

# 1 $\alpha$ -carboxysome formation is mediated by the multivalent 2 and disordered protein CsoS2

3  
4 **Luke M. Oltrogge<sup>a</sup>, Thawatchai Chaijarasphong<sup>a,‡</sup>, Allen W. Chen<sup>b</sup>, Eric R. Bolin<sup>c,d</sup>, Susan  
5 Marqusee<sup>a,b,d</sup>, David F. Savage<sup>a,\*</sup>**

6  
7  
8  
9 <sup>a</sup>Department of Molecular and Cell Biology, University of California, Berkeley, CA 94720;

10 <sup>b</sup>Department of Chemistry, University of California, Berkeley, CA 94720

11 <sup>c</sup>Biophysics Graduate Program, University of California Berkeley, Berkeley, CA 94720

12 <sup>d</sup>California Institute for Quantitative Biosciences, University of California Berkeley, Berkeley, CA  
13 94720

14  
15 <sup>‡</sup>Present address: Department of Biotechnology, Faculty of Science, Mahidol University, Rama  
16 VI Rd., Bangkok 10400, Thailand

17  
18 \*To whom correspondence should be addressed: [savage@berkeley.edu](mailto:savage@berkeley.edu)  
19

20

## 21 **Abstract:**

22 Carboxysomes are bacterial microcompartments that function as the centerpiece of the  
23 bacterial CO<sub>2</sub>-concentrating mechanism, feeding high concentrations of CO<sub>2</sub> to the enzyme  
24 Rubisco for fixation. The carboxysome self-assembles from thousands of individual proteins into  
25 icosahedral-like particles with a dense enzyme cargo encapsulated within a proteinaceous shell.  
26 In the case of the  $\alpha$ -carboxysome, there is little molecular insight into protein-protein interactions  
27 which drive the assembly process. Here we show that the N-terminus of CsoS2, an intrinsically  
28 disordered protein found in the  $\alpha$ -carboxysome, possesses a repeated peptide sequence that  
29 binds Rubisco. X-ray structural analysis of the peptide bound to Rubisco reveals a series of  
30 conserved electrostatic interactions that are only made with properly assembled hexadecameric  
31 Rubisco. Although biophysical measurements indicate this single interaction is weak, its implicit  
32 multivalency induces high-affinity binding through avidity. Taken together, our results indicate  
33 CsoS2 acts as an interaction hub to condense Rubisco and enable efficient  $\alpha$ -carboxysome  
34 formation.  
35

## 36 Introduction:

37 Many carbon-assimilating bacteria possess CO<sub>2</sub>-concentrating mechanisms (CCMs) to  
38 facilitate carbon fixation by the enzyme Rubisco.<sup>1</sup> The centerpiece of the CCM is the  
39 carboxysome, a large protein complex which encapsulates Rubisco and carbonic anhydrase  
40 and is thought to produce locally high concentrations of CO<sub>2</sub>.<sup>2,3</sup> The carboxysome is a large  
41 (100–400 nm diameter) and composite (~10 different protomers) structure comprising both a  
42 virus-like protein shell and cargo enzymes.<sup>4–6</sup> Moreover, carboxysome formation requires  
43 thousands of individual proteins to accurately self-assemble.<sup>7–9</sup> How this mesoscopic complex,  
44 with linear dimensions roughly ten-fold larger than any of its individual components, assembles  
45 with high structural and compositional fidelity remains unknown.

46 Carboxysomes occur in two distinct evolutionary lineages,  $\alpha$  and  $\beta$ , that are functionally  
47 and morphologically similar.<sup>4,10,11</sup> Both enclose a dense enzymatic cargo of Rubisco and  
48 carbonic anhydrase inside the icosahedral shell composed of hexameric and pentameric  
49 proteins. One or more scaffolding proteins serve as interaction hubs, mediating the associations  
50 among the various components.<sup>4</sup>

51 Although the  $\alpha$ -carboxysome was the first to be identified and characterized,<sup>12</sup> the  $\beta$ -  
52 carboxysome assembly process is better understood. Two proteins, CcmM and CcmN, act in  
53 tandem as the scaffold in a hierarchical set of interactions to bridge shell with cargo.<sup>4,13</sup> An  
54 amphipathic encapsulation peptide on CcmN anchors to CcmK, a hexameric shell protein.<sup>14</sup>  
55 CcmN also binds to CcmM, a scaffolding protein which contains three to five tandem repeats of  
56 a Rubisco small subunit like (SSUL) module separated by disordered linkers. SSUL repeats  
57 then interact with Rubisco.<sup>15–18</sup> Contrary to expectations based on sequence homology, SSULs  
58 do not displace the Rubisco small subunit but bind across the interface of two L<sub>2</sub> dimers and a  
59 small subunit.<sup>17</sup>

60 The assembly of  $\alpha$ -carboxysomes—the predominant form among oceanic cyanobacteria  
61 and autotrophic proteobacteria—is, to date, more opaque. One unique component of the  $\alpha$ -  
62 carboxysome is CsoS2, a large (~900 residues) intrinsically disordered protein (IDP), which,  
63 unlike CcmM or CcmN, contains no recognizable domains.<sup>19,20</sup> CsoS2 is indispensable for  
64 carboxysome assembly and thus hypothesized to be a potential scaffolding protein. Knock-outs  
65 in the  $\alpha$ -carboxysome model organism *Halothiobacillus neapolitanus* produce high CO<sub>2</sub>-  
66 requiring phenotypes and result in no observable carboxysomes.<sup>19,21</sup> Pulldown and native  
67 agarose gel-shift assays using purified protein have demonstrated that CsoS2 interacts with  
68 both Rubisco and CsoS1 hexameric shell proteins.<sup>19,22–24</sup> The specific sites of interaction,  
69 however, have not been definitively determined nor is it clear how they collectively give rise to  
70 robust assembly.

71 Here, we show that a repeated peptide motif in the N-terminal domain of CsoS2 interacts  
72 with Rubisco to facilitate encapsulation into the carboxysome. Using a fusion of this peptide with  
73 Rubisco we obtained a structure of the binding site which revealed a predominantly electrostatic  
74 interaction interface mediated by highly conserved residues. This binding site lies at a  
75 conjunction of Rubisco subunits uniquely present in the complete L<sub>8</sub>S<sub>8</sub> oligomer, thus ensuring  
76 the encapsulation of only the functional holoenzyme. Energetic characterization indicated that  
77 the individual peptide/Rubisco interaction is very weak and relies on the engagement of multiple  
78 binding sites to increase its interaction strength. Bioinformatic analysis and expression of  
79 CsoS2-truncated heterologous carboxysomes implicate the multivalency of this interaction as an

80 essential feature of the assembly process. Our data suggest that CsoS2 acts as a protein  
81 interaction hub which gathers Rubisco to nascent carboxysome shell facets through branching  
82 low-affinity interactions that collectively give rise to efficient and robust cargo accumulation.

83

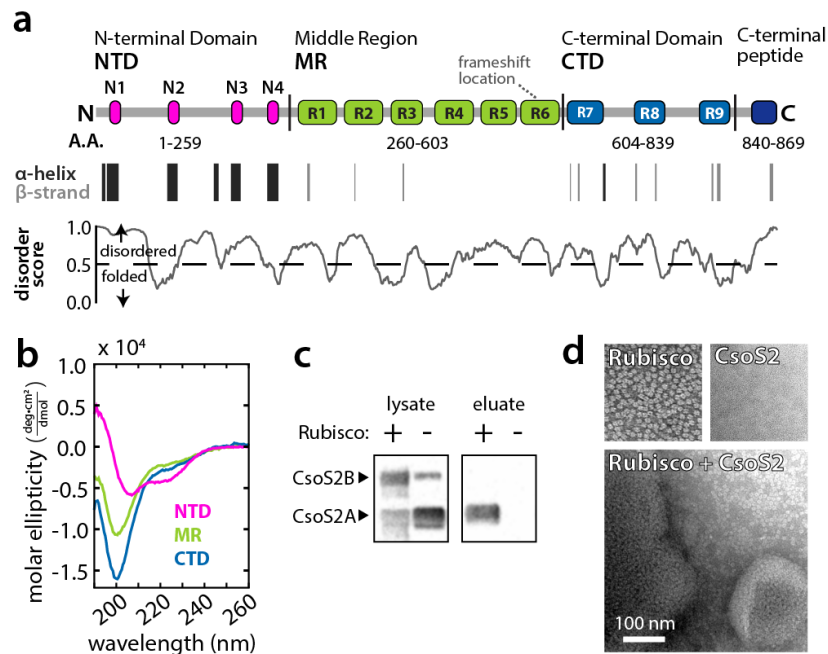
84

85 **Results:**

86 *CsoS2 interacts with Rubisco*

87 We and others have demonstrated the essentiality of CsoS2 to  $\alpha$ -carboxysome  
 88 formation.<sup>19,21</sup> This fact, in combination with CsoS2's unique sequence characteristics,<sup>20</sup> led us  
 89 to consider whether it is the scaffolding protein driving assembly of the  $\alpha$ -carboxysome. CsoS2  
 90 is a repetitive IDP.<sup>19,25</sup> It can be divided into three major domains, the N-terminal domain (NTD),  
 91 Middle region (MR), and C-terminal domain (CTD), based on sequence self-similarity of the  
 92 repeated motifs contained therein.<sup>19</sup> The full protein has a high PONDR-FIT disorder score<sup>26</sup>  
 93 throughout (average = 0.63, >0.5 predicts disordered) and is only predicted to possess  
 94 secondary structure within the repeats of the NTD (hereafter generically referred to as the 'N-  
 95 peptide' or specifically by numbers, e.g. N1 through N4; Fig. 1a).<sup>27</sup> Circular dichroism (CD)  
 96 spectra indicated that only the NTD has  $\alpha$ -helical content (Fig. 1b). However, the repeat  
 97 sequences in the NTD do not necessarily coincide with regions of greater predicted order. It is  
 98 thus possible that the N-peptides are in dynamic equilibrium between helical and unstructured  
 99 conformations.

