

1 **1. Extended Data**

2

3

Figure #	Figure title One sentence only	Filename This should be the name the file is saved as when it is uploaded to our system. Please include the file extension. i.e.: <i>Smith_ED Fig1.jpg</i>	Figure Legend If you are citing a reference for the first time in these legends, please include all new references in the Online Methods References section, and carry on the numbering from the main References section of the paper.
Extended Data Fig. 1	Distribution of Form I' Chloroflexi genomes.	Shih_EDFig1.eps	Maximum likelihood phylogenetic tree of Chloroflexi using ribosomal protein S3 (rpS3) as a marker gene. To map the distribution of Form I' Rubisco genes onto genomes, all MAGs were scanned for presence of both rpS3 and Form I' Rubisco. MAGs containing Form I' Rubisco are highlighted in orange. The scaffolds that encode RbcL vary in size substantially, ranging up to ~106 kbp in length (available as Supplementary Data). At least partial genomic context could be determined in most cases and the gene for phosphoribulokinase was adjacent. In some cases, additional CBB Cycle pathway genes were present in an operon with Rubisco, strongly supporting the function of Rubisco in this pathway. In a subset of cases, other pentose phosphate pathway genes were co-encoded. In no case was there evidence for RbcS, either on the scaffold or in the draft genome bin (where a bin was available). Gene predictions were established via a standard annotation pipeline ^{50,51} and augmented by HMM-based profiling and domain analysis.
Extended Data Fig. 2	In Form I-containing Chloroflexi operons, <i>rbcL</i> and <i>rbcS</i> are always found next to each	Shih_EDFig2.eps	Fragment operons from an example set of 10 Form I Rubisco-containing Chloroflexi genomes shows that <i>rbcS</i> is always found next to <i>rbcL</i> , similar to Form I Rubisco found in Cyanobacteria and Proteobacteria ¹¹ . Form I' Rubisco-

	other, unlike Form I'-containing Chloroflexi operons that lack <i>rbcS</i>.		containing Chloroflexi genomes do not contain small subunit <i>rbcS</i> (Fig. 1b). Scaffold names are shown to the right of their corresponding genome fragments.
Extended Data Fig. 3	PAGE analyses.	Shih_EDFig3.eps	a , Non-denaturing PAGE gel with a molecular weight marker (M, lane 1), and purified proteins of all three candidate Form I' Rubisco (<i>P. breve</i> , 241187, and 170907) with (+) or without (-) prior activation and incubation with 10-fold molar excess of 2CABP. 241187 and 170907 denotes scaffolds B_1_S1_170907_scaffold_241187_5_Tax=RBG_16_Chloroflexi_63_12 and S_p2_S4_170907_scaffold_85440 Rubisco, respectively. b , SDS-PAGE analysis of crude cell lysate from 1) over-expression of untagged <i>P. breve</i> Rubisco with co-expression of GroEL/ES from pBAD33EL/ES, 2) over-expression of His ₁₄ -bdSUMO-tagged <i>P. breve</i> Rubisco with co-expression of GroEL/ES from pBAD33EL/ES, and 3) over-expression of His ₁₄ -bdSUMO-tagged <i>P. breve</i> Rubisco without overexpression of GroEL/ES (background GroEL/ES expression from <i>E. coli</i>). Without GroEL/ES overexpression, untagged RbcL comprises 8 ± 1.0 ($n = 3$) of the total soluble protein, which improves to 14 ± 0.5 ($n = 3$) when GroEL/ES overexpression is induced (see Methods). When the His ₁₄ -bdSUMO tag is included on the N-terminal end of RbcL, soluble expression is 7 ± 0.8 ($n = 3$) and 14 ± 0.8 ($n = 3$) of the total soluble protein, without and with GroEL/ES overexpression, respectively. Reported values collected from n separate experiments (separately grown <i>E. coli</i> cultures) reflect the mean \pm standard deviation.
Extended Data Fig. 4	Form I Rubisco possess a unique RbcL C-terminal extension that interacts with RbcS, which is not found in Form I' Rubisco.	Shih_EDFig4.tiff	a , Sequence alignment of representative Rubisco RbcL sequences from Forms I, I', II, II/III, IIIA and IIIb. Strictly conserved residues have a red background, residues well conserved within a group are indicated by red letters, and the remaining residues are in black letters. Gaps are represented by dots. Residue numbering along the top refers to <i>P. breve</i> RbcL. Symbols above blocks of sequences correspond to the

			<p>secondary structure of <i>P. breve</i> RbcL: α, α-helix; β, β-strand; η, 310-helix. The secondary structure elements were named according to Knight et al., 1990⁵². The positions of loop 6 (black dotted lined), the Form II/III-specific Rubisco assembly domain (cyan line), and the Form I-specific C-terminal extension (purple line) are indicated. The RbcX binding domain-specific to Form IB Rubisco is boxed in pink. The sequence alignment was created using the UniProt RbcL sequences P22859 (<i>Allochromatium vinosum</i>), O85040 (<i>Halothiobacillus neapolitanus</i>), A0A4D4IZ26 (<i>Zea mays</i>), P00880 (<i>Syn6301</i>), Q1QH22 (<i>Nitrobacter hamburgensis</i>), Q3IYC2 (<i>Rhodobacter sphaeroides</i>), P51226 (<i>Porphyra purpurea</i>), Q9GGQ2 (<i>Vaucheria litorea</i>), E1IGS1 (<i>Oscillochloris trichoides</i>), A0A0P9FAF0 (<i>Kouleothrix aurantiaca</i>), A4WW35 (<i>Rhodobacter sphaeroides</i>), P04718 (<i>Rhodospirillum rubrum</i>), Q12TQ0 (<i>Methanococcoides burtonii</i>), A0A1L3Q3Y6 (<i>Methanohalophilus halophilus</i>), B5IH56 (<i>Aciduliprofundum boonei</i>), O93627 (<i>Thermococcus kodakarensis</i>), J1ANE7 (<i>Methanofollis liminatans</i>), and Q2FSY4 (<i>Methanospirillum hungatei</i>). The sequences for representative Form I' homologs are presented in this study (Supplementary Data 1). b, Overlay of amino acid residues 408-458 of <i>Syn6301</i> Rubisco (tan) with residues 415-453 of <i>P. breve</i> Rubisco (blue) depicting the unique RbcL C-terminal extension found in Form I enzymes, but not in Rubisco homologs that do not possess RbcS. Residues R428, N429, and E430 of <i>Syn6301</i> RbcL contact residues N29 and Y32 at the interface of <i>Syn6301</i> RbcS (purple).</p>
Extended Data Fig. 5	Negative-staining electron microscopy 2D images of <i>P. breve</i> Rubisco.	Shih_EDFig5.eps	Images reflect the highest resolution data collected with activated <i>P. breve</i> Rubisco in phosphate buffer. The experiment was performed once ($n = 1$).
Extended Data Fig. 6	Extended SEC-SAXS-MALS data.	Shih_EDFig6.eps	Experimental SAXS profiles (black) of <i>P. breve</i> Rubisco in the absence (purple) or presence (blue) of bound 2CABP is

			displayed with the calculated scattering from the atomistic models shown in Fig. 3c. Inset shows the Guinier plot of experimental SAXS profiles with the linear fit in the $q \times R_g < 1.6$ limits.
Extended Data Fig. 7	Amino acid sequence alignment of <i>Syn6301</i> RbcL and <i>P. breve</i> RbcL.	Shih_EDFig7.eps	a , Structure-based sequence alignment was originally made using PROMALS3D ⁵³ using 1RBL and 6URA structures, then aligned with the complete RbcL sequences using MAFFT ⁵⁴ . Darker shades indicate higher sequence conservation between amino acids. <i>Syn6301</i> and <i>P. breve</i> RbcL residues involved in dimer-dimer interactions are highlighted in green and blue, respectively. <i>Syn6301</i> RbcL residues involved in RbcS contacts are annotated with red stars. All contact residues were identified using CCP4 CONTACTS ⁵⁵ . b-c , Cross-section depictions of 1RBL, without RbcS, and <i>P. breve</i> Rubisco highlighting dimer-dimer interactions as in panel a. d , Map of <i>Syn6301</i> RbcL residues involved in RbcS interactions, highlighted in red as in panel a.
Extended Data Fig. 8	Mutating key amino acid residues at the dimer-dimer interface of <i>P. breve</i> Rubisco disrupts octameric oligomeric assembly.	Shih_EDFig8.eps	Native PAGE gel of recombinant WT, K150A, D161A, W165A, D220A, and Y224A <i>P. breve</i> Rubisco. Native Mark protein ladder denoted by "M". Site directed mutants destabilize the interface between RbcL dimers leading to break down of higher-order (i.e., L_8) oligomers into Rubisco species with variable oligomeric state and conformations, which results in a variety of lower molecular weight migration patterns within the Native PAGE gel. Experiment was performed once ($n = 1$).
Extended Data Fig. 9	Site directed mutagenesis of <i>Syn6301</i> dimer-dimer interface residues imparts marginal stability in the absence of RbcS.	Shih_EDFig9.eps	a , Protein thermal shift data displaying the mean fluorescent signal collected from four separate trials for WT <i>Syn6301</i> RbcL, three separate mutant proteins, L158W, V154D, D349R and a combined four mutant protein, 4SDM (L158W, V154D, F217Y, and D349R). Mutations were designed to reflect homologous dimer-dimer interface residues present in <i>P. breve</i> Rubisco. The peaks corresponding to thermal denaturation of L_8 quaternary structure are boxed, and analysis statistics are presented in the below table. T_m values represent the mean

			and standard deviation of n number of experiments conducted with the same protein sample. Two-tailed P-values for unpaired t test with Welch's corrections are reported in the last column using WT <i>Syn6301</i> RbcL as the reference comparison. n = number of technical replicates conducted in experiment. ns = not significant. ** $P < 0.005$, *** $P < 0.0005$. b , Native gel of purified recombinant WT and mutant <i>Syn6301</i> proteins used in experiment.
--	--	--	---

4 **2. Supplementary Information:**

5

6 **A. Flat Files**

7

8

Item	Present?	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. The extension must be .pdf	A brief, numerical description of file contents. <i>i.e.: Supplementary Figures 1-4, Supplementary Discussion, and Supplementary Tables 1-4.</i>
Supplementary Information	Yes	Supplementary Information.pdf	Supplementary Note, Supplementary Tables 1-3.
Reporting Summary	Yes	nr-reporting-summary_PMS_20200710.pdf	

9

10

11 **B. Additional Supplementary Files**

12

13

Type	Number If there are multiple files of the same type this should be the numerical indicator. i.e. "1" for Video 1, "2" for Video 2, etc.	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. i.e.: <i>Smith_Supplementary Video 1.mov</i>	Legend or Descriptive Caption Describe the contents of the file
Supplementary Data	1	Supp. Data 1.txt	Fasta file containing protein amino acid sequences for Form I' enzymes identified from MAGs.
Supplementary Data	2	Supp. Data 2.xlsx	Representative MAG genbank scaffolds.
Supplementary Data	3	Supp. Data 3.xlsx	Site-directed mutagenesis primers and synthesized candidate Form I' <i>rbcl</i> gene sequences.

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52

Novel Bacterial Clade Reveals Origin of Form I Rubisco

Douglas M. Banda^{1,2}, Jose H. Pereira^{3,4,†}, Albert K. Liu^{1,2,†}, Douglas J. Orr⁵, Michal Hammel⁴,
Christine He⁶, Martin A.J. Parry⁵, Elizabete Carmo-Silva⁵, Paul D. Adams^{3,4}, Jillian F.
Banfield^{6,7,8,9,*} & Patrick M. Shih^{1,2,10,11,*}

¹ Department of Plant Biology, University of California, Davis, CA, USA.

² Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA.

³ Technology Division, Joint BioEnergy Institute, Emeryville, CA, USA.

⁴ Molecular Biophysics and Integrated Bioimaging Division, Lawrence Berkeley National Laboratory, Berkeley, CA, 94720, USA.

⁵ Lancaster Environment Centre, Lancaster University, Lancaster, LA1 4YQ, UK.

⁶ Department of Earth and Planetary Science, University of California, Berkeley, Berkeley, CA, USA.

⁷ Department of Environmental Science, Policy, and Management, University of California, Berkeley, Berkeley, CA, USA.

⁸ Innovative Genomics Institute, University of California, Berkeley, Berkeley, CA, USA.

⁹ Chan Zuckerberg Biohub, San Francisco, CA, USA.

¹⁰ Feedstocks Division, Joint BioEnergy Institute, Emeryville, CA, USA.

¹¹ Genome Center, University of California, Davis, Davis, CA, USA.

† Contributed equally

*Correspondence to: pmsih@ucdavis.edu or jbanfield@berkeley.edu

Abstract

Rubisco sustains the biosphere through the fixation of CO₂ into biomass. In plants and cyanobacteria, Form I Rubisco is structurally comprised of large and small subunits, whereas all other Rubisco Forms lack small subunits. Thus, the rise of the Form I complex through the innovation of small subunits represents a key, yet poorly understood, transition in Rubisco's evolution. Through metagenomic analyses, we discovered a previously uncharacterized clade sister to Form I Rubisco that evolved without small subunits. This clade diverged prior to the evolution of cyanobacteria and the origin of the small subunit; thus, it provides a unique reference point to advance our understanding of Form I Rubisco evolution. Structural and kinetic data presented here reveal how a proto-Form I Rubisco assembled and functioned without the

53 structural stability imparted from small subunits. Our findings provide insight into a key
54 evolutionary transition of the most abundant enzyme on Earth and the predominant entry point
55 for nearly all global organic carbon.

56 **Main Text:**

57 Of all known enzymes, few have been more integral in linking the evolution of life with
58 the geochemical cycles of our planet than Rubisco (D-ribulose 1,5-bisphosphate
59 carboxylase/oxygenase)¹. Rubisco sources nearly all organic carbon to the biosphere through the
60 fixation of atmospheric CO₂ with ribulose 1,5-bisphosphate (RuBP) into biomass, thus sustaining
61 our entire food supply. Rubisco also possesses competing oxygenase activity, which is thought to
62 be a vestige of its evolution in a young, oxygen-depleted atmosphere; yet it has co-evolved with
63 Rubisco's carboxylase activity over billions of years. Although there are several distinct Forms
64 of Rubisco found across all three domains of life^{2,3}, the vast majority of carbon fixation on Earth
65 is driven specifically by Form I Rubisco (found in plants, cyanobacteria, algae, and select
66 bacteria phyla); thus, the evolution of this unique Form of Rubisco has profoundly shaped the
67 trajectory of our planet.

68 Structurally, all Forms of Rubisco are composed of at least two large subunits (RbcL, ~50
69 kDa) which assemble head-to-tail as catalytically active dimers. From this rudimentary dimeric
70 scaffold (found in Form II and III homologs), Rubisco has evolved to function in higher-order
71 structures of large subunits including hexamers (Form II), octamers (Form I), and decamers
72 (Form III). Form I homologs, however, are structurally unique from their divergent Form II and
73 Form III counterparts due to the presence of additional small subunits (RbcS, ~13-17 kDa),
74 which cap either end of a central octameric RbcL assembly to form a hexadecameric (L₈S₈)
75 holoenzyme. Thus, understanding the origins of RbcS is part and parcel to investigating the
76 evolution of Form I Rubisco.

