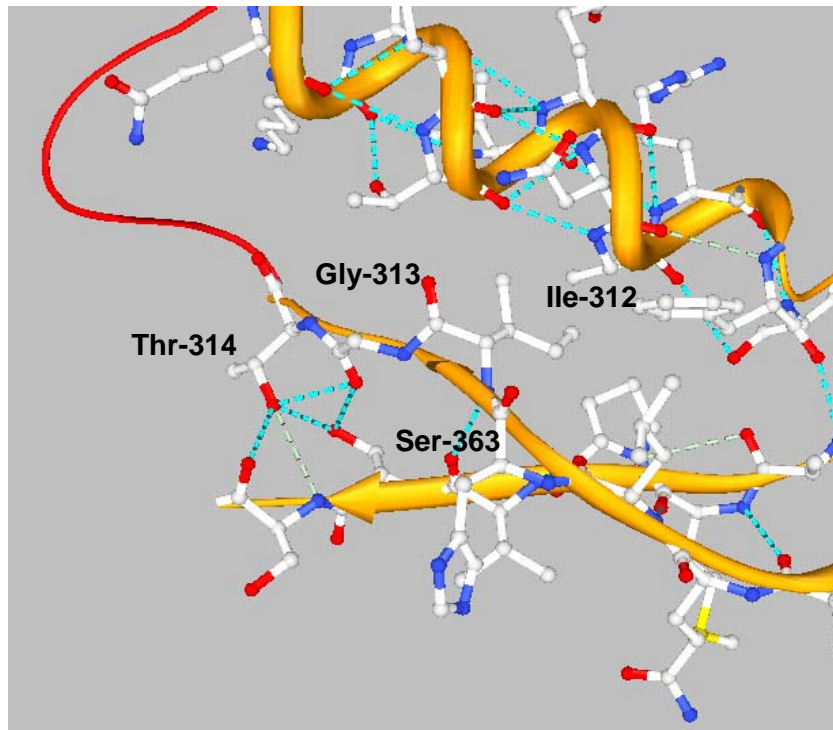




Fig. 7

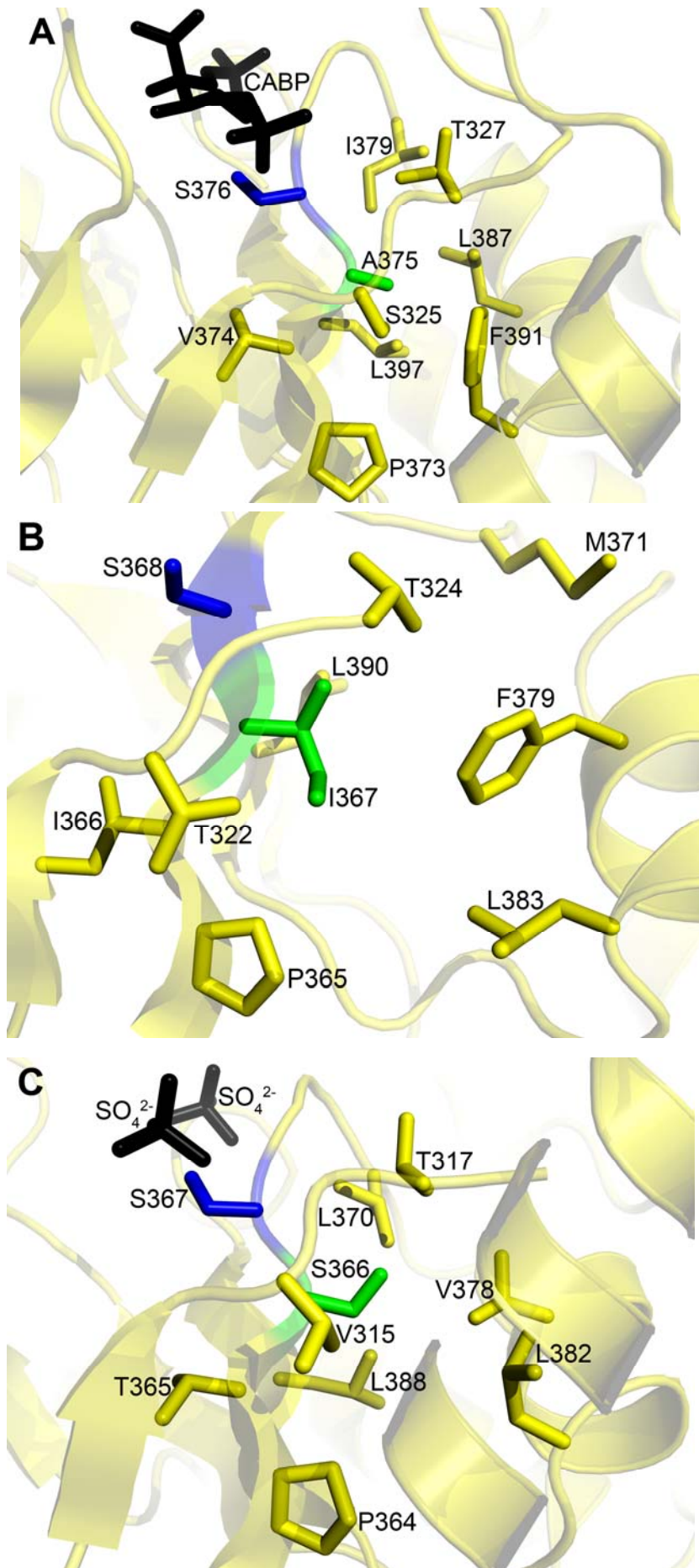


Fig. 8