100



101

102 **Figure 1**

103 **a**, Repeat structure of *H. neapolitanus* CsoS2 with secondary structure prediction and disorder scores.  
 104 "Frameshift location" indicates the site of a programmed ribosomal frameshift which results in expression  
 105 of about 50% prematurely truncated protein (CsoS2A) and 50% full-length protein (CsoS2B).<sup>21</sup> **b**, Circular  
 106 dichroism spectra of each of the CsoS2 domains. **c**, Anti-His Western blot against His-tagged CsoS2  
 107 expressed +/- Strep-tagged Rubisco in the raw lysate and following Strep affinity purification (eluate). **d**,  
 108 Negative stain TEM micrographs of purified Rubisco, CsoS2, and the aggregates observed when mixed.

109

110 Rubisco and CsoS2 together constitute a significant fraction of the cargo mass in  
111 purified carboxysomes and have complementary isoelectric points (5.9 and 9.1, respectively)  
112 suggesting a possible electrostatic association.<sup>5</sup> We therefore tested whether these two proteins  
113 physically interact via pull-down analysis. As hypothesized, affinity purification of Strep-tagged  
114 Rubisco pulled down a 6xHis-tagged CsoS2 when visualized by anti-His Western blotting (Fig.  
115 1c). This result pointed toward a direct interaction between CsoS2 and Rubisco and  
116 corroborated prior evidence.<sup>19</sup> Furthermore, we observed dense aggregates of CsoS2 and  
117 Rubisco by transmission electron microscopy (TEM) when the two proteins were co-incubated  
118 (Fig. 1d).

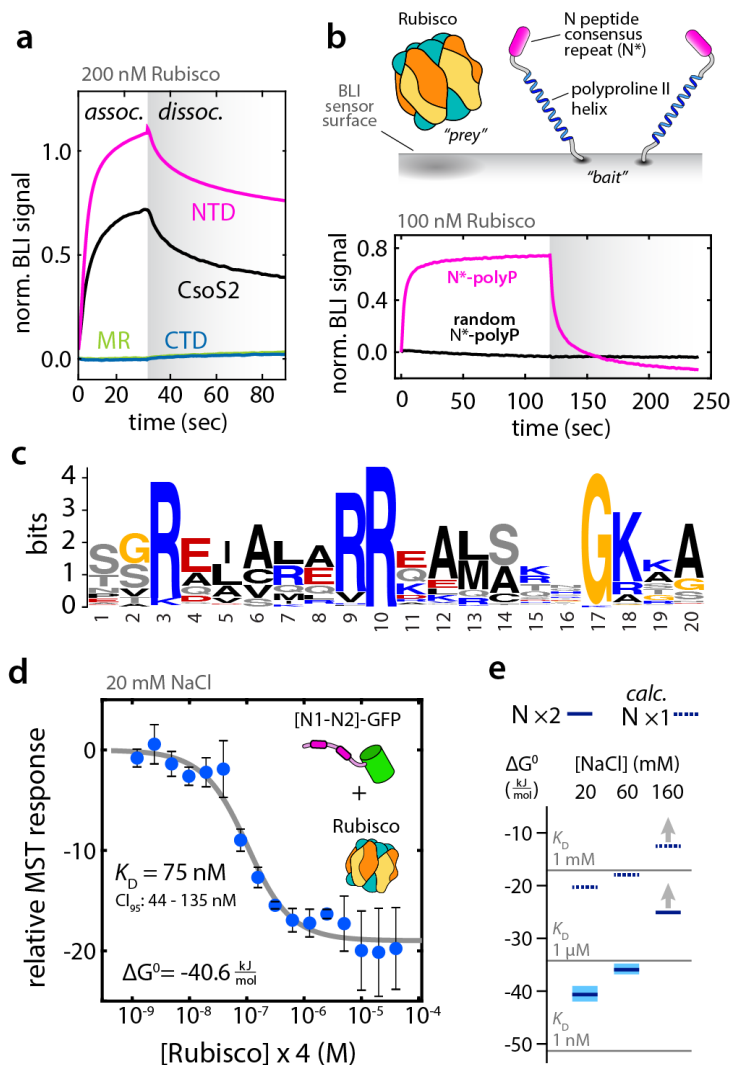
119

### 120 *Repeated NTD motif binds Rubisco with low affinity*

121 We next sought to identify the specific element of CsoS2 capable of interacting with  
122 Rubisco. This was carried out using bio-layer interferometry (BLI)—a label-free optical  
123 technique that monitors recruitment of a “prey” protein by a surface-immobilized “bait.”<sup>28</sup> BLI  
124 analysis on CsoS2 and its various fragments revealed that binding activity resided in the NTD  
125 (Fig. 2a). IDPs often interact with their targets through short linear motifs<sup>29,30</sup> and further  
126 analysis demonstrated that a single peptide derived from the consensus sequence of N1-N4,  
127 which we term N\* (with sequence GRDLARREALSQQGKAAV), was capable of interacting  
128 with Rubisco. A randomized sequence of N\* (GRRKGLRAAGRALQVEQADSRA) did not bind  
129 (Fig. 2a,b), nor did any of the other conserved peptides from the MR or CTD (Fig. S2),  
130 suggesting that the interaction was indeed sequence specific and not, for example, due to  
131 generic charge-charge attraction.

132 The interaction appeared to be driven by a specific sequence of positively charged  
133 residues. We analyzed a set of 231 CsoS2 sequences from  $\alpha$ -cyanobacteria and proteobacteria  
134 with  $\alpha$ -carboxysomes to identify the pan-species consensus N-peptide motif (Fig. 2c),  
135 recapitulating previous results.<sup>19</sup> Notably, among the most highly conserved positions in the N-  
136 peptide motif are basic residues at positions 3, 9, 10, and 18, implying that the interaction likely  
137 has significant ionic character. R to A mutations were made for positions 3 and 10 in all of the  
138 four repeats in the NTD and entirely eliminated the binding in BLI (Fig. S3). Furthermore, a  
139 retrospective statistical examination of CsoS2 peptide array binding data from Cai et al.<sup>19</sup>  
140 revealed a significant enrichment of Rubisco binding to peptides matching the N-peptide  
141 arginine motif (Fig. S7).

142 In principle, the binding energy between Rubisco and the N-peptide should be calculable  
143 from fitting the association and dissociation kinetics. However, due to the inherently high  
144 valency of the  $L_8S_8$  Rubisco complex and the surface-induced avidity of neighboring bait  
145 proteins, it was difficult to obtain reliable fits to a simple binding model (Fig. S1). For this reason,  
146 the solution-phase technique microscale thermophoresis (MST) was used to measure binding in  
147 an alternative fashion. Unexpectedly, while the implied dissociation constants ( $K_D$ 's) from BLI  
148 were in the tens of nM regime, MST revealed no apparent binding under the same conditions  
149 (pH 7.5, 150 mM NaCl) (e.g. Fig S5a). Decreasing the salt to 20 mM NaCl, however, resulted in  
150 robust binding of a tandem N-peptide-GFP species, [N1-N2]-GFP, to Rubisco with a  $K_D$  of 75  
151 nM on a stoichiometric binding site basis (i.e. one [N1-N2]-GFP binds to two of eight sites per  
152 Rubisco) (Fig. 2d).