77 Although not in direct participation with the active site, RbcS is accepted as an
78 indispensable structural component of Form I Rubisco⁴⁻⁶. For example, cyanobacterial Rubisco
79 from *Synechococcus* sp. strain PCC 6301 (*Syn6301*) retains approximately 1% of its carboxylase
80 activity in the absence of RbcS⁴, suggesting that active site structural integrity is compromised.
81 Furthermore, Form I Rubisco from *Rhodobacter sphaeroides* relies on RbcS to correctly arrange
82 RbcL geometry for proper activity⁷, and plant Rubisco RbcL form insoluble aggregates when
83 expressed without RbcS *in planta*^{8,9}. Despite its demonstrated significance in Rubisco catalysis,
84 the structural role RbcS has played in the evolution of Form I Rubisco has long been debated⁶.
85 This quandary, in part, stems from the fact that we have not identified Form I Rubisco that
86 function without small subunits. Thus, the identification and characterization of a small subunit-
87 less Form I Rubisco would provide the necessary reference point from which to better examine
88 the evolutionary role of RbcS. Towards this end, we searched metagenomic datasets for a
89 “missing link” between the evolution of the Form I clade and all other Forms of Rubisco. Here,
90 we report the discovery of a Form I Rubisco with octameric oligomeric assembly that evolved
91 without RbcS, thus challenging our understanding of the structural properties that govern the
92 activity of the most prominent Form of Rubisco.

93 **Discovery of Form I Rubisco that lack small subunits**

94 To determine whether Form I Rubisco lacking small subunits occur in nature, we
95 analyzed a diverse set of metagenomic datasets derived from environmental communities of
96 largely uncultivated bacteria. Our analyses specifically targeted the identification of
97 uncharacterized bacterial *rbcL* genes, which are usually found within operons encoding other key
98 Calvin-Benson-Bassham (CBB) cycle genes¹⁰. Through this process, we identified 24 *rbcL* genes
99 with gene products that share high sequence homology (52-61%) to known Form I Rubisco.

100 Notably, the average amino acid sequence identity between different Forms of Rubisco is
101 approximately 30%, thus it is possible that the identified *rbcL* genes were either within the Form
102 I clade, or within a close sister clade². Further phylogenetic analyses confirmed that the newly
103 discovered *rbcL* sequences indeed form a monophyletic clade sister to Form I Rubisco. Given
104 the unique phylogenetic proximity to Form I, we named this new clade Form I' to distinguish it
105 from all other *bona fide* Forms of Rubisco (**Fig. 1a**).

106 Where metagenome-assembled contigs were of sufficient length to reveal the genomic
107 context surrounding Form I' *rbcL* genes, all identified operons encoded other CBB cycle genes,
108 including the only other CBB cycle-specific gene, phosphoribulokinase (PRK) (**Fig. 1b**). Closer
109 inspection of metagenome-assembled genomes (MAGs) containing Form I' *rbcL* genes indicated
110 the absence of *rbcS* upstream or downstream of *rbcL*. Notably, bacterial Form I *rbcL* and *rbcS*
111 genes are always found within one or two genes of another in operons^{11,12}. Given that Form I'
112 Rubisco lacks RbcS similar to all other non-Form I Rubisco found in various bacteria and
113 archaea, this suggests that the Form I' clade represents a distinct Form of Rubisco that likely
114 diverged from the Form I clade prior to the origin of RbcS.

115 Surprisingly, all Form I' genes identified from MAGs were found exclusively in a single
116 order of the Chloroflexi phylum, Anaerolineales (**Extended Data Figs. 1 and 2**). Although
117 Chloroflexi are commonly known for their phototrophic members in the order Chloroflexales,
118 the majority of the phylum is composed of phenotypically-diverse filamentous bacteria that are
119 non-phototrophic, such as the Anaerolineales¹³. Of the known phototrophic examples of
120 Chloroflexi within the order Chloroflexales, most perform carbon fixation via the 3-
121 hydroxypropionate bicycle (*e.g. Chloroflexus sp.*), or with Form I Rubisco via the CBB cycle
122 (*e.g. Oscillochloris trichoides, Chlorothrix halophila, and Kouleothrix aurantiaca*)¹⁴. Form I'-

123 containing MAGs were not found to contain characteristic 3HP bicycle genes such as propionyl-
124 CoA synthetase, malonyl-CoA reductase/3-hydroxypropionate dehydrogenase, and malonyl-
125 CoA/succinyl-CoA reductase, suggesting that the bacteria consortium from which MAGs were
126 derived use the CBB cycle for autotrophy. Though some examples of phototrophic Chloroflexi
127 have recently been described in clades sister to the Anaerolinea (e.g. the class-level clade
128 *Candidatus Thermofonsia*)¹⁵, none possessed carbon fixation pathway genes and were presumed
129 to be photoheterotrophic. Studies have demonstrated that phototrophy within Chloroflexi may be
130 driven by horizontal gene transfer^{15,16}; however, the tight phylogenetic distribution of Form I'
131 genes within the order Anaerolineales suggests otherwise, albeit future studies may reveal
132 genomes outside of Anaerolineales that possess Form I' genes.

133 **Form I' Rubisco is functional despite lack of small subunits**

134 To characterize genes discovered from MAGs, representative Form I' Rubisco homologs
135 were recombinantly expressed and purified (**Extended Data Fig. 3a**) from *E. coli*
136 overexpressing the bacterial chaperonin system GroEL-GroES (homologous to Cpn60-Cpn10/20
137 in plants), a necessary component of Rubisco biogenesis^{17,18}. The assembly of hexadecameric
138 Form I homologs in cyanobacteria and plants require auxiliary chaperones such as RbcX and
139 RafI, which aid in the stabilization of the octameric RbcL core before the addition of small
140 subunits^{19,20}. Other Form I homologs, however, do not require homologous assembly factors but
141 instead rely on RbcS for efficient assembly, which has been demonstrated for Rubisco from the
142 photosynthetic proteobacterium *Rhodobacter sphaeroides*⁷. RbcX was not found in Form I'
143 containing MAGs (**Fig. 1b**). Consistent with this finding, all Form I' sequences do not possess
144 the C-terminal binding domain for RbcX^{8,21} (**Extended Data Fig. 4a**). Furthermore, Form I'
145 homologs identified to date do not possess small subunits, precluding the necessity of chaperones

146 involved in the assembly of hexadecameric Rubisco¹⁹. Some archaeal Rubisco possess an extra
147 C-terminal domain that is proposed to aid in RbcL core assembly²², but this unique insertion is
148 not found within the described representative homologs of the Form I' clade (**Extended Data**
149 **Fig. 4a**). Notably, *Syn6301* Rubisco expressed in *E. coli* makes up ~1-2% of the total soluble
150 protein, but this number improves to ~6% with the associated overexpression of GroEL/ES²³. In
151 comparison, Rubisco from *R. sphaeroides* comprises ~16% of the total soluble protein when
152 heterologously expressed in *E. coli*, which jumps to 33% with the overexpression of GroEL/ES⁷.
153 With the system outlined in this work, Form I' Rubisco was found to express at ~7-8% of the
154 total soluble protein in BL21(DE3) *E. coli*, which improves to ~14-15% when overexpressed
155 with GroEL/ES (**Extended Data Fig. 3b**). Currently, it is unknown whether the expression
156 levels of Form I' Rubisco in *E. coli* are intrinsic to its amino acid sequence alone, or if auxiliary
157 chaperone factors are necessary for higher expression. Though the Chloroflexi from which these
158 sequences are derived may possess a unique assembly factor that aids in Rubisco biogenesis, no
159 such protein was identified from the metagenomic datasets presented in this work.

160 To assess the catalytic activity of a representative Form I' homolog, we performed
161 detailed enzyme kinetic measurements on Form I' Rubisco from the mesophilic Chloroflexi
162 species "*Candidatus Promineofilum breve*" (*P. breve*) using the method of Parry *et al.*²⁴. At
163 saturating substrate concentrations, Rubisco proteins exhibit maximal rates of catalysis (V_C and
164 V_O for carboxylation and oxygenation, respectively), generally at the expense of the
165 concentration of substrate necessary to achieve a maximal rate (represented by the Michaelis
166 constants K_C and K_O for carboxylation and oxygenation, respectively, which can be considered
167 conceptually as pseudo-dissociation constants for the binding of either CO_2 or O_2)^{25,26}.

168 P. breve Rubisco demonstrated relatively slow V_C and about average K_C when compared
169 to the reported measurements of Form I enzymes at 25 °C²⁶ (**Table 1, Fig. 2**). Conversely, the
170 enzyme demonstrated slightly above average V_O and below average K_O . This is consistent with
171 the discovery of the Form I' clade within the order Anaerolineales, which is typically comprised
172 of obligate anaerobes²⁷, although genomic signatures of aerobic respiration have recently been
173 discovered in some examples of Anaerolineae^{28,29}. Together, these kinetic parameters culminated
174 in a specificity for CO₂ over O₂ (represented by the specificity factor (S_{CO}), a measure of the
175 catalytic efficiency of the carboxylation reaction over the oxygenation reaction) that is lower
176 relative to values reported for Form I enzymes, but higher than Form II and Form III homologs
177 (**Supplementary Table 1**). It is unclear at this time whether the high oxygenase specificity of P.
178 breve Rubisco is linked to the absence of RbcS. Notably, Form I' and Form I Rubisco lineages
179 diverged before the evolution of cyanobacteria suggesting that Form I' enzymes may have
180 evolved in anaerobic conditions.

181 **Form I' Rubisco is octameric, reminiscent of Form I Rubisco**

182 The Form I clade is structurally characterized by two features distinct from other Forms
183 of Rubisco: 1) the presence of RbcS, and 2) the oligomeric assembly of RbcL into octamers.
184 Given the close phylogenetic placement to the Form I clade, we hypothesized that Form I'
185 homologs may possess octameric oligomeric assembly of RbcL, which has not been previously
186 observed for Rubisco in nature. Size exclusion chromatography (SEC) and non-denaturing
187 PAGE analyses revealed that recombinant P. breve RbcL dimers (~100-110 kDa) oligomerized
188 into a higher-order structure (**Fig. 3d**). Previous studies have demonstrated that the addition of
189 the Rubisco-specific transition-state analog, 2-carboxyarabinitol 1,5-bisphosphate (2CABP), may
190 influence the oligomeric state of the enzyme³⁰. Incubation of magnesium-bound and CO₂-

191 activated *P. breve* Rubisco with 2CABP resulted in an observed structural compaction, evident
192 from both later elution in SEC traces, as well as slower migration in non-denaturing gels (**Fig. 3**).

193 To more rigorously characterize the solution-state oligomeric assembly of *P. breve*
194 Rubisco, we performed SEC coupled to small-angle X-ray scattering (SAXS) and multiangle
195 light scattering (MALS) (SEC-SAXS-MALS) experiments³¹ with activated *P. breve* Rubisco in
196 the presence or absence of 2CABP. Protein molecular weights determined by MALS (~400-440
197 kDa) supported the oligomerization of *P. breve* Rubisco as an L₈ complex (theoretical octamer
198 M.W. ~409 kDa), similar to the octameric assembly of RbcL in related Form I enzymes (**Fig. 3**).
199 These observations were corroborated by negative-staining electron microscopy images
200 (**Extended Data Fig. 5**). Experimentally determined pair-distribution, or P(r), functions
201 displayed significant broadening and elongation of *P. breve* Rubisco in the absence of 2CABP
202 relative to the 2CABP-bound protein (**Fig. 3b**). This observation agrees well with the larger
203 radius of gyration (R_g) values of the 2CABP-bound (R_g ~ 46.8 ± 0.4 Å) versus unbound (R_g ~
204 45.0 ± 0.5 Å) protein.

205 In the absence of substrate, Form I Rubisco proteins exist in an “open” conformation that
206 is structurally characterized, in part, by an extended C-terminal domain that is disordered and
207 positioned away from the active site³². Upon active site binding of RuBP, the extended C-
208 terminal domain flips down over the active site with Loop 6 to produce a compact “closed”
209 conformation primed for catalysis. In order to account for observed differences in the radius of
210 gyration between 2CABP-bound and unbound structures, we generated theoretical SAXS data
211 from computational models of octameric *P. breve* Rubisco either in a compact “closed” state
212 (i.e., bound to 2CABP) or an “open” state with disordered C-terminal domains (**Fig. 3c**). Indeed,
213 theoretical SAXS data produced from these models matched well with the experimentally

214 determined P(r) functions (**Fig. 3b**) and SAXS profiles (**Extended Data Fig. 6, Supplementary**
215 **Table 2**, $\chi^2 = 1.8$ and 1.4 for closed and open models, respectively).

216 Overall, the combination of SEC-SAXS-MALS and electron microscopy experiments
217 support an L₈ oligomerization of Form I' Rubisco reminiscent of the L₈S₈ Form I Rubisco.
218 Because no other Form of Rubisco has been convincingly demonstrated to express as octamers in
219 nature (see **Supplementary Note**), the most parsimonious history consistent with our data
220 suggests that the common ancestor of Form I and Form I' clades evolved an octameric core
221 assembly prior to the evolution of RbcS.

222 **Form I' Rubisco structure yields insight into Form I Rubisco evolution**

223 To obtain higher molecular resolution of *P. breve* Rubisco, we solved a 2.2 Å crystal
224 structure of the activated enzyme in complex with 2CABP (**Fig. 4, Supplementary Table 3**).
225 Superposition of *P. breve* RbcL onto the structure of *Syn6301* L₈S₈ Rubisco (PDB ID: 1RBL)³³
226 resulted in a C α RMSD of 0.68 Å between 435 pruned atom pairs (97.5% of *P. breve* RbcL
227 amino acid sequence), with a Q-score of 0.87³⁴. As with all other *bona fide* Rubisco, all key
228 active site residues^{35,36} were positioned in an $\alpha\beta$ -barrel (TIM-barrel) domain (residues 158-405).

229 Many of the characteristic Form I hydrophobic RbcL residues at the interface of large and
230 small subunits³⁷ were either functionally substituted on the surface of *P. breve* Rubisco (~31%)
231 or completely absent (~4%), based on sequence homology to *Syn6301* RbcL (**Extended data**
232 **Fig. 7**). RbcL surface residues between the two structures displayed strikingly similar
233 electrostatic characteristics (**Fig. 4**), which was unexpected given that *P. breve* Rubisco had not
234 evolved to interact with RbcS, unlike its closely related *Syn6301* homolog. Because of this
235 observation and the close phylogenetic relationship between the Form I and Form I' clades, a

236 competing hypothesis is that Form I' evolved *from* Form I homologs and subsequently lost
237 RbcS, as opposed to the hypothesis that Form I' and Form I Rubisco diverged from a common
238 ancestor. To explore this further, we investigated the observation that Form I homologs possess
239 an RbcL “C-terminal extension” (residues 430-442 of *Syn6301* Rubisco, **Extended Data Fig.**
240 **4a**) not found in Rubisco that lack RbcS (*i.e.*, all other Forms of Rubisco). This unique C-
241 terminal extension has evolved in Form I lineages to stabilize key RbcL interactions with RbcS³⁸
242 (**Extended Data Fig. 4b**). The Form I' enzymes identified in this study do not possess this
243 unique C-terminal extension important for RbcS interactions, supporting the hypothesis that
244 Form I' and Form I Rubisco diverged from a common ancestor. This is in accordance with the
245 parsimonious observation that all non-Form I Rubisco lack RbcS, suggesting that the common
246 ancestor to both Form I and Form I' clades most likely lacked RbcS.

