

1	<b>Supplementary Information:</b>	
2		
3		
4	<i>Comments on bio-layer interferometry (BLI)</i>	2
5	<i>BLI of select CsoS2 peptides with Rubisco</i>	4
6	<i>BLI of mutant NTD</i>	5
7	<i>Rubisco binding site mutants</i>	6
8	<i>Microscale thermophoresis (MST) fitting and analysis</i>	7
9	<i>Rubisco-N* fusion characterization</i>	9
10	<i>X-ray crystallography refinement statistics</i>	10
11	<i>Enrichment of binding motif from peptide array data</i>	11
12	<i>Hydrogen / deuterium exchange of Rubisco inside and outside carboxysomes</i>	13
13	<i>Structural comparison to CcmM/Rubisco interaction</i>	14
14	<i>CsoS2 NTD truncations and carboxysome formation</i>	18
15	<i>Protein sequences</i>	19
16	<i>Plasmids used</i>	25
17		
18		

19

20 *Comments on bio-layer interferometry (BLI)*

21

22

23

24

25

26

27

Bio-layer interferometry (BLI) played a critical role in this study as a robust, medium-throughput way to narrow down the CsoS2/Rubisco binding activity to the N-peptide and as a qualitative test of the effects of binding site mutations on the activity. In principle this method can yield specific information on binding energies. However, in this case, the high valency of the interaction (8 for Rubisco and 4 for CsoS2) combined with surface avidity effects make this kind of energetic determination infeasible.

28

29

30

31

32

For any surface-based binding measurement, false positives due to non-specific binding are a concern. Our confidence in the qualitative binding results is born of two observations: one, there is very little signal accumulation on unloaded biosensors (e.g. Fig. S2, black trace) and, two, minor targeted modifications to the bait or the prey could entirely abolish binding (e.g. randomizing the N\*-peptide sequence eliminated all Rubisco binding; see Fig. 2b).

33

34

35

36

37

38

39

40

41

42

43

44

The BLI data for every bait construct with binding activity towards Rubisco (i.e. full length CsoS2, NTD, N\*-polyPro) had clear signatures of surface avidity. This situation is essentially unavoidable because Rubisco has a valency of 8 and the BLI biosensors require high surface densities relative to Rubisco's size. In the case of CsoS2 and NTD, they both contain four N-peptides and the avidity effects are particularly acute as evidenced by the heterogeneous dissociation kinetics (Fig. 2a) which likely arise from density variations of bait and result in microscopic surface sites of differing affinities. N\*-polyPro, with just a single N-peptide, demonstrates simpler fully reversible binding but is nevertheless still poorly modeled as a 1:1 binding interaction (Fig. S1). The implied dissociation constant is around 100 nM but it is not clear, for example, how many individual interactions are involved. In any event, we make no attempt to quantify the binding energetics from these results but rather took them as qualitative indications of binding.

45

46

The data in Fig. S1 was globally fit with the following piece-function:

47

$$\text{Eq. S1} \quad 0 < t \leq t_{dissoc}, \quad S(t, c_{Ru}) = S_{sat} * \left( \frac{1 - e^{-(k_a c_{Ru} + k_d) t}}{1 + \frac{k_d}{k_a c_{Ru}}} \right)$$

48

49

50

$$t_{dissoc} < t, \quad S(t, c_{Ru}) = S(t_{dissoc}, c_{Ru}) * e^{-k_d (t - t_{dissoc})}$$