153

## 154 Figure 2

155 **a**, Bio-layer interferometry (BLI) Rubisco binding response normalized to the bait loading signal for full-  
 156 length CsoS2 and each of the domains. **b**, Upper panel: schematic of the BLI sensor surface with the N\*-  
 157 peptide displayed on an extended polyproline II helix as the bait and Rubisco as the prey species. Lower  
 158 panel: BLI response shows active binding of Rubisco by N\* but not by a scrambled version. **c**, Weblogo  
 159 conservation of the N-peptide motif calculated by MEME<sup>31</sup> from 231 CsoS2 sequences which contained  
 160 901 N-peptide occurrences. **d**, Microscale thermophoresis (MST) binding isotherm with the first two *H.*  
 161 *neapolitanus* CsoS2 N-peptide repeats fused to GFP, [N1-N2]-GFP, as the target and Rubisco as the  
 162 ligand. The abscissa represents the concentration of binding sites for [N1-N2]-GFP, i.e. four per Rubisco.  
 163 Error bars indicate +/- one standard deviation for measurements performed in triplicate. 95% confidence  
 164 interval ( $CI_{95}$ ) estimated by bootstrap analysis. **e**, Standard free energies of binding for the reaction in (d)  
 165 calculated from binding isotherms at 20, 60, and 160 mM NaCl. Solid dark blue lines are measured for  
 166 [N1-N2]-GFP with light blue spanning the 95% confidence interval. Dashed blue lines are calculated  
 167 estimates of the binding energy of a single repeat to a single Rubisco binding site. At 160 mM NaCl, no  
 168 binding could be detected and the lines represent lower limits of the  $K_D$ .

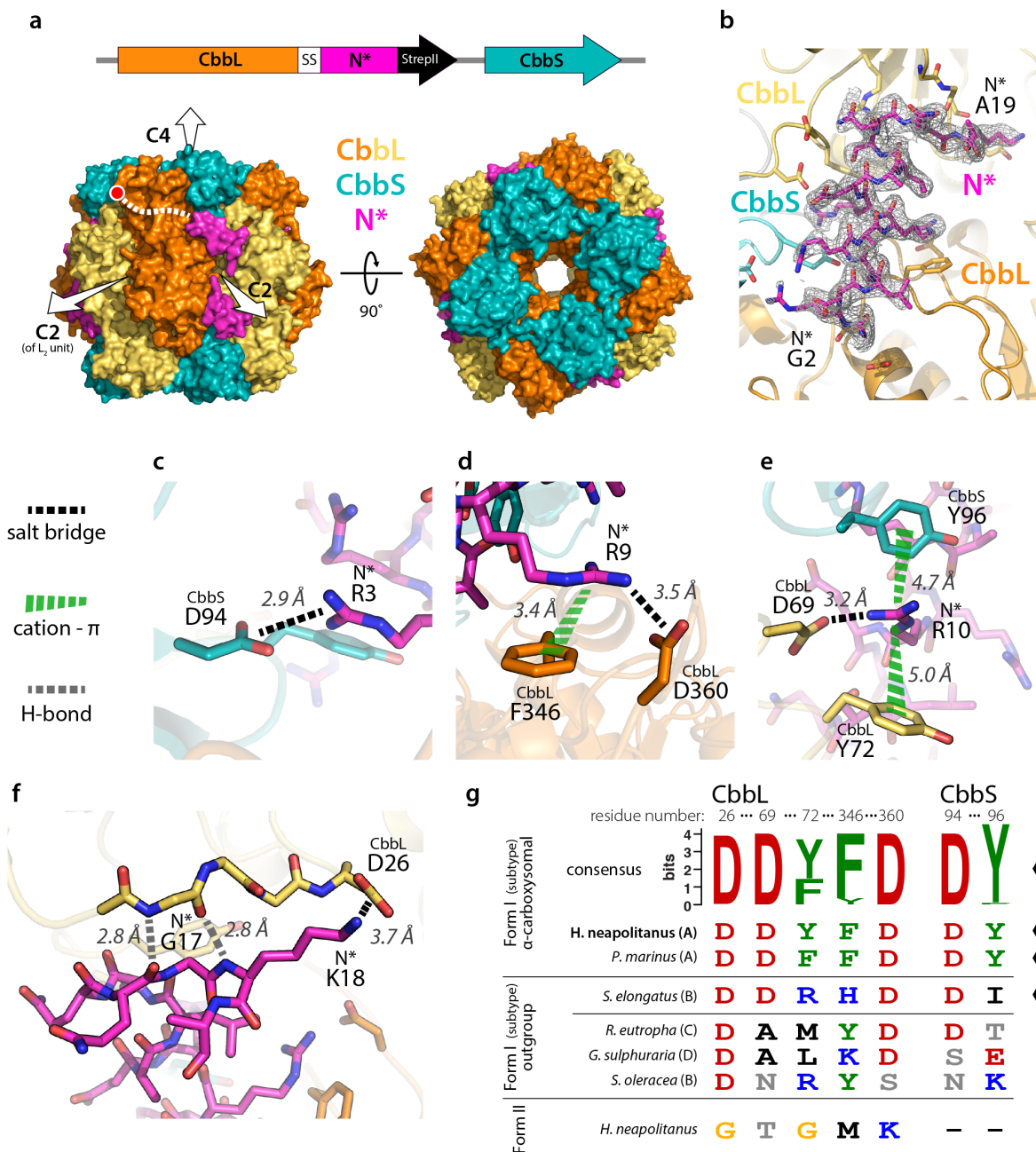
169  
170 MST indicated the N-peptide/Rubisco interaction is highly sensitive to salt concentration.  
171 Increasing NaCl from 20 mM to 60 mM showed a substantial increase in the  $K_D$  from 75 nM to  
172 500 nM (Fig. 2e). Further increasing NaCl to 160 mM—near physiological ionic strength<sup>32</sup> —  
173 weakened the binding beyond detection. Assuming a linear free energy relationship, we can  
174 estimate the binding energy for the individual N-peptide to Rubisco to be half of the  $\Delta G^0$  for the  
175 [N1-N2]-GFP construct leading to  $K_D \sim 250 \mu\text{M}$  at 20 mM NaCl (see SI, *MST fitting and*  
176 *analysis*). Indeed, MST of a single N-peptide-GFP, [N1]-GFP, showed no discernable binding  
177 over the same concentration range (Fig. S5b).

178 Taken together, these data present two puzzling observations. First, the individual N-  
179 peptide/Rubisco interaction alone appears too weak to drive carboxysome cargo encapsulation,  
180 particularly when approaching realistic intracellular ionic strength. Second, the relatively tight  
181 binding of Rubisco by a single N-peptide construct at 150 mM NaCl on BLI stands in apparent  
182 contradiction to the negative binding results obtained from MST under similar conditions. A  
183 mechanistic reconciliation of these issues is presented in the Discussion.

#### 184 185 *Structural determination of the N-peptide/Rubisco complex*

186 We next sought to obtain a structure of the N-peptide/Rubisco complex in order to locate  
187 the binding sites and to establish the nature of the specific molecular contacts. The NTD is  
188 largely disordered and its four N-peptide repeats could, in principle, adopt heterogeneous  
189 arrangements among the eight Rubisco binding sites. Furthermore, the binding of a single N-  
190 peptide is weak and salt sensitive. Disorder, structural heterogeneity, and partial occupancy  
191 therefore all pose significant challenges for co-crystallization. To circumvent these problems, we  
192 fused the N\* consensus peptide to the C-terminus of the Rubisco large subunit (CbbL) via a  
193 short linker, -SS-, (Fig. 3a) to insure high local concentrations and saturation of all putative  
194 binding sites. This fusion protein was readily expressed, purified and was confirmed by size  
195 exclusion chromatography to be of the correct  $L_8S_8$  oligomerization state (Fig. S6a). BLI  
196 measurements revealed no significant interaction of the Rubisco-N\* fusion (prey) to surface N\*  
197 peptide (bait) suggesting that Rubisco-N\* self-passivates its binding site (Fig. S6b,c).

198 After screening and optimization of crystallization conditions, diffraction quality crystals  
199 were obtained (Table S1). X-ray diffraction data were collected and the structure was solved by  
200 molecular replacement using an existing model from Kerfeld and Yeates of *H. neapolitanus*  
201 Rubisco (PDB: 1SVD). The space group was  $C_2$  with four CbbL-N\* and CbbS chains in the  
202 asymmetric unit. The Rubisco structure itself was essentially indistinguishable from wild-type  
203 with an average  $C\alpha$  RMSD of 0.27 Å. Clear unmodeled electron density was observed along the  
204 groove at the interface between two CbbL subunits (spanning separate  $L_2$  dimers) and a CbbS  
205 subunit (Fig. 3a). The N\*-peptide was found to adopt a helical conformation and an all-atom  
206 model was manually built into the experimental density, which was sufficiently clear for  
207 unambiguous assignment of both the peptide direction and sequence registration. Following  
208 several rounds of refinement, the real-space cross-correlation for the modeled portion of N\*  
209 (res. 2-19, Fig. 2c) was 90% or greater for each of the four N\*-peptides in the asymmetric unit  
210 (Fig. 3b). All of the binding sites are occupied, indicating that the neighboring sites are not  
211 mutually occluding. Thus, it is likely that the  $L_8S_8$  biological assembly possesses eight possible  
212 CsoS2 interaction sites.