247 In the absence of RbcS, we hypothesized that *P. breve* Rubisco must possess fortified
248 interactions at the RbcL dimer-dimer interface to support octameric assembly. Indeed, *P. breve*
249 Rubisco possesses an extensive network of hydrogen bonds and salt bridges at the interdimer
250 interface that is not present in *Syn6301* Rubisco (**Fig. 5a**). Site-directed mutagenesis of key
251 amino acid residues within this network (Lys150, Asp161, Trp165, Asp220, and Tyr224) to
252 alanine abolished *P. breve* Rubisco's octameric assembly (**Extended Data Fig. 8**),
253 demonstrating their importance in maintaining holoenzyme stability in the absence of RbcS.
254 Notably, homologous amino acid positions to Asp161, Trp165, and Tyr224 within *Syn6301*
255 (Val154, Leu158, and Phe217, respectively) are incapable of forming a similar electrostatic
256 network due to their side-chain physicochemical properties, necessitating interactions with RbcS
257 for complex stability (**Extended Data Fig. 7**).

258 To quantitatively evaluate how subunit interactions within *Syn6301* and *P. breve* Rubisco
259 affect the thermal stability of the complex quaternary structure, we employed a protein thermal
260 shift assay³⁹ (**Fig. 5b**). In the absence of RbcS, *Syn6301* Rubisco displayed a two-phase melting
261 profile; the first phase ($T_m = 58.6 \pm 0.2$ °C) resulting from quaternary structure disassembly (*i.e.*,
262 the dissociation of octamers into dimers), and the second phase ($T_m = 70.6 \pm 0.2$ °C)
263 corresponding to the simultaneous denaturation of RbcL dimers and RbcL secondary structure⁴⁰.
264 In the presence of RbcS, *Syn6301* Rubisco was significantly stabilized such that L_8S_8
265 disassembly was shifted by more than 15 °C relative to *Syn6301* L_8 ($T_m = 75.5 \pm 0.1$ °C).
266 Interestingly, *P. breve* Rubisco disassembly displayed a modest increase in T_m (82.6 ± 0.1 °C)
267 relative to *Syn6301* L_8S_8 , but a significant increase when compared to the T_m measured for
268 *Syn6301* in the absence of RbcS, consistent with the predicted added stability due to interdimer
269 interface interactions. To stabilize *Syn6301* in the absence of RbcS, we mutated RbcL residues
270 known to interact with RbcS to mimic part of the electrostatic network stabilizing *P. breve*
271 oligomeric assembly (**Extended Data Fig. 9**). This effort yielded modest improvement in
272 stability, highlighting the complexity of forming octamers in the absence of RbcS.

273 **Discussion**

274

275 Accrued evidence from investigations into the evolutionary adaptability of proteins supports a
276 common trend: the catalytic promiscuity of an enzyme is inversely proportional to its
277 conformational stability⁴¹⁻⁴³. In line with previous observations⁶, the data presented in this work
278 suggests that the innovation of a distinct structural subunit (*i.e.*, RbcS) imparted structural
279 stability to Rubisco during the evolution of its carboxylase and oxygenase activities towards
280 “Pareto optimality”⁴⁴. Form I’ Rubisco from *Ca. P. breve* demonstrated high oxygenase activity

281 and lower specificity when compared to Form I homologs (**Fig. 2, Table 1**), likely stemming
282 from the anaerobic lifestyle of the Anaerolineales order of Chloroflexi from which sequences
283 were discovered. Furthermore, the divergence of Form I' and Form I Rubisco from a common
284 ancestor predates the origin of cyanobacteria; thus it is likely that Form I' Rubisco originated
285 during the Archean Eon when atmospheric oxygen was scarce. Collectively, these observations
286 suggest that the appearance of RbcS and the evolutionary transition from L₈ to L₈S₈ may have
287 been an evolutionary response to the rise of oxygen ~2.4 Ga. This environmental transition may
288 have provided a strong selective pressure to L₈-containing autotrophs (e.g., stem-group
289 cyanobacteria) that necessitated a tradeoff between conformational rigidity (i.e., enhanced
290 interactions at the dimer-dimer interface of octameric Rubisco) and active site plasticity. The
291 selective pressure driving this tradeoff likely stemmed from an increased demand for improved
292 carboxylation activity to drive flux through carbon metabolism during a rapidly changing
293 paleoatmosphere^{45,46}. To evolve this conformational dynamism while maintaining an optimized
294 oligomeric state (i.e., L₈), we posit that RbcS evolved to facilitate the adaptive evolution of
295 Rubisco's catalytic activity, effectively buffering the cost of destabilizing mutations and
296 allowing the sampling of higher genetic diversity during the random walk through sequence
297 space.

298 In addition to the evolutionary insight gleaned from this work, the discovery of the Form
299 I' clade from MAG's may offer alternative means to explore Rubisco engineering efforts in
300 plants. Notably, Form I Rubisco has long been recalcitrant to directed evolution experiments for
301 improved carbon fixation, with notable exceptions⁴⁷, in part due to challenges associated with
302 effectively exploring the sequence space of two genes (i.e., RbcL and RbcS) simultaneously;
303 thus, the absence of RbcS in Form I' enzymes may streamline such future efforts. Overall,

304 performing directed evolution experiments^{47,48} with *P. breve* Rubisco in conjunction with the
305 continued characterization of the Form I' clade will offer novel opportunities to advance our
306 understanding of Rubisco evolution.

307
308
309

310

311

312

313 **Methods**

314

315 **Metagenomic and phylogenetic analysis.** All metagenomes were sequenced using 150 bp,
316 paired-end Illumina reads and assembled into scaffolds using either IDBA-UD or Megahit.
317 Scaffolds were binned based on GC content, coverage, presence of ribosomal proteins,
318 presence/copies of single copy genes, tetranucleotide frequency, and patterns of coverage across
319 samples. Bins were manually curated, dereplicated, and filtered for completeness and
320 contamination. Genes were predicted using hidden Markov models (HMMs) based on Pfam,
321 TIGRfams, KEGG, and custom databases. Phylogeny of bins containing Rubisco genes was
322 identified using overall scaffold gene content as well as maximum likelihood phylogenetic trees
323 of 16 concatenated ribosomal protein sequences. Rubisco gene sequences were dereplicated at
324 97% amino acid identity using CD-Hit, aligned using MAFFT (default parameters), and columns
325 with >95% gaps were removed using TrimAI. A maximum-likelihood phylogenetic tree was
326 constructed using RAxML-HPC BlackBox (v. 8.2.10) as implemented on ciperes.org (default
327 parameters with LG model). To construct Figure 1A, branches with bootstrap values of <0.65
328 were collapsed. Both the alignment file and the tree file with bootstrap values are available on
329 figshare (DOI: 10.6084/m9.figshare.9980630).

330

331 **Plasmids, cloning, and site-directed mutagenesis.** Representative Form I' *rbcL* genes were
332 synthesized by Twist Biosciences (San Francisco, CA) (sequences available as supplementary
333 data) and cloned into a pET28 vector with an N-terminal His₁₄-bdSUMO tag⁵⁹. Plasmids
334 pSF1389⁵⁹, pET11a-Syn6301-*rbcLS*, pET11A-Syn6301-*rbcL*, pBADES/*EL*, and pG-KJE8²¹ were
335 gifts. Site-directed mutagenesis (SDM) was conducted using an Agilent QuikChange SDM kit
336 and standard procedures. Primers were designed using the Agilent QuikChange Primer Design
337 tool (available as supplementary data).

338

339 **Expression and purification of recombinant proteins.** *Brachypodium distachyon* SUMO-
340 specific protease (bdSENP1) was prepared by transforming pSF1389 into chemically competent

341 BL21 DE3 Star *E. coli* cells (Macrolab, QB3-Berkeley, CA). Cells were grown to mid-log phase
342 at 37 °C (OD₆₀₀ ~ 0.6) and induced with 0.3 mM IPTG for 3 hours. Cells were resuspended in
343 pH 7.0 Lysis Buffer (20 mM sodium phosphate, 300 mM NaCl, 10 mM imidazole, 5% glycerol,
344 2 mM MgCl₂) with ~5 mM PMSF and subject to a freeze-thaw cycle before lysis by use of a
345 Microfluidizer high pressure homogenizer (Microfluidics, Westwood, MA), and centrifugation
346 (15,000 RCF, 20 min). Soluble protein was 0.2/0.8 µm filtered and applied to Ni-NTA Resin
347 (Thermo Fisher, Waltham, MA) and batch bound according to the manufacturer's protocols.
348 Columns were washed thoroughly before elution. TEV protease (MilliporeSigma, Burlington,
349 MA) was added to the eluted fraction according to the manufacturer's suggestion and rocked
350 gently overnight at 4 °C to facilitate His tag cleavage. The flow-through from TEV protease
351 reactions was buffer exchanged into pH 7.0 Ni Equilibration buffer (20 mM sodium phosphate,
352 300 mM NaCl, 10 mM imidazole, 10% glycerol) and passed over Ni-NTA resin again to
353 separate cleaved His tag from the target protein. bdSENP1-containing flow-through was
354 analyzed by SDS-PAGE for purity and stored at -80 °C in storage buffer (20 mM sodium
355 phosphate pH 7.0, 300 mM NaCl, 1 mM DTT, 10% glycerol).

356 *P. breve* Rubisco was prepared by co-transforming plasmids containing His₁₄-bdSUMO-
357 tagged *P. breve* RbcL into chemically competent BL21 DE3 Star *E. coli* with pBADES/EL
358 plasmid. Cells were grown to mid-log phase at 30 °C (OD₆₀₀ ~ 0.6) and overexpression of
359 GroEL/ES was induced by the addition of 0.2% w/v arabinose, and further incubation for 2
360 hours. Cells were resuspended in fresh LB media (without arabinose) with 300 mM NaCl and 20
361 mM L-proline and shaken for 16 hours at 16 °C. Pelleted cells were resuspended in pH 8.0 Lysis
362 Buffer (20 mM sodium phosphate, 300 mM NaCl, 10 mM imidazole, 5% glycerol, 2 mM
363 MgCl₂) with ~5 mM PMSF and subject to a freeze-thaw cycle at -80 °C before lysis by use of a
364 Microfluidizer high pressure homogenizer. The soluble fraction was collected by centrifugation
365 (15,000 RCF, 20 min) and 0.2/0.8 µm filtered. Clarified cell lysate was batch-bound to pre-
366 equilibrated Ni-NTA resin as described above. Columns were washed thoroughly before
367 resuspension in bdSENP1 Reaction Buffer (20 mM sodium phosphate pH 8.0, 300 mM NaCl, 1
368 mM DTT, 10% glycerol). Purified bdSENP1 was added to resuspended columns and rocked
369 gently overnight at 4 °C to facilitate cleavage of the His₁₄-bdSUMO tag from the target protein.
370 Flow-through from the bdSENP1 reaction was applied to a 5 mL HiTrap Q FF column
371 equilibrated in Q Buffer A (100 mM HEPES pH 8.0). Protein was eluted off the column over a
372 linear NaCl gradient from 5 mM to 1 M. Eluted fractions were analyzed by SDS-PAGE prior to
373 concentration and separation by size exclusion chromatography using a Superose 6 Increase
374 10/300 GL column (GE Healthcare Life Sciences, Marlborough, MA) equilibrated in SEC Buffer
375 (50 mM sodium phosphate pH 8.0, 300 mM NaCl, 25 mM MgCl₂, 1 mM DTT, 5 mM NaHCO₃).
376 Eluted SEC fractions were analyzed by SDS-PAGE and Native PAGE for Rubisco content and
377 purity. Samples were stored in 20 mM sodium phosphate pH 8.0, 150 mM NaCl, 10 mM MgCl₂,
378 10 mM NaHCO₃ at -80 °C.

379 *Syn6301 RbcLS* was prepared in a similar fashion to previous reports^{21,40}. Plasmids
380 *Syn6301-rbcLS-pET11A* and *pBADES/EL* were co-transformed into BL21 DE3 Star *E. coli*
381 cells. Cells were grown to mid-log phase at 30 °C (OD₆₀₀ ~ 0.6) and overexpression of
382 GroEL/ES was induced by 0.4% w/v arabinose for 1.5 hours. Cells were resuspended in fresh
383 media (without arabinose) and induced with 1 mM IPTG for 16 hours at 16 °C. Cells were lysed
384 by using a Microfluidizer high pressure homogenizer and centrifuged (15,000 RCF, 20 minutes).
385 Soluble protein from whole-cell lysate was 0.2/0.8 µm filtered and subject to ammonium sulfate
386 precipitation at the 30-40% cut (where the protein is soluble at 30% w/v ammonium sulfate, but
387 precipitates at 40% saturation. Precipitated protein was resuspended in pH 8.0 Lysis Buffer,
388 desalted, and applied to a MonoQ 10/100 GL column (GE Healthcare Life Sciences,
389 Marlborough, MA) equilibrated in Q Buffer A. Protein was eluted off the column over a linear
390 NaCl gradient from 5 mM to 1 M. Eluted fractions were analyzed by SDS-PAGE prior to
391 concentration and size exclusion chromatography as described for *P. breve* Rubisco. Samples
392 were stored in 20 mM sodium phosphate pH 8.0, 150 mM NaCl, 10 mM MgCl₂, 10 mM
393 NaHCO₃ at -80 °C.

394 *Syn6301 RbcL* expressed without RbcS was prepared in a similar fashion to previous
395 reports^{21,40}. Plasmids *Syn6301-rbcL-pET11A* and *pG-KJE8* were co-transformed into BL21 DE3
396 Star *E. coli* cells. Cells were grown to mid-log phase at 30 °C (OD₆₀₀ ~ 0.6) and overexpression
397 of *dnaK/dnaJ/grpE* was induced by 0.4% w/v arabinose for 2 hours. Cells were resuspended in
398 fresh media (without arabinose) and induced with 1 mM IPTG for 16 hours at 16 °C. Cells were
399 lysed and centrifuged as described for *Syn6301 RbcLS*. Soluble protein from whole-cell lysate
400 was subject to ammonium sulfate precipitation at the 50-60% cut. Precipitated protein at 60%
401 saturation was resuspended in lysis buffer and purified via anion exchange and size exclusion
402 chromatography, then stored at -80 °C as described for *Syn6301 RbcLS*.

403 **PAGE analyses.** Rubisco samples were activated with excess NaHCO₃ and incubated with 10-
404 fold molar excess 2-carboxyarabinitol 1,5-bisphosphate (2CABP) as described previously³⁰.
405 2CABP was synthesized according to previously described methods^{60,61}. SDS-PAGE samples
406 were prepared according to standard procedures in Laemmli Sample Buffer (Bio-rad, Hercules,
407 CA) with 2-mercaptoethanol, and heated at 98 °C for 5 minutes, followed by centrifugation in a
408 benchtop centrifuge at maximum speed for 1 minute. Samples were resolved on 12% Mini-
409 PROTEAN® TGX™ precast protein gels (Bio-rad) in 1x Tris/Glycine/SDS buffer (Bio-Rad)
410 and stained in AcquaStain (Bulldog Bio, Portsmouth, NH). Non-denaturing PAGE samples were
411 prepared by mixing protein with Native Sample Buffer (Bio-Rad) at 4 °C. Samples were
412 resolved at 4 °C on 4-15% Mini PROTEAN® TGX™ precast protein gels (Bio-rad) in 1x
413 Tris/Glycine buffer (Bio-Rad) and visualized by staining with AcquaStain.