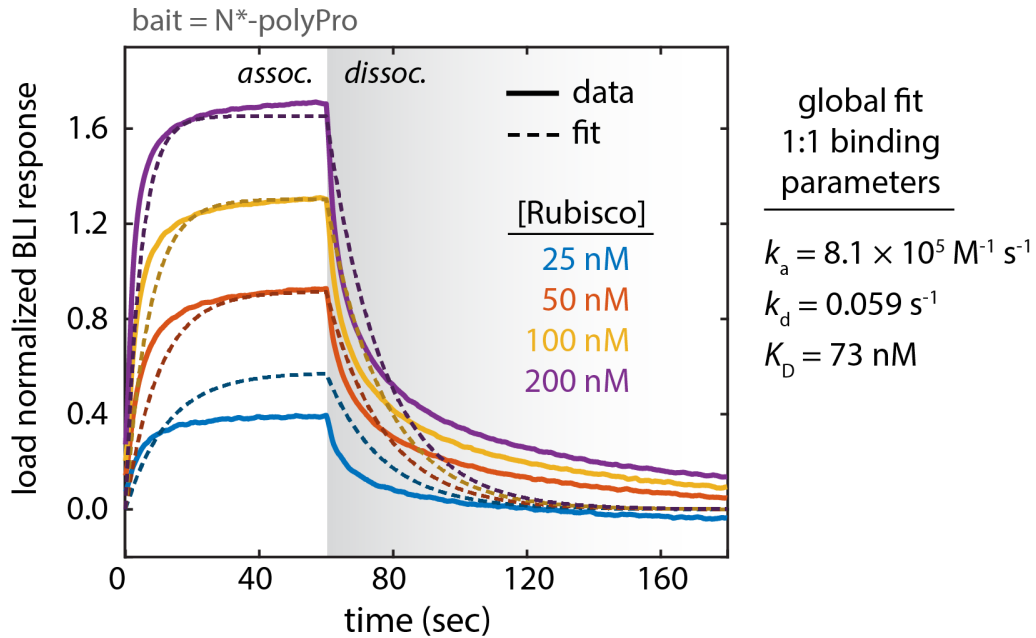
51

52

53

54

where  $t$  is time,  $t_{dissoc}$  is the time of biosensor transfer to dissociation buffer,  $c_{Ru}$  is the Rubisco concentration,  $k_a$  is the bimolecular kinetic association rate constant,  $k_d$  is the unimolecular dissociation rate constant, and  $S_{sat}$  is the BLI signal at saturation.

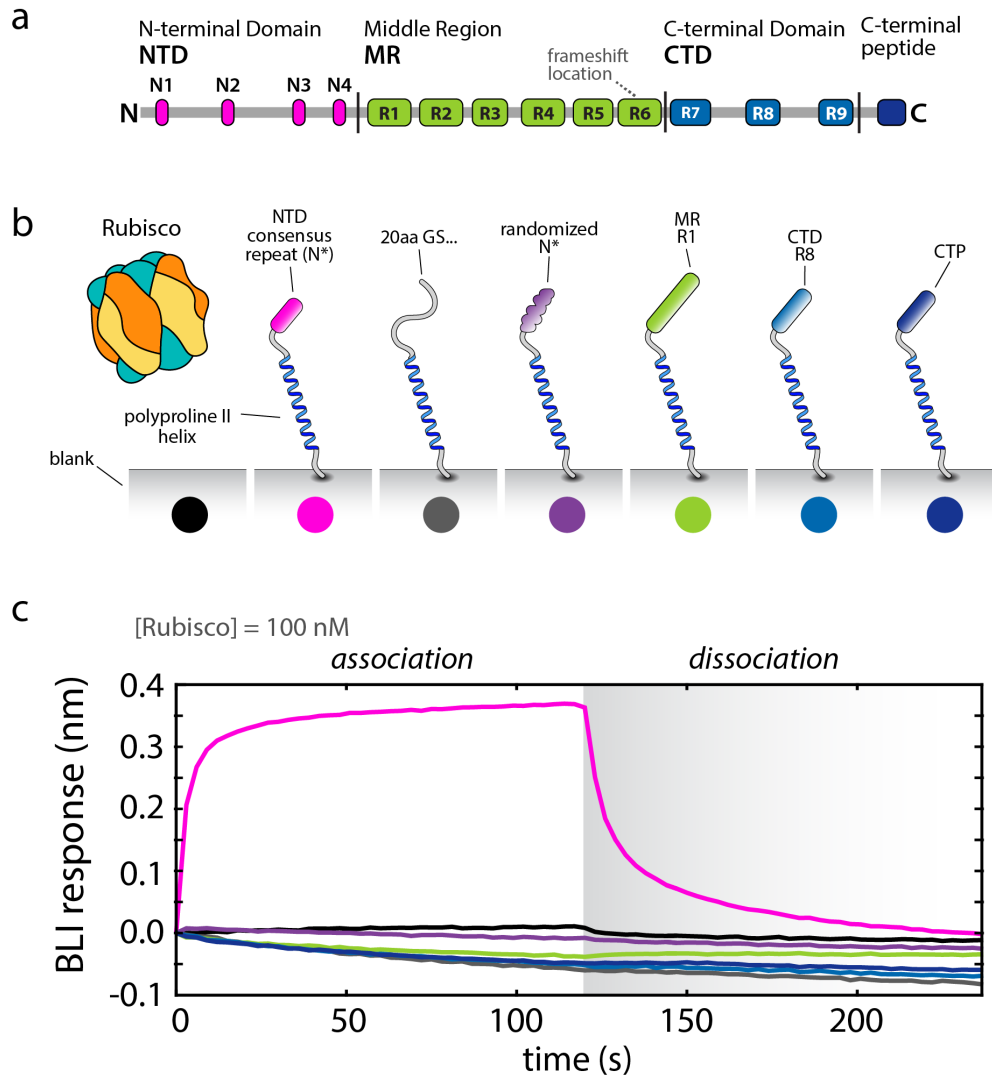


55  
56  
57

**Figure S1**

58 BLI response curves for Rubisco binding to N\*-polyPro. The 1:1 binding model global fit to Eq. S1 is  
59 shown. Due to the high valency of Rubisco, this model only applies in the limit of very low surface density  
60 of the monovalent bait (N\*-polyPro). The binding profiles exhibit deviations from this idealization which  
61 are consistent with significant surface avidity effects. It therefore does not lend itself to simple  
62 deconvolution of the energetics.