213

### 214 Figure 3

215 **a**, Schematic of the Rubisco-N\* fusion construct and side and top views of a surface representation of the  
 216  $L_8S_8$  biological assembly with bound N\*-peptide. CbbL and CbbS are the large and small Rubisco  
 217 subunits, respectively. The molecular symmetry axes are indicated by white arrows. The yellow and  
 218 orange CbbLs are identical; the coloring is to highlight the  $L_2$  dimer units. The red dot is at the last  
 219 structured residue of CbbL, while the dashed white line indicates the probable linkage to N\*. **b**, Zoomed  
 220 view of binding site with  $2F_o-F_c$  map at  $\sigma = 0.8$  carved within 1.6 Å of N\*. The first and last structured  
 221 residues of N\* are labeled. **c-f**, Molecular interactions of each of the five highly conserved residues of the  
 222 N-peptide motif: R3, R9, R10, G17, and K18. Salt bridges, cation- $\pi$  interactions, and select hydrogen



223 bonds are specifically highlighted. The specific interactions were characterized with the PDBePISA<sup>33</sup> and  
224 CaPTURE<sup>34</sup> web servers. **g**, Rubisco sequence comparison at the N\*-peptide interaction site. The  
225 Weblogo conservation sequence is from 231  $\alpha$ -carboxysomal Form IA Rubiscos. Two specific  
226 representatives, *H. neapolitanus* (used in this study) and *Prochlorococcus marinus* MIT 9313, are shown.  
227 Below are various outgroup Form I Rubiscos and the *H. neapolitanus* Form II Rubisco. Participation in  
228 carboxysomes ( $\alpha$  or  $\beta$ ) is indicated along the right of the table. Note that the residues are non-sequential  
229 and are numbered according to the *H. neapolitanus* sequence.

230

231 The structure of the bound N\*-peptide is largely  $\alpha$ -helical, consistent with the secondary  
232 structure predictions and CD data (Fig. 1a,b). The last clearly structured residue of CbbL is at  
233 position 455, which is typical of structures of non-activated Form I Rubisco.<sup>35</sup> The remainder of  
234 the CbbL C-terminus and the -SS- linker preceding N\* are not observed in the electron density.  
235 Although lack of density complicates the assignment of N\*/CbbL pairings, the structured portion  
236 of N\* begins near CbbL helix 6 and the fusion thus likely originates from the C-terminus of this  
237 same subunit. This also agrees with previous structural models of other Rubiscos, in which the  
238 C-terminus extends over the so-called loop 6 in the same direction as the N\* binding site (Fig  
239 3a, dashed white).<sup>35</sup> From there, the N\* helix makes contacts with CbbS, spans the boundary to  
240 the neighboring L<sub>2</sub> dimer, and finishes by breaking out of the helix at the N-terminal domain of  
241 the second CbbL. A noteworthy quality of the N\*/Rubisco binding site is that, by contacting both  
242 CbbL and CbbS and bridging the L<sub>2</sub> dimer interface, it exists only on the L<sub>8</sub>S<sub>8</sub> Rubisco  
243 holoenzyme. This fact implies that only fully assembled Rubisco would be admitted into the  
244 carboxysome.

245 Each one of the highly conserved N\* motif residues (Fig. 2c) is observed to make key  
246 binding contacts along the Rubisco interface. R3 is salt-bridged with CbbS D94 (Fig. 3c). R9  
247 forms a salt-bridge with CbbL D360 and cation- $\pi$  interaction with F346 (Fig. 3d). R10 has a salt-  
248 bridge to CbbL D69 and dual cation- $\pi$  interactions with CbbL Y72 and CbbS Y96 (Fig. 3e). G17  
249 appears to play a critical role in breaking the N\* helix by facilitating backbone hydrogen bonds  
250 with CbbL and adopting glycine-specific  $\psi$ - $\phi$  angles. Finally, K18 makes a salt bridge with CbbL  
251 D26 (Fig. 3f). All together the interactions are predominantly ionic and offer a structural  
252 explanation as to the energetic sensitivity to salt.

253 Amino acid residues involved in these electrostatic interactions are conserved for  $\alpha$ -  
254 carboxysomal Form IA Rubisco. However, these residues were, in general, not conserved  
255 among an outgroup of various other Form I Rubiscos and the *H. neapolitanus* Form II Rubisco  
256 (Fig. 3g). To assay if these evolutionary observations are significant, two binding site mutants  
257 were made to test disruption of the binding interface. In one, each of the cation- $\pi$  aromatics was  
258 mutated to alanine (CbbL Y72A, F346A; CbbS Y96A). In the other, a mutation was selected to  
259 resemble the  $\beta$ -carboxysomal Rubisco and to perturb the binding environment of N\* R10 (CbbL  
260 Y72R). Neither mutant interacted with N\* (Fig. S4).

261

### 262 *Structural comparison to CcmM/Rubisco*

263 The general binding site of N\*/Rubisco significantly overlaps with that of the recently  
264 determined CcmM/Rubisco interaction from the  $\beta$ -carboxysome, however, the specific molecular  
265 details are distinct.<sup>17</sup> While CcmM binds with multiple regions across the SSUL domain,<sup>17</sup> N\* has  
266 a smaller footprint as a single  $\alpha$ -helix (Fig. S10). In both cases, salt bridges—with the positive  
267 charge contributed by the scaffolding protein—are key parts of the interactions. A notable

268 feature of the N\*/Rubisco interaction, but absent in CcmM, are the prominent cation- $\pi$   
269 interactions.<sup>34</sup> The complete conservation of the aromatics in the Rubisco binding site and the  
270 lack of binding when mutated to alanines suggest that the cation- $\pi$  interactions indeed  
271 contribute meaningfully to the binding energy and specificity. Interestingly, cation- $\pi$  contacts are  
272 a particularly common interaction modality among IDPs involved in protein liquid-liquid phase  
273 separation.<sup>36–38</sup>

274

#### 275 *Hydrogen/deuterium exchange of carboxysomal versus purified Rubisco*

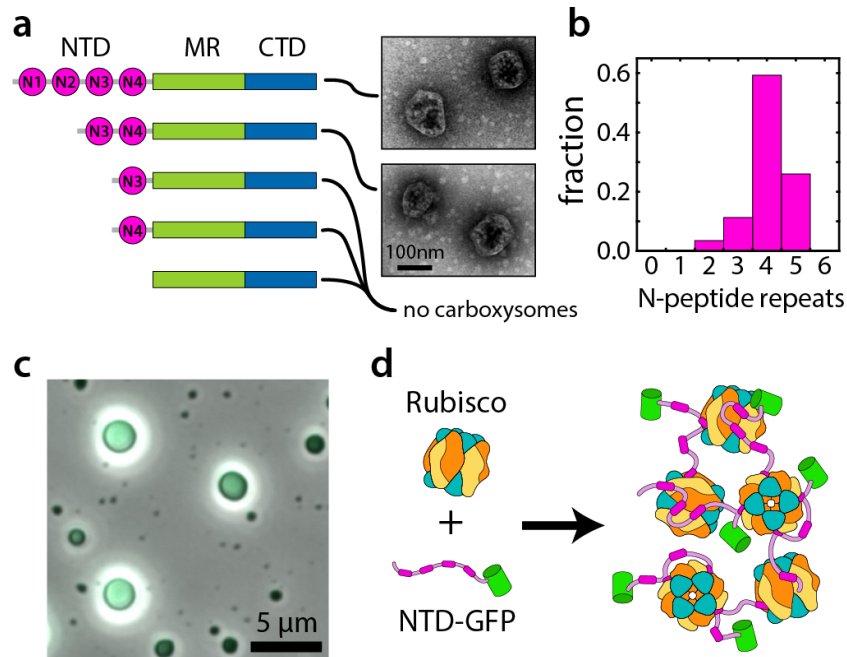
276 To interrogate the CsoS2/Rubisco interaction in a native context, hydrogen/deuterium  
277 exchange (HDX) mass spectrometry experiments were performed in order to identify regions of  
278 Rubisco possessing differential protection when encapsulated within carboxysomes. HDX  
279 analysis of purified Rubisco versus carboxysomal Rubisco revealed a majority of peptides had  
280 nearly identical HDX rates. The most notable exception was CbbL 328-341 on helix 6 which  
281 experienced significantly greater protection inside carboxysomes (Fig. S8). This peptide, while  
282 not directly contacting N\*, is connected through water-bridged hydrogen bond networks (Fig.  
283 S9). Although the NTD interaction does not apparently alter the crystal structure of Rubisco, it is  
284 possible that peptide binding may affect the dynamics of Rubisco structural elements.