414
415 **Crystallization, X-ray data collection, and structure determination.** For crystallography, *P.*
416 *breve* Rubisco was prepared as described above, but with a final buffer composition of 100 mM
417 HEPES-OH pH 8.0, 100 mM NaCl, 25 mM MgCl₂, 1 mM DTT, 5 mM NaHCO₃. Samples at 10-

418 15 mg/mL were activated as described above. Samples crystallized in the presence of 2CABP
419 were incubated for 1 hr at ambient temperature in the presence of a 10-fold molar excess of
420 2CABP before setting up crystal trays. *P. breve* Rubisco protein was screened using the
421 crystallization screens: Berkeley Screen⁶², Crystal Screen, SaltRx, PEG/Ion, Index and PEGRx
422 (Hampton Research, Aliso Viejo, CA). The crystals of *P. breve* Rubisco were found in 0.1 M
423 Tris pH 8.0 and 30 % Polyethylene glycol monomethyl ether 2,000 obtained by the sitting-drop
424 vapor-diffusion method with drops consisting of a mixture of 0.2 μ L of protein solution and 0.2
425 μ L of reservoir solution.

426
427 A crystal of *P. breve* Rubisco was placed in a reservoir solution containing 20% (v/v) glycerol,
428 then flash-cooled in liquid nitrogen. The X-ray data sets for *P. breve* Rubisco were collected at
429 the Berkeley Center for Structural Biology beamline 8.2.2 of the Advanced Light Source at
430 Lawrence Berkeley National Laboratory (LBNL). The diffraction data were recorded using an
431 ADSC-Q315r detector. The data sets were processed using the program Xia2⁶³.

432
433 The *P. breve* Rubisco crystal structure was determined by the molecular-replacement
434 method with the program PHASER⁶⁴ within the Phenix suite^{65,66}, using as a search model the
435 structure of a Rubisco from *Thermosynechococcus elongatus* (PDB code 2YBV), which shows
436 57 % sequence identity to the target. The atomic positions obtained from molecular replacement
437 and the resulting electron density maps were used to build the *P. breve* Rubisco structure and
438 initiate crystallographic refinement and model rebuilding. Structure refinement was performed
439 using the phenix.refine program⁶⁶. Translation-libration-screw (TLS) refinement was used, with
440 each protein chain assigned to a separate TLS group. Manual rebuilding using COOT⁶⁷ and the
441 addition of water molecules allowed construction of the final model. The final model of *P. breve*
442 Rubisco has an R factor of 18.8 % and an R_{free} of 22.5 %. Root-mean-square deviation
443 differences from ideal geometries for bond lengths, angles and dihedrals were calculated with
444 Phenix. The stereochemical quality of the final model of *P. breve* Rubisco was assessed by the
445 program MOLPROBITY⁶⁸.

446
447 **Small-angle X-ray-scattering (SAXS) data collection and analysis.** Small-angle X-ray
448 scattering (SAXS) coupled with multi-angle light scattering (MALS) in line with size-exclusion
449 chromatography (SEC) experiments were performed with 50 μ L samples containing 4.6 mg/mL
450 of *P. breve* Rubisco incubated with or without 2CABP prepared in 20 mM HEPES-OH (pH 8.0),
451 300 mM NaCl, 10 mM MgCl_2 , 10 mM NaHCO_3 . SEC-SAXS-MALS data were collected at the
452 ALS beamline 12.3.1 at Lawrence-Berkeley National Lab⁶⁹. The X-ray wavelength was set at
453 $\lambda=1.127 \text{ \AA}$ and the sample-to-detector distance was 2100 mm resulting in scattering vectors (q)
454 ranging from 0.01 \AA^{-1} to 0.4 \AA^{-1} . The scattering vector is defined as $q = 4\pi\sin\theta/\lambda$, where 2θ is the
455 scattering angle. All experiments were performed at 20 °C and the data was processed as
456 described⁷⁰. Briefly, a SAXS flow cell was directly coupled with an online 1260 Infinity HPLC
457 system (Agilent, Santa Clara, CA) using a Shodex KW804 column (Showa Denko, Tokyo,

458 Japan). The column was equilibrated with running buffer (20 mM HEPES-OH (pH 8.0), 300 mM
459 NaCl, 10 mM MgCl₂, 10 mM NaHCO₃) with a flow rate of 0.5 mL/min. 90 μL of sample was
460 separated by SEC, and three second X-ray exposures were collected continuously during a 30
461 min elution. The SAXS frames recorded prior to sample analysis were subtracted from all other
462 frames. The subtracted frames were investigated by radius of gyration (R_g) derived by the
463 Guinier approximation, $I(q) = I(0) \exp(-q^2 R_g^2/3)$ with the limits $qR_g < 1.6$. The elution peak was
464 mapped by comparing integral of ratios to background and R_g relative to the recorded frame
465 using the program SCATTER. Uniform R_g values across an elution peak represent a
466 homogenous assembly. Final merged SAXS profiles, derived by integrating multiple frames
467 across the elution peak, were used for further analysis including Guinier plot which determined
468 aggregation free state. The program SCATTER was used to compute the pair distribution, or
469 P(r), functions presented in Figure 3B. P(r) functions were normalized based on the molecular
470 weight determined by SCATTER using volume of correlation V_c⁴⁹ (Supplementary Table 2).
471 Eluent was subsequently split 3:1 between the SAXS line and a series of UV detectors at 280 and
472 260 nm, a MALS detector, a quasi-elastic light scattering (QELS) detector, and a refractometer
473 detector. MALS experiments were performed using an 18-angle DAWN HELEOS II light
474 scattering detector connected in tandem to an Optilab refractive index concentration detector
475 (Wyatt Technology, Goleta, CA). System normalization and calibration was performed with
476 bovine serum albumin using a 45 μL sample at 10 mg/mL in SEC Buffer and a dn/dc value of
477 0.19. The light scattering experiments were used to perform analytical scale chromatographic
478 separations for M.W. determination of the principal peaks in the SEC analysis. UV, MALS, and
479 differential refractive index data was analyzed using Wyatt ASTRA 7 software to monitor the
480 homogeneity of the sample across the elution peak complementary to the above-mentioned SEC-
481 SAXS signal validation.

482
483 **SAXS modeling.** The atomistic model of P. breve Rubisco in the open conformation was
484 prepared based on the crystal structure of the closed conformation presented in this study by
485 including missing N- and C-terminal residues using the program MODELLER⁷¹. Different
486 extensions and compactions of the unfolded tails were built to screen conformational variability.
487 The experimental SAXS profiles were then compared to theoretical scattering curves generated
488 from these atomistic models using FoXS^{57,58}. Theoretical scattering profiles were used to
489 calculate P(r) functions and further compared to experimental P(r) functions to validate solution
490 state conformations of P. breve Rubisco.

491
492 **Negative-staining electron microscopy.** 3 μL of 1 mg/mL P. breve Rubisco in SEC Buffer were
493 applied to a glow-discharged carbon grid (30 mA, 30 sec) and incubated for 1 min at room
494 temperature. Five drops of 2% uranyl acetate were then sequentially applied and blotted off for
495 negative staining. 50 images were taken on a JEOL 2100F at x40,000 nominal magnification,
496 200 kV, with 1.48 Å/pixel sampling on a DE-20 detector. 4062 particles were selected and 2-D
497 classified using cisTEM.

498

499 **Rubisco activity assays.** Rubisco specificity was determined using the method of Parry *et al.*²⁴,
500 with the exception that the activation buffer included 250 mM NaCl to enhance the solubility of
501 *P. breve* Form I' Rubisco, and pKa of 6.11 was used for calculations. Measurements using *T.*
502 *aestivum* (bread wheat) Rubisco were used for normalization as previously described²⁴, and
503 results from testing with *T. aestivum* Rubisco showed no effect of NaCl in the activation buffer.
504 Purified Rubisco was used to determine catalytic properties as described previously⁷², with the
505 following alterations to protein desalting and activation: an aliquot of concentrated Rubisco was
506 diluted with an activation mix containing 100 mM Bicine-NaOH pH 8.0, 20 mM MgCl₂, 250
507 mM NaCl, 10 mM NaHCO₃, and 1 % (v/v) Plant Protease Inhibitor cocktail (Sigma-Aldrich,
508 UK). This was then incubated on ice for 20 min before used to assay at CO₂ concentrations of
509 20, 40, 60, 120, 280, and 400 μM. These were combined with O₂ concentrations of either 0, 21,
510 40, or 70 % (v/v) to determine K_O . V_O was calculated from measured parameters using the
511 equation $S_{C/O} = (V_C/K_C)/(V_O/K_O)$. V_C was determined using measurements with 0% O₂. An
512 aliquot of the activated protein was used for determination of Rubisco active sites via ¹⁴C-CABP
513 binding using the method of Sharwood *et al.*⁷³ with 250 mM NaCl, instead of the typical 75 mM,
514 in the activation buffer.

515

516 **Protein thermal shift (PTS) assay.** The PTS assay was conducted using a Protein Thermal
517 Shift™ kit (Thermo Fisher, Waltham, MA). Samples were prepared with 1 mg/mL protein in 1x
518 PTS phosphate buffer, and 4x PTS dye in Thermo Fisher MicroAmp Optical 8-Tube Strips.
519 Assay was conducted on an Applied Biosciences QuantStudio 3 RT-PCR machine. The assay
520 consisted of initial cooling and hold at 16 °C for 1 minute, followed by an 0.05°C/s increase to
521 95 °C, and a final hold at 95 °C for 1 minute. Data was analyzed in Protein Thermal Shift™
522 Software.

523

524 **Other software.** Structure-based sequence alignments were conducted using PROMALS3D⁵³
525 and MAFFT⁵⁴. Analyses of protein amino acid contacts and subunit interface thermodynamics
526 were performed using CCP4 CONTACTS⁵⁵, and PISA^{74,75}, respectively. UCSF Chimera⁷⁶ was
527 utilized for the visualization of protein models, generating electrostatic potential maps, and the
528 preparation of manuscript figures.

529

530 **Data availability.** Form I' RbcL amino acid sequences are included as a supplementary file
531 (Supplementary data 1). Sequences used to generate Fig. 1a were uploaded to figshare (DOI:
532 10.6084/m9.figshare.9980630) along with the associated phylogenetic tree. Representative MAG
533 genbank scaffolds are included as a supplementary file (Supplementary data 2). Site-directed
534 mutagenesis primers and synthesized candidate Form I' *rbcL* genes are included as a
535 supplementary file (Supplementary data 3). The structural coordinates of 2CABP-bound P. breve
536 Rubisco have been deposited in the PDB under the accession ID 6URA. The crystal structure of
537 Syn6301 Rubisco can be found on the PDB under the accession ID 1RBL. Publicly available
538 databases used in this study include: PDB (www.rcsb.org), pfam (www.pfam.xfam.org),
539 TIGRFams (www.tigrfams.jcvi.org), and KEGG database (www.genome.jp/kegg.html). Two
540 Chloroflexi genomes identified in this study are available at:

541

542 https://ggkbase.berkeley.edu/Chloroflexi_Rubisco_PatrickShih/organisms.

543

544 **Materials & correspondence.** Correspondence and material requests should be addressed to
545 P.M. Shih and J.F. Banfield

546

547 **Author information**

548

549 **Author contributions.** D.M.B, A.K.L., and P.M.S. designed experiments. D.M.B and A.K.L.
550 prepared all protein samples, performed all PAGE analyses, and protein thermal shift
551 experiments. M.H. performed all SEC-SAXS-MALS experiments and data analysis. J.H.P.
552 performed X-ray crystallography data acquisition, image processing, and structure determination.
553 D.M.B. performed all structural analyses. A.K.L. performed all site-directed mutagenesis
554 experiments. D.J.O. performed all Rubisco activity and kinetic measurements. C. H. and J.F.B
555 performed all metagenomic and phylogenetic analyses. All authors participated in writing and
556 manuscript preparation.

557

558 **Acknowledgments.** D.M.B, A.K.L., and P.M.S. acknowledge support from a Society in
559 Science–Branco Weiss fellowship from ETH Zurich. J.H.P., P.D.A., and P.M.S. acknowledge
560 support from the Joint BioEnergy Institute which is supported by the US Department of Energy,
561 Office of Science, Office of Biological and Environmental Research under Contract No. DE-
562 AC02-05CH11231 between Lawrence Berkeley National Laboratory and the US Department of
563 Energy. C.H. and J.F.B. thank Adi Lavy and Allison Sharrar for providing unpublished Rubisco
564 sequences, Jacob West-Roberts for assistance, the Rifle IFRC/SFA 2.0 Metagenomics and
565 Proteomics Data Analysis Project, the Allen Foundation, the Chan Zuckerberg Biohub, and the
566 Innovative Genomics Institute for support. C.H. acknowledges the Camille and Henry Dreyfus
567 Foundation for a postdoctoral fellowship, and the Joint Genome Institute CSP for sequencing.
568 M.H. acknowledge support from the Department of Energy BER Integrated Diffraction Analysis
569 Technologies (IDAT) program, NIGMS grant P30 GM124169-01, and ALS-ENABLE for SAXS

570 data collection at SIBYLS. D.J.O., M.A.J.P., and E.C.S. acknowledge support from the UK
571 Biotechnology and Biological Sciences Research Council (BBSRC; grant number
572 BB/I024488/1). We would like to thank Manajit Hayer-Hartl (Max Planck Institute of
573 Biochemistry, Martinsried, Germany) for the kind donation of the *Syn6301-rbcL*-pET11a,
574 *Syn6301-rbcLS*-pET11a, pG-*KJE8*, and pBAD33*ES/EL* plasmids used in this study.
575 Additionally, we would like to thank Noam Prywes for the kind donation of the pET28-His₁₄-
576 bdSUMO and pSF1389 plasmids. We would also like to thank Fei Guo and the UC Davis
577 BioEM core facility for EM images, and the laboratory of Justin Siegel (UC Davis Genome
578 Center) for use of their qPCR machine for protein thermal shift experiments. We are grateful to
579 Avi Flamholz for collecting publicly available Form I Rubisco kinetic data used in this study,
580 and to Alyssa Marinas and Rick Vermon Callado for assisting with enzyme purifications. We
581 thank Kasey Markel for his edits and suggestions on the manuscript.

582

583 **Declarations of interest**

584

585 The authors declare no competing interests.

586

587 **References**

- 588 1. Nisbet, E. G. *et al.* The age of Rubisco: the evolution of oxygenic photosynthesis.
589 *Geobiology* **5**, 311–335 (2007).
- 590 2. Tabita, F. R. *et al.* Function, Structure, and Evolution of the RubisCO-Like Proteins and
591 Their RubisCO Homologs. *Microbiol. Mol. Biol. Rev.* **71**, 576–599 (2007).
- 592 3. Tabita, F. R., Satagopan, S., Hanson, T. E., Kreel, N. E. & Scott, S. S. Distinct form I, II,
593 III, and IV Rubisco proteins from the three kingdoms of life provide clues about Rubisco
594 evolution and structure/function relationships. *J. Exp. Bot.* **59**, 1515–1524 (2007).
- 595 4. Andrews, T. J. Catalysis by cyanobacterial ribulose-bisphosphate carboxylase large subunits
596 in the complete absence of small subunits. *J. Biol. Chem.* **263**, 12213–12219 (1988).
- 597 5. Morell, M. K., Wilkin, J. M., Kane, H. J. & Andrews, T. J. Side reactions catalyzed by
598 ribulose-bisphosphate carboxylase in the presence and absence of small subunits. *J. Biol.*
599 *Chem.* **272**, 5445–5451 (1997).
- 600 6. Spreitzer, R. J. Role of the small subunit in ribulose-1,5-bisphosphate

- 601 carboxylase/oxygenase. *Archives of Biochemistry and Biophysics* vol. 414 141–149 (2003).
- 602 7. Joshi, J., Mueller-Cajar, O., Tsai, Y.-C. C., Hartl, F. U. & Hayer-Hartl, M. Role of Small
603 Subunit in Mediating Assembly of Red-type Form I Rubisco. *J. Biol. Chem.* **290**, 1066–
604 1074 (2015).
- 605 8. Liu, C. *et al.* Coupled chaperone action in folding and assembly of hexadecameric Rubisco.
606 *Nature* **463**, 197–202 (2010).
- 607 9. Grabsztunowicz, M., Górski, Z., Luciński, R. & Jackowski, G. A reversible decrease in
608 ribulose 1,5-bisphosphate carboxylase/oxygenase carboxylation activity caused by the
609 aggregation of the enzyme's large subunit is triggered in response to the exposure of
610 moderate irradiance-grown plants to low irradiance. *Physiol. Plant.* **154**, 591–608 (2015).
- 611 10. Kusian, B. & Bowien, B. Organization and regulation of cbb CO₂ assimilation genes in
612 autotrophic bacteria. *FEMS Microbiol. Rev.* **21**, 135–155 (1997).
- 613 11. Tabita, F. R. Microbial ribulose 1,5-bisphosphate carboxylase/oxygenase: A different
614 perspective. *Photosynth. Res.* **60**, 1–28 (1999).
- 615 12. Whitney, S. M. & Andrews, T. J. The gene for the ribulose-1,5-bisphosphate
616 carboxylase/oxygenase (Rubisco) small subunit relocated to the plastid genome of tobacco
617 directs the synthesis of small subunits that assemble into Rubisco. *Plant Cell* **13**, 193–205
618 (2001).
- 619 13. Bryant, D. A. & Liu, Z. Chapter Four - Green Bacteria: Insights into Green Bacterial
620 Evolution through Genomic Analyses. in *Advances in Botanical Research* (ed. Beatty, J. T.)
621 vol. 66 99–150 (Academic Press, 2013).
- 622 14. Shih, P. M., Ward, L. M. & Fischer, W. W. Evolution of the 3-hydroxypropionate bicycle
623 and recent transfer of anoxygenic photosynthesis into the Chloroflexi. *Proc. Natl. Acad. Sci.*

- 624 *U. S. A.* **114**, 10749–10754 (2017).
- 625 15. Ward, L. M., Hemp, J., Shih, P. M., McGlynn, S. E. & Fischer, W. W. Evolution of
626 Phototrophy in the Chloroflexi Phylum Driven by Horizontal Gene Transfer. (2018)
627 doi:10.3389/fmicb.2018.00260.
- 628 16. Fischer, W. W., Hemp, J. & Johnson, J. E. Evolution of Oxygenic Photosynthesis. *Annu.*
629 *Rev. Earth Planet. Sci.* **44**, 647–683 (2016).
- 630 17. Roy, H. Rubisco assembly: a model system for studying the mechanism of chaperonin
631 action. *Plant Cell* **1**, 1035–1042 (1989).
- 632 18. Hayer-Hartl, M. From chaperonins to Rubisco assembly and metabolic repair. *Protein Sci.*
633 **26**, 2324–2333 (2017).
- 634 19. Aigner, H. *et al.* Plant RuBisCo assembly in *E. coli* with five chloroplast chaperones
635 including BSD2. *Science* **358**, 1272–1278 (2017).
- 636 20. Wilson, R. H. & Hayer-Hartl, M. Complex Chaperone Dependence of Rubisco Biogenesis.
637 *Biochemistry* **57**, 3210–3216 (2018).
- 638 21. Saschenbrecker, S. *et al.* Structure and Function of RbcX, an Assembly Chaperone for
639 Hexadecameric Rubisco. *Cell* **129**, 1189–1200 (06/2007).
- 640 22. Gunn, L. H., Valegård, K. & Andersson, I. A unique structural domain in
641 *Methanococcoides burtonii* ribulose-1,5-bisphosphate carboxylase/oxygenase (Rubisco)
642 acts as a small subunit mimic. *J. Biol. Chem.* **292**, 6838–6850 (2017).
- 643 23. Goloubinoff, P., Christeller, J. T., Gatenby, A. A. & Lorimer, G. H. Reconstitution of active
644 dimeric ribulose bisphosphate carboxylase from an unfolded state depends on two
645 chaperonin proteins and Mg-ATP. *Nature* **342**, 884–889 (1989).
- 646 24. Parry, M. A. J., Keys, A. J. & Gutteridge, S. Variation in the Specificity Factor of C3

- 647 Higher Plant Rubiscos Determined by the Total Consumption of Ribulose-P2. *J. Exp. Bot.*
648 **40**, 317–320 (1989).
- 649 25. Tcherkez, G. G. B., Farquhar, G. D. & Andrews, T. J. Despite slow catalysis and confused
650 substrate specificity, all ribulose biphosphate carboxylases may be nearly perfectly
651 optimized. *Proceedings of the National Academy of Sciences* **103**, 7246–7251 (2006).
- 652 26. Flamholz, A. I. *et al.* Revisiting Trade-offs between Rubisco Kinetic Parameters.
653 *Biochemistry* **58**, 3365–3376 (2019).
- 654 27. Yamada, T. & Sekiguchi, Y. Cultivation of uncultured chloroflexi subphyla: significance
655 and ecophysiology of formerly uncultured chloroflexi ‘subphylum i’ with natural and
656 biotechnological relevance. *Microbes Environ.* **24**, 205–216 (2009).
- 657 28. Hemp, J., Ward, L. M., Pace, L. A. & Fischer, W. W. Draft Genome Sequence of
658 *Ornatilinea apprima* P3M-1, an Anaerobic Member of the Chloroflexi Class Anaerolineae.
659 *Genome Announc.* **3**, (2015).
- 660 29. Ward, L. M., Hemp, J., Pace, L. A. & Fischer, W. W. Draft Genome Sequence of
661 *Leptolinea tardivitalis* YMTK-2, a Mesophilic Anaerobe from the Chloroflexi Class
662 Anaerolineae. *Genome Announc.* **3**, (2015).
- 663 30. Alonso, H., Blayney, M. J., Beck, J. L. & Whitney, S. M. Substrate-induced assembly of
664 *Methanococcoides burtonii* D-ribulose-1,5-bisphosphate carboxylase/oxygenase dimers into
665 decamers. *J. Biol. Chem.* **284**, 33876–33882 (2009).
- 666 31. Knott, G. J. *et al.* Structural basis for AcrVA4 inhibition of specific CRISPR-Cas12a. *Elife*
667 **8**, (2019).
- 668 32. Duff, A. P., Andrews, T. J. & Curmi, P. M. The transition between the open and closed
669 states of rubisco is triggered by the inter-phosphate distance of the bound bisphosphate. *J.*

- 670 *Mol. Biol.* **298**, 903–916 (2000).
- 671 33. Newman, J., Branden, C. I. & Jones, T. A. Structure determination and refinement of
672 ribulose 1,5-bisphosphate carboxylase/oxygenase from *Synechococcus* PCC6301. *Acta*
673 *Crystallogr. D Biol. Crystallogr.* **49**, 548–560 (1993).
- 674 34. Lu, Z., Zhao, Z. & Fu, B. Efficient protein alignment algorithm for protein search. *BMC*
675 *Bioinformatics* **11 Suppl 1**, S34 (2010).
- 676 35. Cleland, W. W., Andrews, T. J., Gutteridge, S., Hartman, F. C. & Lorimer, G. H.
677 Mechanism of Rubisco: The Carbamate as General Base ^x. *Chem. Rev.* **98**, 549–562 (1998).
- 678 36. Andersson, I. & Backlund, A. Structure and function of Rubisco. *Plant Physiol. Biochem.*
679 **46**, 275–291 (2008).
- 680 37. van Lun, M., van der Spoel, D. & Andersson, I. Subunit interface dynamics in
681 hexadecameric rubisco. *J. Mol. Biol.* **411**, 1083–1098 (2011).
- 682 38. Schneider, G. *et al.* Comparison of the crystal structures of L2 and L8S8 Rubisco suggests a
683 functional role for the small subunit. *EMBO J.* **9**, 2045–2050 (1990).
- 684 39. Huynh, K. & Partch, C. L. Analysis of protein stability and ligand interactions by thermal
685 shift assay. *Curr. Protoc. Protein Sci.* **79**, 28.9.1–14 (2015).
- 686 40. Greene, D. N., Whitney, S. M. & Matsumura, I. Artificially evolved *Synechococcus*
687 PCC6301 Rubisco variants exhibit improvements in folding and catalytic efficiency.
688 *Biochem. J* **404**, 517–524 (2007).
- 689 41. DePristo, M. A., Weinreich, D. M. & Hartl, D. L. Missense meanderings in sequence space:
690 a biophysical view of protein evolution. *Nat. Rev. Genet.* **6**, 678–687 (2005).
- 691 42. Tokuriki, N., Stricher, F., Serrano, L. & Tawfik, D. S. How protein stability and new
692 functions trade off. *PLoS Comput. Biol.* **4**, e1000002 (2008).

- 693 43. Tokuriki, N. & Tawfik, D. S. Protein dynamism and evolvability. *Science* **324**, 203–207
694 (2009).
- 695 44. Erb, T. J. & Zarzycki, J. A short history of RubisCO: the rise and fall (?) of Nature’s
696 predominant CO₂ fixing enzyme. *Curr. Opin. Biotechnol.* **49**, 100–107 (02/2018).
- 697 45. Badger, M. R., Hanson, D. & Dean Price, G. Evolution and diversity of CO₂ concentrating
698 mechanisms in cyanobacteria. *Funct. Plant Biol.* **29**, 161–173 (2002).
- 699 46. Studer, R. A., Christin, P.-A., Williams, M. A. & Orengo, C. A. Stability-activity tradeoffs
700 constrain the adaptive evolution of RubisCO. *Proceedings of the National Academy of*
701 *Sciences* **111**, 2223–2228 (2014).
- 702 47. Zhou, Y. & Whitney, S. Directed Evolution of an Improved Rubisco; In Vitro Analyses to
703 Decipher Fact from Fiction. *Int. J. Mol. Sci.* **20**, (2019).
- 704 48. Wilson, R. H., Alonso, H. & Whitney, S. M. Evolving *Methanococoides burtonii* archaeal
705 Rubisco for improved photosynthesis and plant growth. *Sci. Rep.* **6**, 22284 (2016).
- 706 49. Rambo, R. P. & Tainer, J. A. Accurate assessment of mass, models and resolution by small-
707 angle scattering. *Nature* **496**, 477–481 (2013).
- 708 50. Diamond, S. *et al.* Mediterranean grassland soil C-N compound turnover is dependent on
709 rainfall and depth, and is mediated by genomically divergent microorganisms. *Nat*
710 *Microbiol* **4**, 1356–1367 (2019).
- 711 51. Lavy, A. *et al.* Microbial communities across a hillslope-riparian transect shaped by
712 proximity to the stream, groundwater table, and weathered bedrock. *Ecol. Evol.* **9**, 6869–
713 6900 (2019).
- 714 52. Knight, S., Andersson, I. & Brändén, C. I. Crystallographic analysis of ribulose 1,5-
715 bisphosphate carboxylase from spinach at 2.4 Å resolution. Subunit interactions and active

- 716 site. *J. Mol. Biol.* **215**, 113–160 (1990).
- 717 53. Pei, J., Kim, B.-H. & Grishin, N. V. PROMALS3D: a tool for multiple protein sequence
718 and structure alignments. *Nucleic Acids Res.* **36**, 2295–2300 (2008).
- 719 54. Katoh, K., Rozewicki, J. & Yamada, K. D. MAFFT online service: multiple sequence
720 alignment, interactive sequence choice and visualization. *Brief. Bioinform.* (2017)
721 doi:10.1093/bib/bbx108.
- 722 55. Potterton, E., Briggs, P., Turkenburg, M. & Dodson, E. A graphical user interface to the
723 CCP4 program suite. *Acta Crystallogr. D Biol. Crystallogr.* **59**, 1131–1137 (2003).
- 724 56. Mueller-Cajar, O., Morell, M. & Whitney, S. M. Directed Evolution of Rubisco in
725 *Escherichia coli* Reveals a Specificity-Determining Hydrogen Bond in the Form II Enzyme.
726 *Biochemistry* **46**, 14067–14074 (2007).
- 727 57. Schneidman-Duhovny, D., Hammel, M. & Sali, A. FoXS: a web server for rapid
728 computation and fitting of SAXS profiles. *Nucleic Acids Res.* **38**, W540–4 (2010).
- 729 58. Schneidman-Duhovny, D., Hammel, M., Tainer, J. A. & Sali, A. Accurate SAXS profile
730 computation and its assessment by contrast variation experiments. *Biophys. J.* **105**, 962–974
731 (2013).
- 732 59. Frey, S. & Görlich, D. A new set of highly efficient, tag-cleaving proteases for purifying
733 recombinant proteins. *J. Chromatogr. A* **1337**, 95–105 (2014).
- 734 60. Kane, H. J., Wilkin, J. M., Portis, A. R. & John Andrews T. Potent inhibition of ribulose-
735 biphosphate carboxylase by an oxidized impurity in ribulose-1,5-biphosphate. *Plant*
736 *Physiol.* **117**, 1059–1069 (1998).
- 737 61. Pierce, J., Tolbert, N. E. & Barker, R. Interaction of ribulosebiphosphate
738 carboxylase/oxygenase with transition-state analogues. *Biochemistry* **19**, 934–942 (1980).