285

#### 286 *Effect of N-peptide multivalency on carboxysome formation*

287 We set out to determine the importance of the number of N-peptide repeats on  
288 carboxysome assembly. *H. neapolitanus* CsoS2 contains four copies of the repeat but there is  
289 likely significant natural diversity. To this end, the consensus motif was used to quantify  
290 occurrences throughout the set of 231 CsoS2 genes.<sup>39</sup> Every sequence contained at least two  
291 copies of the motif suggesting that a valency greater than one may be a general requirement for  
292 carboxysome assembly (Fig. 4b). Using a previously developed method whereby carboxysomes  
293 are produced heterologously in *E. coli* by expressing the known genes from a single plasmid  
294 (pHnCB10),<sup>40</sup> we tested the effect of N-peptide repeat number on carboxysome formation. A  
295 series of pHnCB10 constructs were made possessing CsoS2 variants with a decreasing number  
296 of N-peptide repeats and tested for carboxysome expression. Only CsoS2 variants with two or  
297 more repeats were capable of forming carboxysomes (Fig. 4a and Fig. S11), consistent with the  
298 bioinformatic result.

299



300

### 301 **Figure 4**

302 **a**, Truncated CsoS2 proteins with variable numbers of N-peptide repeats and TEM images of the resulting  
 303 carboxysomes if any were formed. **b**, Histogram of N-peptide repeat numbers across 231 CsoS2  
 304 sequences. **c**, Merged GFP fluorescence and phase contrast images of protein liquid-liquid droplets  
 305 formed from a solution of Rubisco and NTD-GFP. **d**, Microscopic model of the phase separated state.  
 306 The branching of interactions due to the multivalency of both components provides the liquid cohesion  
 307 while the relative weakness and exchangeability of the individual interactions confers fluidity.

308

#### 309 *Phase separation of Rubisco and NTD*

310 IDPs are highly represented in systems that undergo protein liquid-liquid phase  
 311 separation. The propensity toward phase separation is promoted by weak individual  
 312 interactions, often salt sensitive, and multivalent association either through well-defined binding  
 313 sites or via less specific interactions related to the general amino acid composition.<sup>41,42</sup> Phase  
 314 separation has recently emerged as a common theme for the organization of Rubisco into CCM  
 315 architectures. In the algal pyrenoid, Rubisco phase separates with EPYC1, a repetitive IDP.<sup>43–45</sup>  
 316 From  $\beta$ -carboxysomes the short form of the scaffold protein CcmM, M35,<sup>46</sup> was shown to demix  
 317 with Rubisco into protein liquid droplets.<sup>17</sup> We hypothesized that CsoS2 and, in particular, the  
 318 NTD may similarly demix with Rubisco. Indeed, when Rubisco and NTD-GFP are combined at  
 319 1.0  $\mu$ M each at low salt (20 mM NaCl) the solution became turbid. Imaging by phase contrast  
 320 and epifluorescence microscopy revealed that round green fluorescent droplets are formed (Fig.  
 321 4c) and are fully re-dissolved upon salt addition up to 150 mM NaCl. No droplets are observed  
 322 with either individual component at the same concentrations.

323

## 324 Discussion:

325 We have characterized in molecular detail the binding interface of Rubisco and CsoS2  
326 which facilitates  $\alpha$ -carboxysome cargo encapsulation. CsoS2, as a large IDP, posed a  
327 significant challenge for structural determination. Through biophysical binding assays we  
328 narrowed down the interaction to a repeated motif within the CsoS2 NTD, fused this fragment  
329 directly to Rubisco, and obtained an x-ray crystal structure of the protein-peptide complex. We  
330 suggest that this workflow might be a valuable general strategy for determining structures of  
331 IDPs interacting with structured proteins since these interactions are often individually weak and  
332 transient.

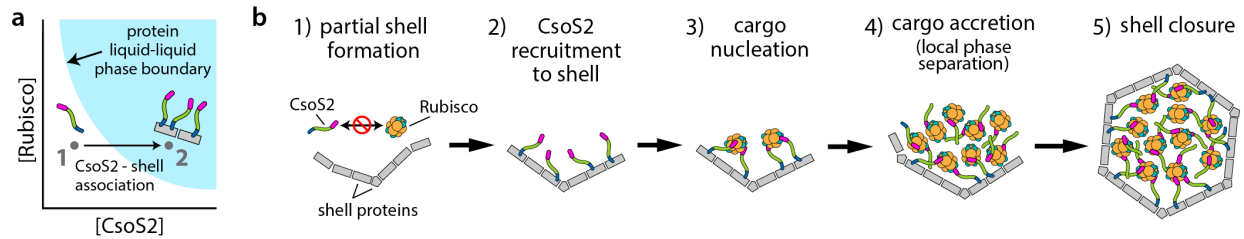
333 Despite no apparent sequence similarity, the CsoS2/Rubisco binding bears striking  
334 parallels to the recently characterized CcmM/Rubisco interaction at the heart of  $\beta$ -carboxysome  
335 assembly.<sup>17</sup> In both cases the scaffold protein binding element has multiple repeats interspersed  
336 by flexible linkers. The binding locations on Rubisco are very similar; both straddle an L<sub>2</sub> dimer  
337 interface while also making critical contacts with a small subunit. This site is only present in the  
338 fully assembled L<sub>8</sub>S<sub>8</sub> Rubisco holoenzyme so Rubisco assembly intermediates, namely L<sub>2</sub> and  
339 (L<sub>2</sub>)<sub>4</sub>, would presumably not be encapsulated prematurely. Notwithstanding this global similarity,  
340 the specific structural details of the binding are distinct, making this an intriguing example of  
341 convergent evolution.

342 Another commonality between the  $\alpha$ - and  $\beta$ -carboxysome scaffold/Rubisco systems is  
343 the propensity to undergo protein liquid-liquid phase separation. Phase separation is  
344 increasingly understood to play an organizational role in eukaryotes in the formation of  
345 membrane-less organelles.<sup>47</sup> These structures and the droplets we observe (Fig. 4c), however,  
346 have at least a thousand-fold greater volume than carboxysomes. Furthermore, they are not  
347 enclosed within protein shells. Therefore, while suggestive of a dense liquid cargo phase, the  
348 role of demixing in the carboxysome assembly process remains unresolved.

349 The N-peptide/Rubisco interface is comprised chiefly of salt bridges and cation- $\pi$   
350 interactions. Consequently, the binding energy is highly sensitive to the solution ionic strength.  
351 Indeed, our solution phase binding measurements with MST indicate that the interaction  
352 dramatically weakens, with single site  $K_D$ 's greater than 1 mM, at near-physiological ionic  
353 strength. Moreover, the phase separated droplets are fully dissolved under the same elevated  
354 salt concentrations. In apparent contradiction, however, the BLI measurements under the same  
355 conditions indicated strong binding ( $K_D \sim 100$  nM).

356 The essential difference is that BLI is a surface-based technique. Since the "prey"  
357 Rubisco has a site valency of eight, it could be simultaneously engaged by multiple "bait" N\*-  
358 peptides in microscopically dense patches on the surface (see SI, *Comments on BLI*). This  
359 surface avidity effect enabled tight Rubisco binding even when the individual interactions were  
360 very weak. We propose that this artificial surface avidity represents a useful analogy to the early  
361 stages of carboxysome assembly. Several experiments have implicated CsoS2 association with  
362 the CsoS1 shell hexamer including native gel shifts<sup>19</sup> and pulldown assays.<sup>22</sup> Furthermore, the  
363 CsoS2 C-terminus was found at the shell<sup>25</sup> and truncation of the CTD precludes carboxysome  
364 formation.<sup>21</sup> Through the shell interaction, multiple CsoS2 molecules could be recruited to  
365 achieve high local concentration and then bind to Rubisco in a multivalent fashion with high  
366 affinity.

367



368

## 369 **Figure 5**

370 **a**, Model phase diagram of the hypothesized Rubisco/CsoS2 phase separation driven by the multivalent  
371 NTD interaction with Rubisco. The blue region represents the joint concentrations at which demixing  
372 occurs. At point 1 the cytosolic concentrations lie within the soluble region and both are fully dissolved.  
373 Through interactions with a nascent carboxysome shell, multiple CsoS2s are brought together, thus  
374 greatly increasing the concentration locally while the Rubisco concentration remains the same (point 2).  
375 This process locally exceeds the phase transition threshold and leads to local phase separation in the  
376 immediate vicinity of the shell. **b**, Model of  $\alpha$ -carboxysome assembly in which the specific accumulation of  
377 cargo on the shell proceeds via the mechanism described in (**a**).  
378

378

379 Our data have led us to the following speculative model of  $\alpha$ -carboxysome assembly: At  
380 physiological ionic strength and the likely free concentrations of Rubisco and CsoS2 the  
381 interaction is insufficiently strong to drive significant association or demixing (Fig. 5a, point 1).  
382 However, in the presence of shell proteins, CsoS2 is gathered to high local concentration via  
383 interaction to the nascent shell surface and facilitates phase separation with Rubisco in the  
384 immediate vicinity of the shell (Fig. 5a, point 2). Eventually more shells with cargo droplets  
385 coalesce until the structure is fully enclosed.