- 739 62. Pereira, J. H., McAndrew, R. P., Tomaleri, G. P. & Adams, P. D. Berkeley Screen: a set of
740 96 solutions for general macromolecular crystallization. *J. Appl. Crystallogr.* **50**, 1352–
741 1358 (2017).
- 742 63. Winter, G., Lobley, C. M. C. & Prince, S. M. Decision making in xia2. *Acta Crystallogr. D*
743 *Biol. Crystallogr.* **69**, 1260–1273 (2013).
- 744 64. McCoy, A. J. *et al.* Phaser crystallographic software. *J. Appl. Crystallogr.* **40**, 658–674
745 (2007).
- 746 65. Adams, P. D. *et al.* PHENIX: a comprehensive Python-based system for macromolecular
747 structure solution. *Acta Crystallogr. D Biol. Crystallogr.* **66**, 213–221 (2010).
- 748 66. Afonine, P. V. *et al.* Towards automated crystallographic structure refinement with
749 phenix.refine. *Acta Crystallogr. D Biol. Crystallogr.* **68**, 352–367 (2012).
- 750 67. Emsley, P. & Cowtan, K. Coot: model-building tools for molecular graphics. *Acta*
751 *Crystallogr. D Biol. Crystallogr.* **60**, 2126–2132 (2004).
- 752 68. Davis, I. W. *et al.* MolProbity: all-atom contacts and structure validation for proteins and
753 nucleic acids. *Nucleic Acids Res.* **35**, W375–83 (2007).
- 754 69. Dyer, K. N. *et al.* High-throughput SAXS for the characterization of biomolecules in
755 solution: a practical approach. *Methods Mol. Biol.* **1091**, 245–258 (2014).
- 756 70. Hura, G. L. *et al.* Robust, high-throughput solution structural analyses by small angle X-ray
757 scattering (SAXS). *Nat. Methods* **6**, 606–612 (2009).
- 758 71. Sali, A. & Blundell, T. L. Comparative protein modelling by satisfaction of spatial
759 restraints. *J. Mol. Biol.* **234**, 779–815 (1993).
- 760 72. Prins, A. *et al.* Rubisco catalytic properties of wild and domesticated relatives provide scope
761 for improving wheat photosynthesis. *J. Exp. Bot.* **67**, 1827–1838 (03/2016).

- 762 73. Sharwood, R. E., Ghannoum, O. & Whitney, S. M. Prospects for improving CO₂ fixation in
763 C₃-crops through understanding C₄-Rubisco biogenesis and catalytic diversity. *Curr. Opin.*
764 *Plant Biol.* **31**, 135–142 (2016).
- 765 74. Krissinel, E. & Henrick, K. Inference of macromolecular assemblies from crystalline state.
766 *J. Mol. Biol.* **372**, 774–797 (2007).
- 767 75. Krissinel, E. Crystal contacts as nature’s docking solutions. *J. Comput. Chem.* **31**, 133–143
768 (2010).
- 769 76. Pettersen, E. F. *et al.* UCSF Chimera--a visualization system for exploratory research and
770 analysis. *J. Comput. Chem.* **25**, 1605–1612 (2004).

771

772

773

774

775

776

777

778 **Figure Legends**

779

780 **Fig. 1. Metagenomics-enabled identification of a novel clade of Form I Rubisco that lack**
781 **small subunits.** **a**, Maximum likelihood phylogeny of Rubisco RbcL. By including recently
782 discovered metagenome-assembled genomes (MAGs) from Chloroflexi, the emergence of a
783 *bona fide*, well-supported clade of Rubisco was identified (Form I’). Black circles indicate
784 bootstrap values of 100 and white circles indicate bootstrap values >90. **b**, Example Chloroflexi
785 operons with Form I’ Rubisco (dark blue) reveal no presence of a *rbcS*, a defining feature of
786 Form I Rubisco, which are almost always found immediately neighboring *rbcL* in bacteria;
787 however, other CBB cycle-related genes are found in the operon (light blue). White, other
788 enzymes; gray, hypothetical protein. Annotated loci (i-v) represent Scaffolds 211530, 92,
789 509483, 467972, and 172446, respectively. For the full annotation information see
790 Supplementary Data 2. GAPDH, glyceraldehyde-3-phosphate dehydrogenase; *cbbT*,

791 transketolase; PRK, phosphoribulokinase; FBP, fructose bisphosphate; TBP, tagatose
792 bisphosphate; cbbF, fructose 1,6-bisphosphatase.
793

794 **Fig. 2. Comparison of *P. breve* kinetic data to reported values of Form I Rubisco.** Scatter
795 plots of reported Form I Rubisco kinetic data (black circles) collected at 25 °C²⁶ against *P. breve*
796 Form I' Rubisco (green dots), including maximum rates of carboxylation and oxygenation of
797 RuBP (V_C and V_O , respectively), the catalytic efficiency of carboxylation over oxygenation
798 ($S_{C/O}$), and Michaelis constants for carboxylation and oxygenation of RuBP (K_C and K_O ,
799 respectively). Gray dotted lines represent the median for collected Form I Rubisco kinetic data.
800

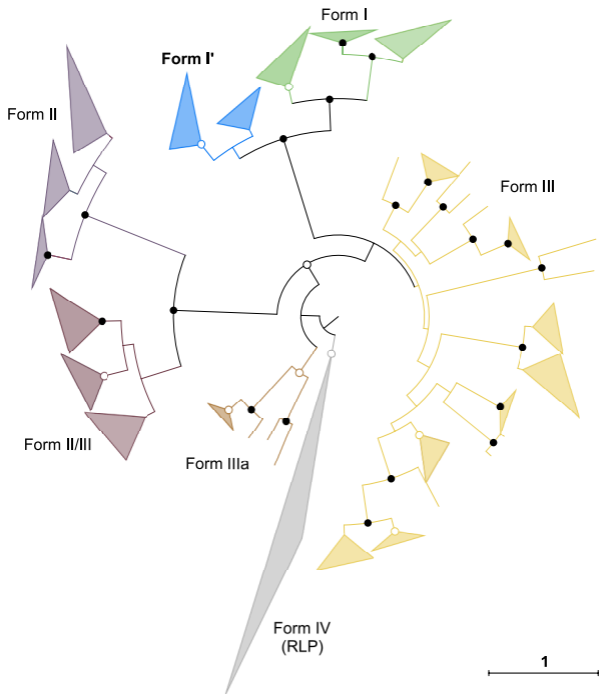
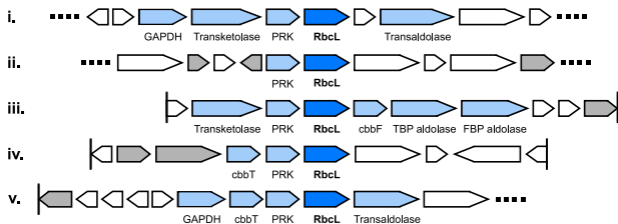
801 **Fig. 3. Solution-state characterization of Form I' oligomerization reveals an octameric**
802 **holoenzyme reminiscent of canonical Form I Rubisco.** **a**, SEC-SAXS-MALS chromatograms
803 of the separation of activated *P. breve* Rubisco in the absence (top) or presence (bottom) of
804 bound 2CABP. Solid gray lines represent the UV absorbance reading at 280 nm, dashed black
805 lines represent the integrated SAXS signal, while circles represent molecular mass (light blue)
806 data collected from MALS, and R_g values for each SAXS frame (dark blue) versus elution time.
807 **b**, Experimental $P(r)$ functions determined from SAXS profiles (black dashes) of *P. breve*
808 Rubisco in the open conformation (light blue) or bound to 2CABP (dark blue). The area under
809 the $P(r)$ function is normalized relative to the molecular weight estimated by SAXS⁴⁹ and is
810 listed in Supplementary Table 2. Theoretical $P(r)$ functions are calculated from the theoretical
811 SAXS curves of the corresponding models shown in panel C. The radius where $P(r)$ approaches
812 zero intensity identifies the maximal dimension of the macromolecule (dashed arrows). **c**,
813 Surface representation models of *P. breve* Rubisco with extended (open conformation) or
814 compact (closed conformation) C-terminal regions. **d**, A representative non-denaturing PAGE
815 gel demonstrating the migration of *P. breve* Rubisco in the absence (-) or presence (+) of
816 2CABP. M = molecular weight marker. Native gel electrophoresis experiment was performed at
817 $n > 10$.
818

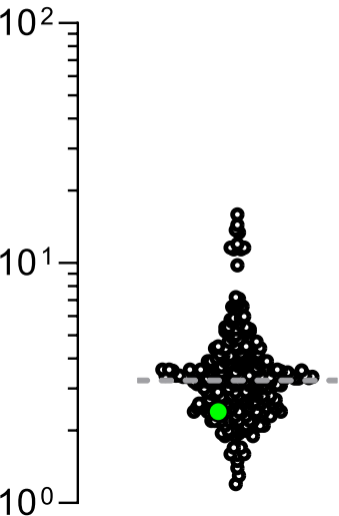
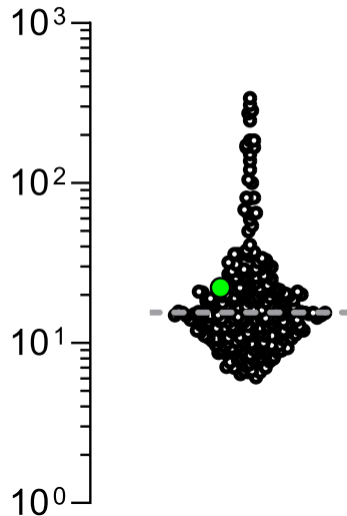
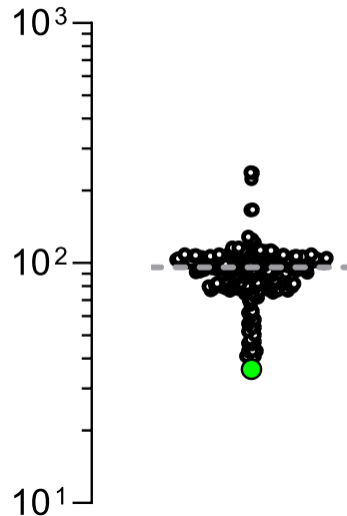
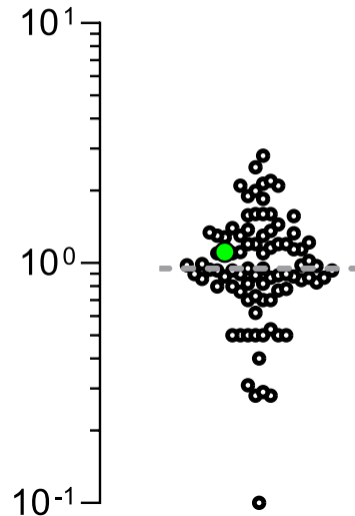
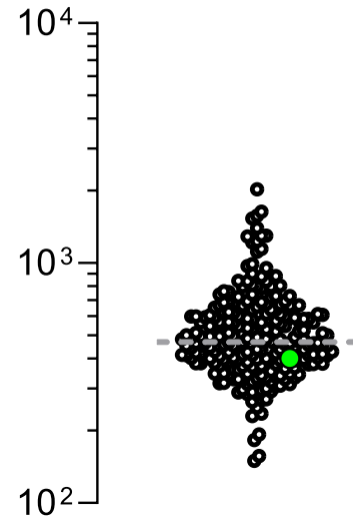
819 **Fig. 4. Crystal structure of Form I' Rubisco compared to cyanobacterial Form I Rubisco.**
820 Comparison of the structural models of **a**, Form I Rubisco from *Synechococcus sp.* strain PCC
821 6301 (PDB ID: 1RBL) RbcL (green) with RbcS (tan), and **b**, Form I' Rubisco from *P. breve*
822 (PDB ID: 6URA, blue) which lacks RbcS. Coulombic electrostatic potential maps of 1RBL
823 (RbcS removed) and *P. breve* Rubisco are illustrated by the charge distributions (negative, red;
824 neutral, white; positive, blue) of the surface residues of either structure.
825
826

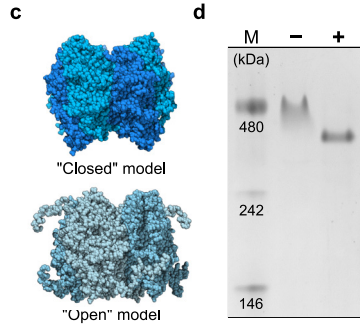
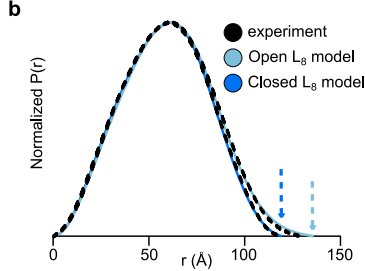
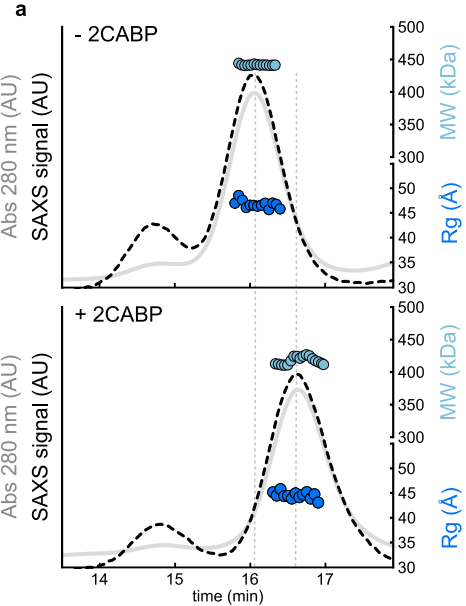
827 **Fig. 5. Structure of Form I' Rubisco suggests that a tradeoff between stability and catalytic**
828 **activity spurred the evolution of the small subunit.** **a**, Salt bridge and hydrogen bond networks
829 present at the dimer-dimer interface of *P. breve* Rubisco mediate holoenzyme stability in the
830 absence of small subunits. Separate RbcL dimers at the dimer-dimer pair are distinguished by

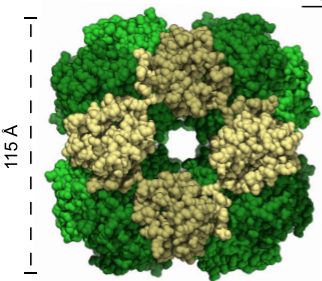
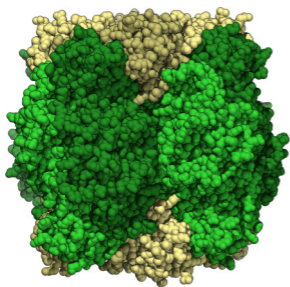
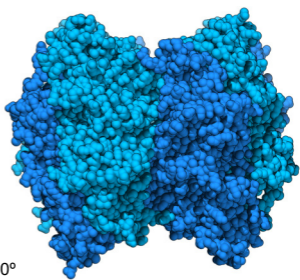
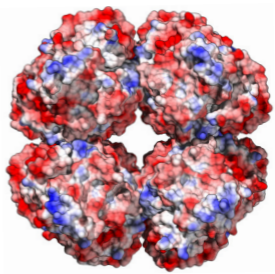
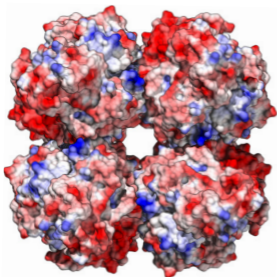
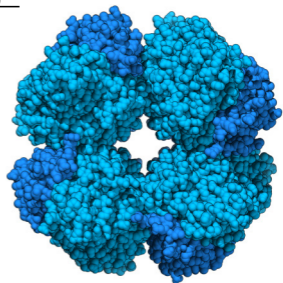
831 two separate shades of blue. **b**, Protein thermal shift assay data with annotated melting
832 temperatures (T_m) for the disassembly of RbcL dimer quaternary structure from wild-type
833 *Syn6301* RbcL, *Syn6301* RbcLS, and *P. breve* RbcL. Reported T_m values represent the average
834 measured from a total of four experiments.
835