386 A full accounting of the interaction partners and the site binding energetics is alone  
387 insufficient to understand the carboxysome assembly process. Multivalency, surface avidity,  
388 protein liquid-liquid phase separation appear to play important roles but their relationships to the  
389 shell and the emergent size regularity remain unclear and warrant further investigation.  
390 Ultimately a detailed understanding of the principles of carboxysome assembly may be  
391 leveraged toward the design of synthetic microcompartments for biotechnological applications.  
392

392

393

394 **Methods:**

395

396 *Protein expression and purification*

397 All proteins used for biochemical assays contained a terminal affinity tag, either a  
398 hexahistidine tag or a Strep-tag II (see SI for complete sequences). Each construct was cloned  
399 via Golden Gate assembly<sup>48</sup> into a pET-14 based destination vector with ColE1 origin, T7  
400 promoter, and carbenicillin resistance. These were transformed into *E. coli* BL21-AI expression  
401 cells. All Rubisco constructs were also co-transformed with pGro7 for expressing GroEL-GroES  
402 to facilitate proper protein folding. Cells were grown at 37°C to OD600 of 0.3-0.5 in 1 L of LB  
403 media before lowering the temperature to 18°C, inducing with 0.1% (w/v) L-arabinose, and  
404 growing overnight.

405 Cultures were harvested by centrifugation at 4,000 *g* and the pellets were frozen and  
406 stored at -80°C. The pellets were thawed on ice and resuspended with ~25 mL of lysis buffer  
407 (50 mM Tris, 150 mM NaCl, pH 7.5) supplemented with 1 mM phenylmethanesulfonyl fluoride  
408 (PMSF), 0.1 mg/mL lysozyme, and 0.01 mg/mL DNaseI. The cells were lysed with three passes  
409 through an Avestin EmulsiFlex-C3 homogenizer and clarified by centrifugation at 12,000 *g* for  
410 30 min. The clarified lysate was then incubated with the appropriate affinity resin for 30 min at  
411 4°C with 2 mL of resin per 1 L of initial culture and transferred to a gravity column. His-tagged  
412 proteins were bound to HisPur Ni-NTA resin (Thermo), washed with lysis buffer with 30 mM  
413 imidazole, and eluted with lysis buffer with 300 mM imidazole. Strep-II-tagged proteins were  
414 bound to Strep-Tactin resin (EMD Millipore), washed with lysis buffer, and eluted with lysis  
415 buffer containing 2.5 mM desthiobiotin. All proteins were buffer exchanged to lysis buffer with  
416 10DG Desalting Columns (Bio-Rad). For storage, proteins were made to 10% (w/v) glycerol,  
417 flash frozen in liquid nitrogen, and stored at -80°C.

418 Protein purity was assessed by SDS-PAGE gel analysis. In general all protein was >90%  
419 the desired product. Size exclusion chromatography was performed analytically to confirm purity  
420 and aggregation state and, if needed, as a final preparative step.

421

422 *Bio-layer interferometry*

423 Protein-protein interactions were measured using bio-layer interferometry (BLI) with an  
424 Octet RED384 (Forte Bio). The “bait” protein was immobilized on Ni-NTA Dip and Read  
425 Biosensors via a terminal His-tag. Typical “bait” concentrations for the sensor loading were 10  
426 µg/mL. The soluble “prey” protein concentrations were varied in the nanomolar to micromolar  
427 range. The buffer used for all loading, association/dissociation, and wash steps was 50 mM Tris,  
428 150 mM NaCl, 0.01% (w/v) Triton X-100, pH 7.5. Sensor regeneration of the Ni-NTA was done  
429 with 50 mM Tris, 150 mM NaCl, 0.05% (w/v) SDS, 300 mM imidazole, pH 7.5. The typical  
430 experimental binding sequence used was: load “bait”, buffer wash, “prey” association, “prey”  
431 dissociation in buffer, sensor regeneration, buffer wash. For the experiments testing the binding  
432 activity of specific peptides (Fig. 2b and Fig. S2), “bait” proteins were designed with a 40 amino  
433 acid proline rich region between the His-tag and the peptide (see SI, *Protein Sequences*). This  
434 insertion is expected to adopt an extended polyproline II helix conformation ~10 nm in length<sup>49</sup>  
435 and was included to limit possible surface occlusion.

436

437

#### 438 *Microscale thermophoresis*

439 Solution protein-protein binding was monitored by microscale thermophoresis (MST)  
440 with a Monolith NT.115 (Nanotemper). The target proteins were portions of the CsoS2 NTD  
441 fused to Superfolder GFP and used at a concentration of 50 nM. Unlabeled Rubisco was used  
442 as the ligand with concentrations varied in two-fold increments from 10  $\mu$ M (as  $L_8S_8$ ) down to 0.3  
443 nM. Experiments were carried out in buffer with 6.7 mM Tris, 0.01% Triton X-100, pH 7.5 and  
444 either 20, 60, or 160 mM NaCl. The samples were loaded into MST Premium Coated Capillaries  
445 (Nanotemper) and analyzed using 20% blue LED power for fluorescence excitation and Medium  
446 infrared laser power for the thermophoresis. Data fitting and bootstrap error estimation was  
447 performed using custom scripts in MATLAB (MathWorks).

448

#### 449 *Crystallization, x-ray diffraction, and refinement*

450 Initial screening of crystallization conditions for CbbL-N\*, CbbS was done using the  
451 Hampton Crystal Screen (HR2-110) with protein at 15 mg/mL combined 1:1 with the screen  
452 mother liquors. Due to the hypothesized ionic nature of the interaction, screen conditions having  
453 lower salt concentrations were prioritized in the follow-up optimization. Ultimately the best  
454 crystals were obtained from a mother liquor of 0.2M  $MgCl_2 \bullet 6H_2O$ , 0.1M HEPES, 30% (v/v)  
455 PEG-400. Protein at 15 mg/mL diluted 1:2 with mother liquor was allowed to equilibrate for one  
456 day by hanging drop vapor diffusion whereupon it was microseeded with pulverized crystals  
457 from more concentrated conditions delivered with a cat whisker.

458 Crystals were looped and directly frozen on the beamline under a 100K nitrogen jet  
459 without additional cryoprotectant. X-ray diffraction was collected with wavelength 1.11  $\text{\AA}$  on a  
460 Pilatus3 S 6M (Dectris) detector with a 50 $\mu$ m beam pinhole at the Advanced Light Source, BL  
461 8.3.1, Berkeley, CA.

462 The data were indexed and integrated with XDS<sup>50</sup> and scaled and merged with  
463 AIMLESS.<sup>51,52</sup> Molecular replacement was carried out in Phenix using the existing wild-type *H.*  
464 *neapolitanus* Rubisco structure (PDB ID: 1SVD) as the search model.<sup>53,54</sup> Cycles of automatic  
465 refinement were performed with Phenix while Coot was used for manual model building.<sup>55</sup> The  
466 final refined structure backbone conformations were 96.0% Ramachandran favored, 3.8%  
467 allowed, and 0.2% outliers.

468

#### 469 *Carboxysome construct generation and purification*

470 Heterologous expression of carboxysomes in *E. coli* was performed following the  
471 methods of Bonacci et al. using the plasmid pHnCB10 which contains genes encoding all ten of  
472 the proteins known to participate in carboxysome formation.<sup>40</sup> Golden Gate assembly was used  
473 to make the truncations of the CsoS2 NTD shown in Fig. 4a.

474 Carboxysomes were purified as previously described.<sup>21</sup> Briefly, the cells were harvested,  
475 resuspended in 25mL TEMB buffer (10 mM Tris, 10 mM  $MgCl_2$ , 1 mM EDTA, and 20 mM  
476  $NaHCO_3$ , pH 8.4), lysed with a homogenizer, and the lysate clarified by centrifugation at 12,000  
477 g for 30 min. The supernatant was further centrifuged at 40,000 g for 30min to pellet the  
478 carboxysomes. The carboxysome pellet was resuspended in 1x Cellytic B (Sigma-Aldrich) in  
479 order to solubilize any residual membrane fragments. The solution was spun a second time at  
480 40,000 g for 30 min to pellet the carboxysomes again. The pellet was resuspended with 3mL of  
481 TEMB, clarified with a 5min spin at 3,000 g, and loaded on top of a 25-mL sucrose step gradient

482 (10, 20, 30, 40, and 50% w/v sucrose). This was ultracentrifuged at 105,000 g for 30 min. The  
483 solution was fractionated and analyzed by SDS-PAGE. Those fractions containing the expected  
484 set of carboxysomal proteins (and which also demonstrated visible Tyndall scattering) were  
485 pooled, pelleted by centrifugation for 90min at 105,000 g, resuspended in 1mL of TEMB, and  
486 stored at 4°C.

487

#### 488 *Negative stain TEM*

489 Purified carboxysomes were visualized by negative stain transmission electron  
490 microscopy. Formvar/carbon coated copper grids were prepared by glow discharge prior to  
491 sample application. The grids were washed with deionized water several times before staining  
492 with 2% (w/v) uranyl acetate. Imaging was performed on a JEOL 1200 EX transmission electron  
493 microscope.

494

#### 495 *Hydrogen/deuterium exchange mass spectrometry*

496 Peptide mass fingerprinting from purified Rubisco and carboxysomes was performed  
497 using on-column pepsin digestion, followed by reversed-phase HPLC, and tandem mass  
498 spectrometry on a Thermo Scientific LTQ Orbitrap Discovery.<sup>56,57</sup> For hydrogen exchange, the  
499 samples were diluted 1:10 in D<sub>2</sub>O buffer (50 mM Tris, 150 mM NaCl, pD 7.5) and then aliquots  
500 removed and quenched in 500 mM glycine, 2 M guanidinium hydrochloride (GdnHCl), pH 2.0  
501 buffer at log-spaced time intervals from 20 seconds to 48 hours. Samples were immediately  
502 frozen in liquid nitrogen upon addition of quenching solution. Deuterated control samples were  
503 prepared by 1:10 dilution in D<sub>2</sub>O, 50 mM Tris, 150 mM NaCl, 6 M GdnHCl, pD 7.5 and  
504 quenching with 500 mM glycine, pH 2.0. Samples were thawed, digested on-column as before,  
505 and analyzed by LCMS. Data analysis was performed with HDEaminer (Sierra Analytics).

506

#### 507 *CD spectroscopy*

508 Purified protein was first exchanged into CD buffer (20 mM sodium phosphate and 20  
509 mM sodium sulfate, pH 7.4) to minimize the background absorbance. From this solution, 300 µL  
510 was transferred to a 1-mm quartz cell. The sample containing only CD buffer was included as a  
511 negative control. Data were collected on a J-815 circular dichroism spectrometer (JASCO).  
512 Spectra were collected from 190 to 260 nm in 0.5 nm steps with the scanning speed of 20  
513 nm/min and signal averaging for 1 s for each step. Each sample was measured 3 times and the  
514 spectra were averaged. Protein concentrations were determined using 280 nm absorbance and  
515 extinction coefficients calculated using ProtParam.

516

#### 517 *Bioinformatics*

518 The CsoS2 secondary structure predictions were made using JPred.<sup>27</sup> The disorder  
519 score was calculated with PONDR-FIT.<sup>26</sup>

520 The candidate  $\alpha$ -carboxysome-associated CsoS2 sequences were selected from the  
521 Integrated Microbial Genomes (IMG) database by searching for the CsoS2 PFAM (PF12288)  
522 within 100kb of loci containing the Rubisco large and small subunits (PF00016 and PF00101),  
523  $\alpha$ -carboxysomal carbonic anhydrase (PF08936), and bacterial microcompartment shell proteins  
524 (PF00936). These sequences (n=231) were aligned with ClustalOmega,<sup>58</sup> truncated to include



525 only the NTD (i.e. all sequence before the first MR repeat), and analyzed with MEME<sup>31</sup> to find  
526 repeated sequence motifs (Fig. 2c). The Motif Alignment and Search Tool (MAST)<sup>39</sup> was used  
527 to locate and count all occurrences of the motif within the full CsoS2 sequences (Fig. 4b).

#### 528 **Acknowledgements:**

529 We thank Cecilia Blikstad for helpful comments on the manuscript. We also thank Peter Huang  
530 for his help with the BLI instrumentation and Cheryl Kerfeld for advice on Rubisco crystallization.  
531 Yinon Bar-On assisted us in gathering the CsoS2 sequences. We acknowledge the staff at the  
532 UC Berkeley Electron Microscope Laboratory for training and assistance with TEM. George  
533 Meigs and James Holton assisted with the x-ray diffraction and we gratefully acknowledge their  
534 input. Whiskers for crystal microseeding were kindly gifted by S.T. Kuhl. Beamline 8.3.1 at the  
535 Advanced Light Source is operated by the University of California Office of the President,  
536 Multicampus Research Programs and Initiatives grant MR-15-328599, the National Institutes of  
537 Health (R01 GM124149 and P30 GM124169), Plexikon Inc. and the Integrated Diffraction  
538 Analysis Technologies program of the US Department of Energy Office of Biological and  
539 Environmental Research. The work was supported by grants from the U.S. Department of  
540 Energy (DE-SC00016240) and the National Institute of General Medical Sciences  
541 (R01GM129241) to D.F.S. and a grant from the National Institute of General Medical Sciences  
542 (R01GM050945) to S.M.

543

544

#### 545 **Competing interests:**

546

547 D.F.S. is a co-founder of Scribe Therapeutics and a scientific advisory board member of Scribe  
548 Therapeutics and Mammoth Biosciences. All other authors declare no competing interests.

549

550 **References:**

- 551
- 552 1. Raven, J. A., Cockell, C. S. & De La Rocha, C. L. The evolution of inorganic carbon  
553 concentrating mechanisms in photosynthesis. *Philos. Trans. R. Soc. Lond. B. Biol. Sci* **363**,  
554 2641–2650 (2008).
- 555 2. Mangan, N. M., Flamholz, A., Hood, R. D., Milo, R. & Savage, D. F. pH determines the  
556 energetic efficiency of the cyanobacterial CO<sub>2</sub> concentrating mechanism. *Proc Natl Acad*  
557 *Sci USA* **113**, E5354-62 (2016).
- 558 3. Espie, G. S. & Kimber, M. S. Carboxysomes: cyanobacterial RubisCO comes in small  
559 packages. *Photosyn. Res.* **109**, 7–20 (2011).
- 560 4. Rae, B. D., Long, B. M., Badger, M. R. & Price, G. D. Functions, compositions, and  
561 evolution of the two types of carboxysomes: polyhedral microcompartments that facilitate  
562 CO<sub>2</sub> fixation in cyanobacteria and some proteobacteria. *Microbiol. Mol. Biol. Rev.* **77**, 357–  
563 379 (2013).
- 564 5. Heinhorst, S., Cannon, G. C. & Shively, J. M. in *Complex intracellular structures in*  
565 *prokaryotes* (ed. Shively, J. M.) **2**, 141–165 (Springer Berlin Heidelberg, 2006).
- 566 6. Kerfeld, C. A. & Melnicki, M. R. Assembly, function and evolution of cyanobacterial  
567 carboxysomes. *Curr. Opin. Plant Biol.* **31**, 66–75 (2016).
- 568 7. Tanaka, S. *et al.* Atomic-level models of the bacterial carboxysome shell. *Science* **319**,  
569 1083–1086 (2008).
- 570 8. Schmid, M. F. *et al.* Structure of *Halothiobacillus neapolitanus* carboxysomes by cryo-  
571 electron tomography. *J. Mol. Biol.* **364**, 526–535 (2006).
- 572 9. Iancu, C. V. *et al.* The structure of isolated *Synechococcus* strain WH8102 carboxysomes  
573 as revealed by electron cryotomography. *J. Mol. Biol.* **372**, 764–773 (2007).
- 574 10. Shih, P. M. *et al.* Biochemical characterization of predicted Precambrian RuBisCO. *Nat.*  
575 *Commun.* **7**, 10382 (2016).
- 576 11. Whitehead, L., Long, B. M., Price, G. D. & Badger, M. R. Comparing the in vivo function of  
577  $\alpha$ -carboxysomes and  $\beta$ -carboxysomes in two model cyanobacteria. *Plant Physiol.* **165**,  
578 398–411 (2014).
- 579 12. Shively, J. M., Ball, F., Brown, D. H. & Saunders, R. E. Functional organelles in  
580 prokaryotes: polyhedral inclusions (carboxysomes) of *Thiobacillus neapolitanus*. *Science*  
581 **182**, 584–586 (1973).
- 582 13. Cameron, J. C., Wilson, S. C., Bernstein, S. L. & Kerfeld, C. A. Biogenesis of a bacterial  
583 organelle: the carboxysome assembly pathway. *Cell* **155**, 1131–1140 (2013).
- 584 14. Kinney, J. N., Salmeen, A., Cai, F. & Kerfeld, C. A. Elucidating essential role of conserved  
585 carboxysomal protein CcmN reveals common feature of bacterial microcompartment  
586 assembly. *J. Biol. Chem.* **287**, 17729–17736 (2012).
- 587 15. Long, B. M., Badger, M. R., Whitney, S. M. & Price, G. D. Analysis of carboxysomes from  
588 *Synechococcus* PCC7942 reveals multiple Rubisco complexes with carboxysomal proteins  
589 CcmM and CcaA. *J. Biol. Chem.* **282**, 29323–29335 (2007).
- 590 16. Ryan, P. *et al.* The small RbcS-like domains of the  $\beta$ -carboxysome structural protein CcmM  
591 bind RubisCO at a site distinct from that binding the RbcS subunit. *J. Biol. Chem.* **294**,  
592 2593–2603 (2019).
- 593 17. Wang, H. *et al.* Rubisco condensate formation by CcmM in  $\beta$ -carboxysome biogenesis.  
594 *Nature* **566**, 131–135 (2019).
- 595 18. Long, B. M., Rae, B. D., Badger, M. R. & Price, G. D. Over-expression of the  $\beta$ -  
596 carboxysomal CcmM protein in *Synechococcus* PCC7942 reveals a tight co-regulation of  
597 carboxysomal carbonic anhydrase (CcaA) and M58 content. *Photosyn. Res.* **109**, 33–45  
598 (2011).
- 599 19. Cai, F. *et al.* Advances in understanding carboxysome assembly in prochlorococcus and