Table 1 Kinetic characterization of Form I' Rubisco at 25 °C.					
Rubisco	V_C (s ⁻¹)	K_C (μM)	$S_{C/O}$	V_O (s ⁻¹)	K_O (μM)
Form I' <i>P. breve</i>	2.23 ± 0.04 (5)	22.2 ± 9.7 (5)	36.1 ± 0.9 (10)	1.11(5)	401 ± 115 (5)
Form I <i>Synechococcus sp.</i> strain PCC 6301	14.3 ± 0.71 (4)	235 ± 20.0 (4)	56.1 ± 1.3 (4)	1.10 (4)	983 ± 81 (4)
<p>V_C and V_O correspond to the maximal rates of the carboxylation and oxygenation reactions, respectively, under saturating substrate concentrations. K_C and K_O are the Michaelis constants (K_M) for the carboxylation and oxygenation reactions, respectively. $S_{C/O} = (V_C/K_C)/(V_O/K_O)$. Values represent the mean ± S.E. with n indicated in parentheses, where n reflects the number of experiments conducted with the same protein sample.</p>					

a**b**

V_C (s^{-1}) K_C (μM) $S_{C/O}$  V_O (s^{-1}) K_O (μM)



a*Syn6301* RbcLS (L8S8)**b***P. breve* RbcL (L8) 30°
↑
↓

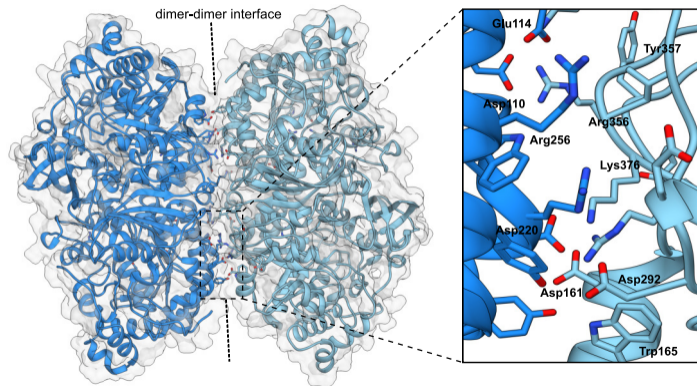
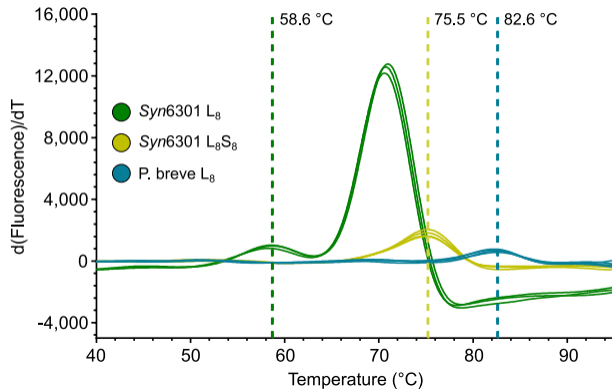
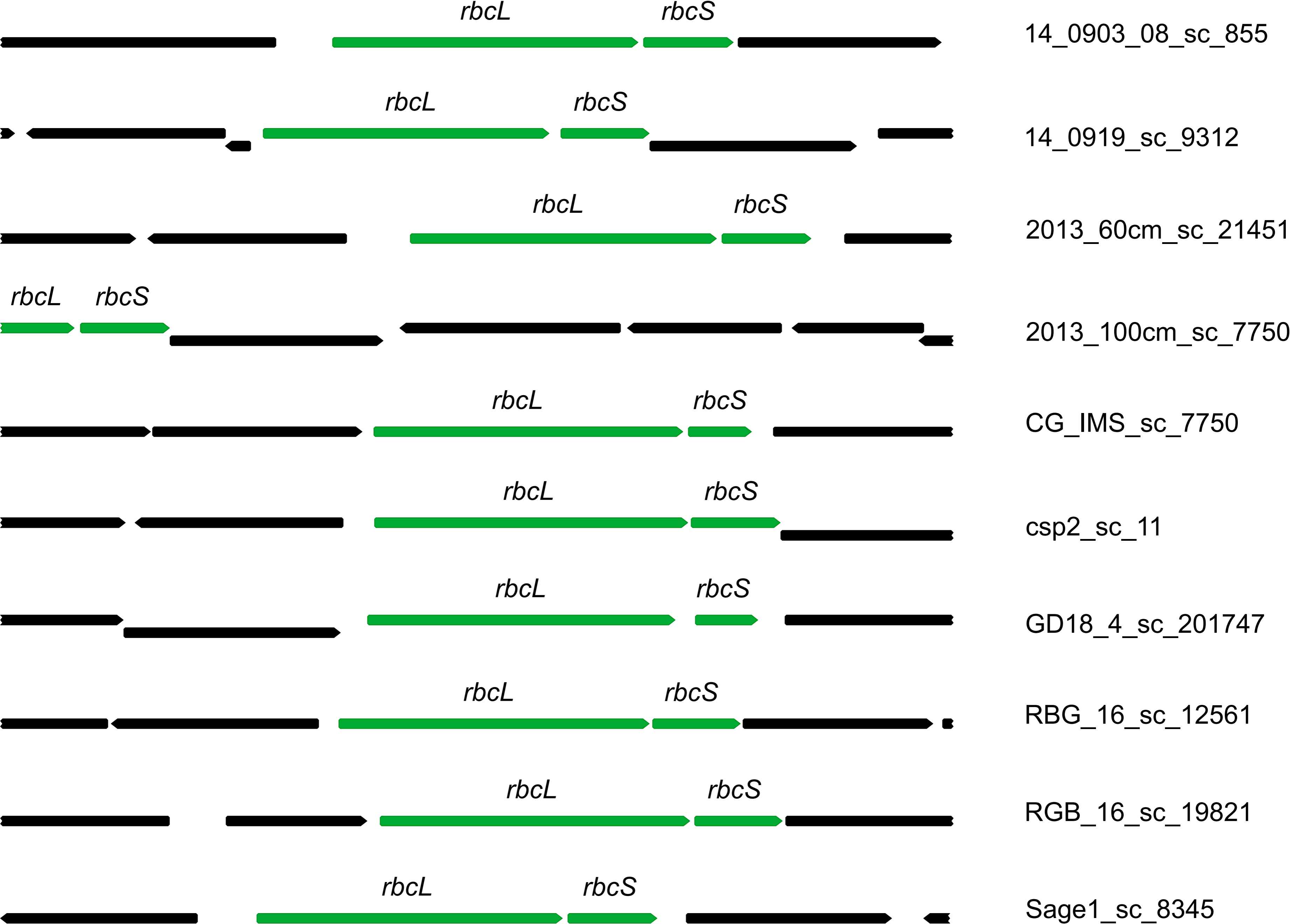
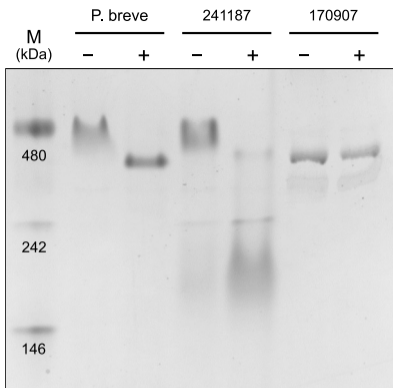
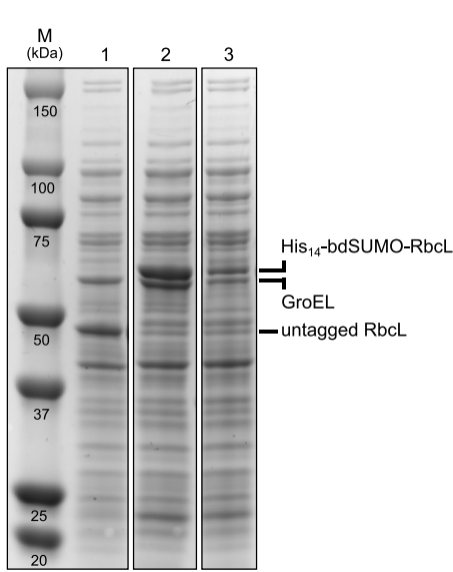
a**b**

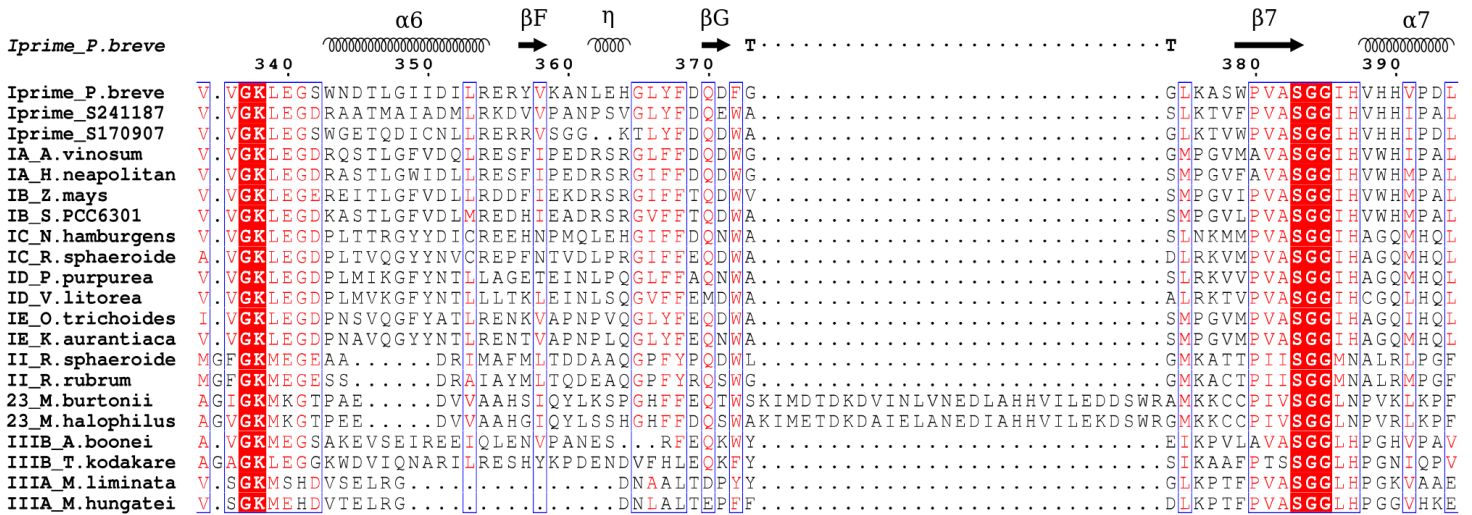
Table 1 | Kinetic characterization of Form I' Rubisco at 25 °C.

Rubisco	V_C (s⁻¹)	K_C (μM)	$S_{C/O}$	V_O (s⁻¹)	K_O (μM)
Form I' P. breve	2.23 ± 0.04 (5)	22.2 ± 9.7 (5)	36.1 ± 0.9 (10)	1.11(5)	401 ± 115 (5)
Form I <i>Synechococcus</i> sp. strain PCC 6301	14.3 ± 0.71 (4)	235 ± 20.0 (4)	56.1 ± 1.3 (4)	1.10 (4)	983 ± 81 (4)

V_C and V_O correspond to the maximal rates of the carboxylation and oxygenation reactions, respectively, under saturating substrate concentrations. K_C and K_O are the Michaelis constants (K_M) for the carboxylation and oxygenation reactions, respectively. $S_{C/O} = (V_C/K_C)/(V_O/K_O)$. Values represent the mean ± S.E. with n indicated in parentheses, where n reflects the number of experiments conducted with the same protein sample.

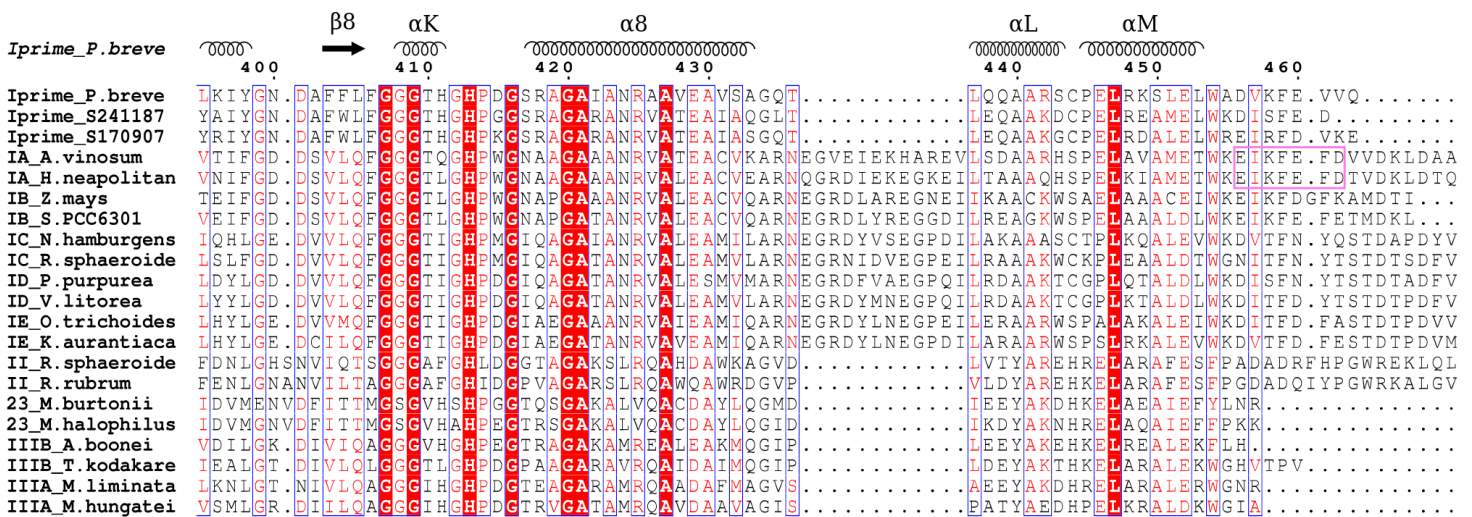


a**b**

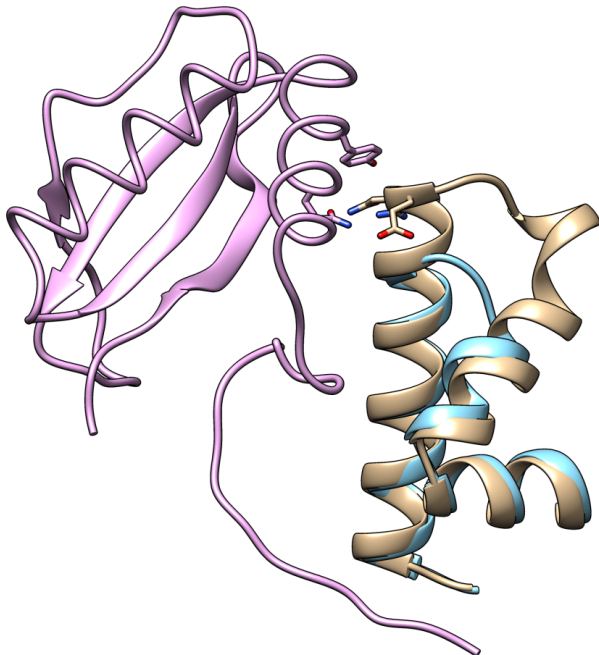
a

loop 6

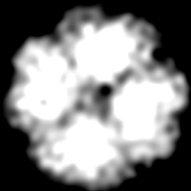
Rubisco assembly domain (Form II/III)



C-terminal ext. (Form I)

b

100 Å



1



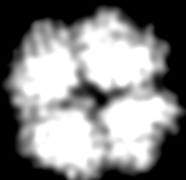
2



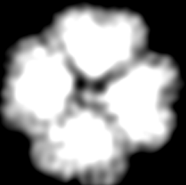
3



4



5



6

Intensity

-2CABP

+2CABP

- experiment
- Open L_8 model
- Closed L_8 model

$\ln(I(q))$

0.000

$q^2 (\text{\AA}^{-2})$

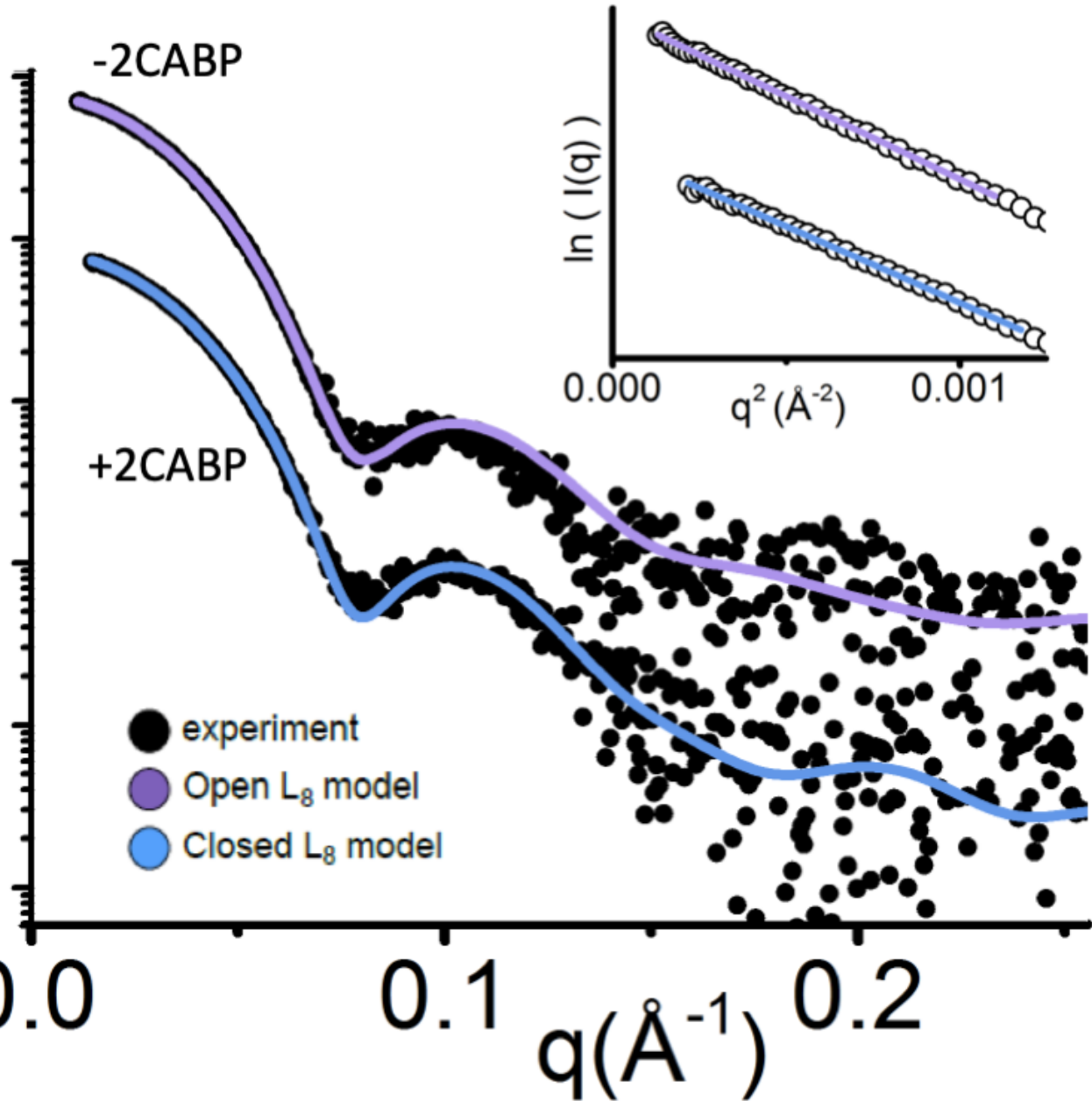
0.001

0.0

0.1

0.2

$q(\text{\AA}^{-1})$



a

Form I <i>Syn6301</i>	1	MP-----KTQSAAGYKAGVKDYKLTYYTPDYTPKDTDLLAAAFREFSPOPGV
Form I' <i>P. breve</i>	1	MAIHNPLAGPKTVKARPTAELSDAYKAGVRAYAVDYVVPDYIPQDTDLLCAFRIQPR-GV
Form I <i>Syn6301</i>	46	PAD [*] EAGAAIAAESSTGTWTTVWTDLLTDM [*] DRYK [*] GKCYHIEPVQGEENSYFAEIAYP [*] LDL [*] LF
Form I' <i>P. breve</i>	60	DMI [*] EAAA [*] AVAAESSTGTWTEVWSN [*] QLTDID [*] FYKAKVYAITG-----DIAIYIAYPLDL [*] LF
Form I <i>Syn6301</i>	106	E [*] EGSV [*] TNILT [*] SIVGNVFGFKAIRSLRLEDI [*] RFP [*] VALVKT [*] FQ [*] GPPH [*] GI [*] Q [*] VER [*] DL [*] LNKY [*] GR [*] P
Form I' <i>P. breve</i>	113	E [*] ENS [*] VVNIMSSIVGNVFGFKAVGALRLED [*] MRIP [*] LALVKT [*] TF [*] GPRV [*] GI [*] YDER [*] VW [*] SNK [*] WDR [*] P
Form I <i>Syn6301</i>	166	MLGCTIKPKLGL [*] SAK [*] NYGRAVYECL [*] RGG [*] LDF [*] TKD [*] DEN [*] INS [*] Q [*] PF [*] OR [*] WR [*] DR [*] FL [*] EVAD [*] AI [*] H [*] KS
Form I' <i>P. breve</i>	173	LIGGT [*] VKPKLGLSP [*] KAY [*] STIIYECL [*] SGGL [*] DT [*] SKD [*] DEN [*] MNS [*] Q [*] PF [*] SR [*] WR [*] DR [*] FM [*] Y [*] AQ [*] EAV [*] DR [*] A
Form I <i>Syn6301</i>	226	QAETGEIKGHY [*] LNVTAPT [*] CEEM [*] K [*] RAE [*] FA [*] K [*] ELG [*] MP [*] I [*] IMH [*] D [*] FLT [*] AG [*] FT [*] ANT [*] TL [*] AK [*] WC [*] RD [*] NG
Form I' <i>P. breve</i>	233	AAETNEFKGHWHNVTAG [*] STEES [*] LR [*] RL [*] E [*] Y [*] AVEL [*] GS [*] RM [*] VM [*] FDEL [*] TAG [*] FA [*] ASAD [*] IF [*] KR [*] AG [*] EL [*] D
Form I <i>Syn6301</i>	286	VLLH [*] I [*] H [*] RAM [*] H [*] AVID [*] RQ [*] R [*] NH [*] GI [*] H [*] ER [*] VL [*] AK [*] CL [*] RL [*] LS [*] GG [*] DHL [*] HS [*] GT [*] VV [*] GK [*] LEG [*] DK [*] AST [*] L [*] GF [*] V [*] DL
Form I' <i>P. breve</i>	293	MIV [*] HC [*] H [*] RAM [*] H [*] AV [*] TR [*] Q [*] AN [*] H [*] GI [*] AM [*] RV [*] AK [*] WL [*] RL [*] T [*] GG [*] DHL [*] HT [*] GT [*] VV [*] GK [*] LEG [*] SW [*] ND [*] TL [*] GI [*] IDI
Form I <i>Syn6301</i>	346	MRED [*] HI [*] EADR [*] SR [*] GV [*] FE [*] TQ [*] DW [*] AS [*] MP [*] GV [*] LP [*] VAS [*] GGI [*] H [*] V [*] WH [*] MP [*] AL [*] VE [*] IE [*] FG [*] DD [*] SV [*] LQ [*] FG [*] GG [*] T [*] LG
Form I' <i>P. breve</i>	353	LR [*] E [*] RY [*] VKAN [*] LEH [*] G [*] LY [*] FD [*] Q [*] DFG [*] CL [*] KAS [*] WP [*] VAS [*] GGI [*] H [*] V [*] H [*] VP [*] DL [*] LKI [*] YG [*] ND [*] AFF [*] LF [*] GG [*] G [*] TH [*] G
Form I <i>Syn6301</i>	406	HPWGNAPGATANRVALEACVQARNEGRDLYREGGDI [*] LR [*] EAC [*] K [*] W [*] SP [*] E [*] LAAALDLWKEIK [*] FE
Form I' <i>P. breve</i>	413	HPD [*] GS [*] RAGAIANRAAVEAVSAG-----QT [*] LQQAARS [*] CP [*] EL [*] RKS [*] LE [*] LW [*] ADV [*] K [*] FE
Form I <i>Syn6301</i>	466	FETMDKL
Form I' <i>P. breve</i>	461	VVQ----

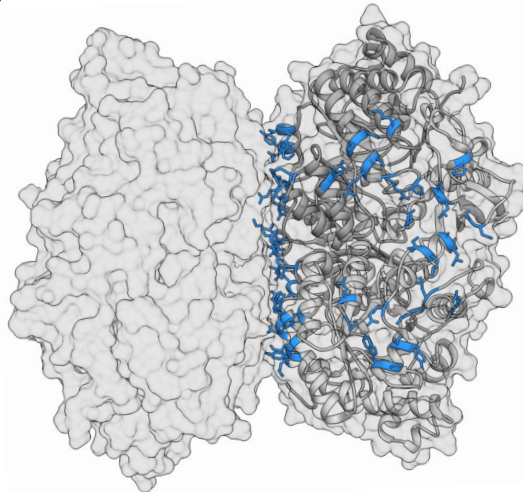
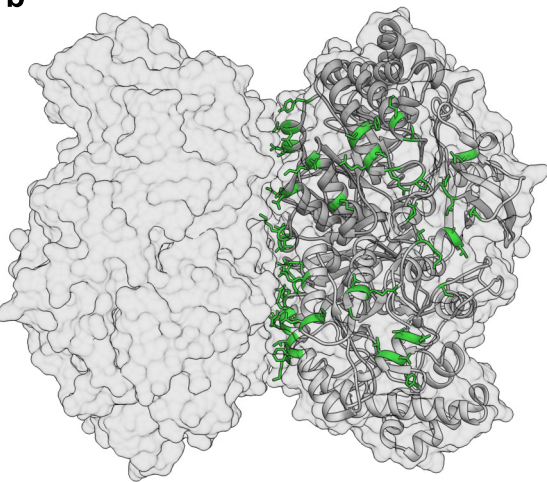
● *Syn6301* dimer-dimer interactions

● *P. breve* dimer-dimer interactions

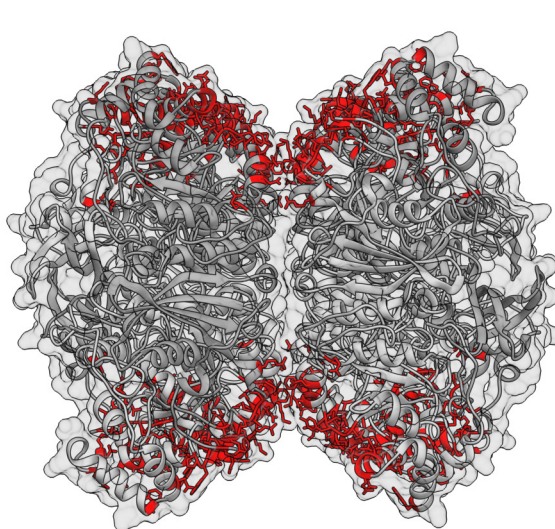
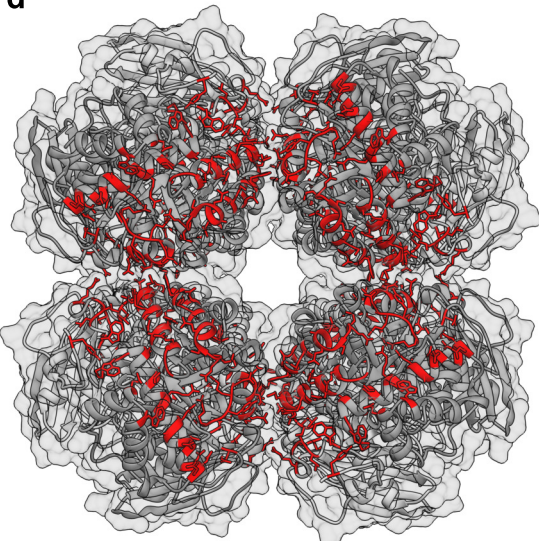
* *Syn6301* RbcL-RbcS interactions

b

c

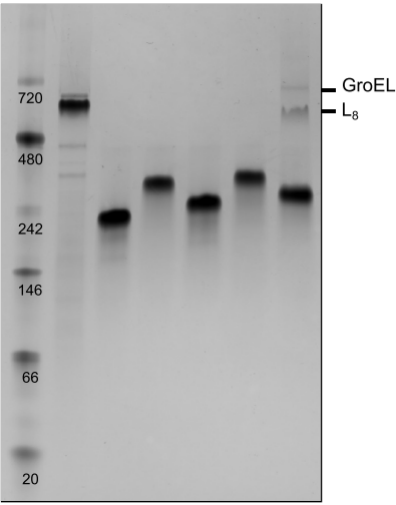


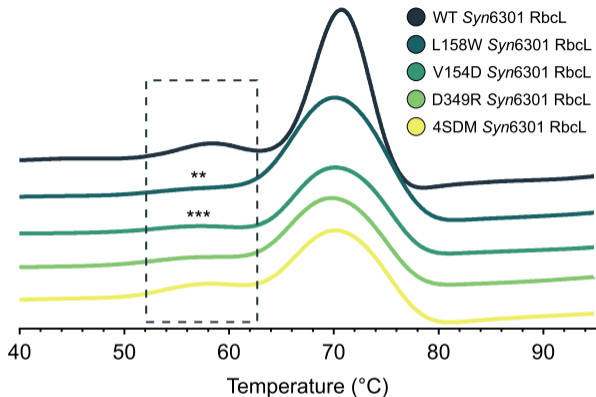
d



M
(kDa)

WT K150A D161A W165A D220A Y224A



a

Rubsico	T _m (°C)	n	p-value
WT <i>Syn6301 RbcL</i>	58.6 ± 0.2	4	ref.
L158W <i>Syn6301 RbcL</i>	57.7 ± 0.3	4	0.0037
V154D <i>Syn6301 RbcL</i>	57.3 ± 0.3	4	0.0007
D349R <i>Syn6301 RbcL</i>	57.9 ± 0.9	4	ns
4SDM <i>Syn6301 RbcL</i>	58.6 ± 0.1	4	ns

b