- 600 synechococcus implicate csos2 as a critical component. *Life (Basel)* **5**, 1141–1171 (2015).
- 601 20. Cannon, G. C. *et al.* Organization of carboxysome genes in the thiobacilli. *Curr. Microbiol.*  
602 **46**, 115–119 (2003).
- 603 21. Chaijarasphong, T. *et al.* Programmed Ribosomal Frameshifting Mediates Expression of  
604 the  $\alpha$ -Carboxysome. *J. Mol. Biol.* **428**, 153–164 (2016).
- 605 22. Williams, E. B. Identification and Characterization of Protein Interactions in the  
606 Carboxysome of *Halothiobacillus neapolitanus*. (2006).
- 607 23. Liu, Y. *et al.* Deciphering molecular details in the assembly of alpha-type carboxysome. *Sci.*  
608 *Rep.* **8**, 15062 (2018).
- 609 24. Gonzales, A. D. *et al.* Proteomic analysis of the CO<sub>2</sub>-concentrating mechanism in the  
610 open-ocean cyanobacterium *Synechococcus* WH8102. *Can. J. Bot.* **83**, 735–745 (2005).
- 611 25. Baker, S. H. *et al.* The correlation of the gene csos2 of the carboxysome operon with two  
612 polypeptides of the carboxysome in *thiobacillus neapolitanus*. *Arch. Microbiol.* **172**, 233–  
613 239 (1999).
- 614 26. Xue, B., Dunbrack, R. L., Williams, R. W., Dunker, A. K. & Uversky, V. N. PONDR-FIT: a  
615 meta-predictor of intrinsically disordered amino acids. *Biochim. Biophys. Acta* **1804**, 996–  
616 1010 (2010).
- 617 27. Drozdetskiy, A., Cole, C., Procter, J. & Barton, G. J. JPred4: a protein secondary structure  
618 prediction server. *Nucleic Acids Res.* **43**, W389–94 (2015).
- 619 28. Abdiche, Y., Malashock, D., Pinkerton, A. & Pons, J. Determining kinetics and affinities of  
620 protein interactions using a parallel real-time label-free biosensor, the Octet. *Anal.*  
621 *Biochem.* **377**, 209–217 (2008).
- 622 29. van der Lee, R. *et al.* Classification of intrinsically disordered regions and proteins. *Chem.*  
623 *Rev.* **114**, 6589–6631 (2014).
- 624 30. Davey, N. E. *et al.* Attributes of short linear motifs. *Mol. Biosyst.* **8**, 268–281 (2012).
- 625 31. Bailey, T. L. & Elkan, C. Fitting a mixture model by expectation maximization to discover  
626 motifs in biopolymers. *Proc. Int. Conf. Intell. Syst. Mol. Biol.* **2**, 28–36 (1994).
- 627 32. Alberty, R. A. *Thermodynamics of biochemical reactions*. (John Wiley & Sons, Inc., 2003).  
628 doi:10.1002/0471332607
- 629 33. Krissinel, E. & Henrick, K. Inference of macromolecular assemblies from crystalline state. *J.*  
630 *Mol. Biol.* **372**, 774–797 (2007).
- 631 34. Gallivan, J. P. & Dougherty, D. A. Cation- $\pi$  interactions in structural biology. *Proc Natl*  
632 *Acad Sci USA* **96**, 9459–9464 (1999).
- 633 35. Schneider, G., Lindqvist, Y. & Brändén, C. I. RUBISCO: structure and mechanism. *Annu.*  
634 *Rev. Biophys. Biomol. Struct.* **21**, 119–143 (1992).
- 635 36. Brangwynne, C. P., Tompa, P. & Pappu, R. V. Polymer physics of intracellular phase  
636 transitions. *Nat. Phys.* **11**, 899–904 (2015).
- 637 37. Nott, T. J. *et al.* Phase transition of a disordered nuage protein generates environmentally  
638 responsive membraneless organelles. *Mol. Cell* **57**, 936–947 (2015).
- 639 38. Qamar, S. *et al.* FUS Phase Separation Is Modulated by a Molecular Chaperone and  
640 Methylation of Arginine Cation- $\pi$  Interactions. *Cell* **173**, 720–734.e15 (2018).
- 641 39. Bailey, T. L. & Gribskov, M. Combining evidence using p-values: application to sequence  
642 homology searches. *Bioinformatics* **14**, 48–54 (1998).
- 643 40. Bonacci, W. *et al.* Modularity of a carbon-fixing protein organelle. *Proc Natl Acad Sci USA*  
644 **109**, 478–483 (2012).
- 645 41. Li, P. *et al.* Phase transitions in the assembly of multivalent signalling proteins. *Nature* **483**,  
646 336–340 (2012).
- 647 42. Boeynaems, S. *et al.* Protein phase separation: A new phase in cell biology. *Trends Cell*  
648 *Biol.* **28**, 420–435 (2018).
- 649 43. Mackinder, L. C. M. *et al.* A repeat protein links Rubisco to form the eukaryotic carbon-  
650 concentrating organelle. *Proc Natl Acad Sci USA* **113**, 5958–5963 (2016).

- 651 44. Wunder, T., Cheng, S. L. H., Lai, S.-K., Li, H.-Y. & Mueller-Cajar, O. The phase separation  
652 underlying the pyrenoid-based microalgal Rubisco supercharger. *Nat. Commun.* **9**, 5076  
653 (2018).
- 654 45. Freeman Rosenzweig, E. S. *et al.* The Eukaryotic CO<sub>2</sub>-Concentrating Organelle Is Liquid-  
655 like and Exhibits Dynamic Reorganization. *Cell* **171**, 148-162.e19 (2017).
- 656 46. Long, B. M., Tucker, L., Badger, M. R. & Price, G. D. Functional cyanobacterial beta-  
657 carboxysomes have an absolute requirement for both long and short forms of the CcmM  
658 protein. *Plant Physiol.* **153**, 285–293 (2010).
- 659 47. Hyman, A. A., Weber, C. A. & Jülicher, F. Liquid-liquid phase separation in biology. *Annu.*  
660 *Rev. Cell Dev. Biol.* **30**, 39–58 (2014).
- 661 48. Engler, C., Kandzia, R. & Marillonnet, S. A one pot, one step, precision cloning method with  
662 high throughput capability. *PLoS ONE* **3**, e3647 (2008).
- 663 49. Schuler, B., Lipman, E. A., Steinbach, P. J., Kumke, M. & Eaton, W. A. Polyproline and the  
664 “spectroscopic ruler” revisited with single-molecule fluorescence. *Proc Natl Acad Sci USA*  
665 **102**, 2754–2759 (2005).
- 666 50. Kabsch, W. Integration, scaling, space-group assignment and post-refinement. *Acta*  
667 *Crystallogr. D Biol. Crystallogr.* **66**, 133–144 (2010).
- 668 51. Collaborative Computational Project, Number 4. The CCP4 suite: programs for protein  
669 crystallography. *Acta Crystallogr. D Biol. Crystallogr.* **50**, 760–763 (1994).
- 670 52. Evans, P. R. & Murshudov, G. N. How good are my data and what is the resolution? *Acta*  
671 *Crystallogr. D Biol. Crystallogr.* **69**, 1204–1214 (2013).
- 672 53. McCoy, A. J. *et al.* Phaser crystallographic software. *J. Appl. Crystallogr.* **40**, 658–674  
673 (2007).
- 674 54. Adams, P. D. *et al.* PHENIX: a comprehensive Python-based system for macromolecular  
675 structure solution. *Acta Crystallogr. D Biol. Crystallogr.* **66**, 213–221 (2010).
- 676 55. Emsley, P. & Cowtan, K. Coot: model-building tools for molecular graphics. *Acta*  
677 *Crystallogr. D Biol. Crystallogr.* **60**, 2126–2132 (2004).
- 678 56. Lim, S. A., Bolin, E. R. & Marqusee, S. Tracing a protein’s folding pathway over  
679 evolutionary time using ancestral sequence reconstruction and hydrogen exchange. *elife* **7**,  
680 (2018).
- 681 57. Samelson, A. J. *et al.* Kinetic and structural comparison of a protein’s cotranslational folding  
682 and refolding pathways. *Sci. Adv.* **4**, eaas9098 (2018).
- 683 58. Sievers, F. & Higgins, D. G. Clustal Omega for making accurate alignments of many protein  
684 sequences. *Protein Sci.* **27**, 135–145 (2018).