63 *BLI of select CsoS2 peptides with Rubisco*



64  
65

66 **Figure S2**

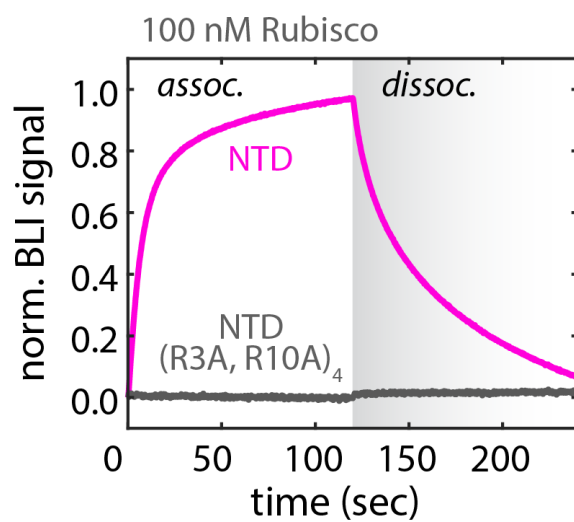
67  
68

69 **a**, Primary sequence of CsoS2 with each of the repeated or conserved elements. **b**, Schematic  
70 representation of a set of BLI experiments testing the specificity of the Rubisco - CsoS2 interaction. Each  
71 of a series of CsoS2 elements and control sequences was fused to polyproline II helices which were  
72 surface immobilized to a Ni-NTA functionalized biosensor surface via an N-terminal hexahistidine tag. **c**,  
73 BLI traces of the constructs from **(b)** when incubated with 100 nM Rubisco. The trace colors match the  
74 dots in **(b)**. Only the N\*-peptide demonstrates any specific binding activity.

75 *BLI of mutant NTD*

76

77



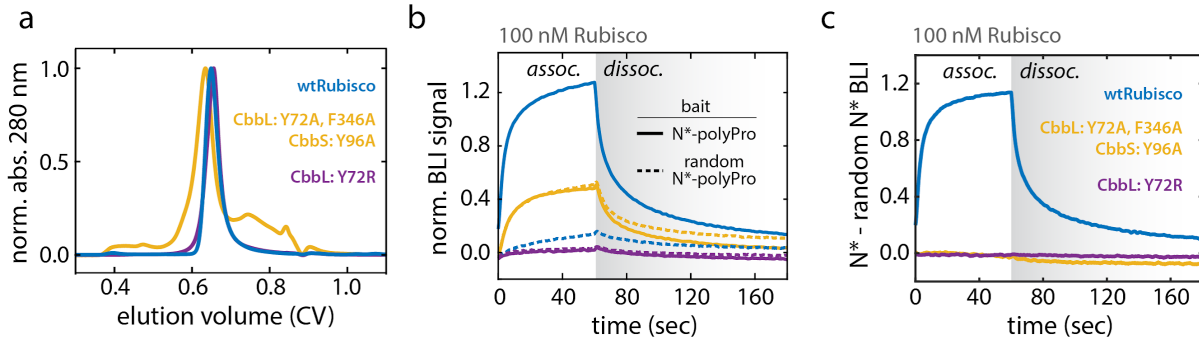
78

79 **Figure S3**

80 BLI response towards 100 nM Rubisco with bait of either the NTD or the NTD with R3A, R10A mutations  
81 made within all four of the N-peptide repeats. Removing those conserved arginines entirely eliminates the  
82 binding.

83 Rubisco binding site mutants

84



85

86

87 **Figure S4**

88

89 **a**, Size exclusion chromatograms of wild-type *H. neapolitanus* Rubisco (wtRubisco), a mutant with all  
90 cation- $\pi$  aromatics mutated to alanines (CbbL: Y72A, F346A; CbbS: Y96A), and a salt bridge disrupting  
91 mutation (CbbL: Y72A). All species eluted at a volume consistent with the  $L_8S_8$  structure. **b**, Each Rubisco  
92 species was tested for binding activity by BLI to N\*-polyPro (solid lines) and the randomized N\*-polyPro  
93 negative control (dashed lines). Only the wild-type Rubisco had specific binding activity to N\*-polyPro  
94 over the randomized N\*-peptide control. The aromatic removal mutant (yellow) had some non-specific  
95 binding to both baits but showed no preference for the real N\*-peptide sequence. **c**, Differential BLI  
96 binding signal of each Rubisco species to N\*-polyPro relative to random N\*-polyPro. Both Rubisco  
97 binding site mutants clearly possess no specific association.

99

100 MST data were, in general, collected for a series of 16 ligand concentrations as serial 2-  
 101 fold dilutions. The target concentrations were 50 nM of the GFP fusion species. The isotherms  
 102 were fit to a 1:1 binding model according to the law of mass action,

103

104

105 Eq. S2

106

107

$$S(c_L) = S_{unbound} + (S_{bound} - S_{unbound}) \frac{c_L - c_T + K_D - \sqrt{(c_L + c_T + K_D)^2 - 4 c_L c_T}}{2 c_T}$$

108

109 where  $S$  is the observed MST response,  $S_{unbound}$  and  $S_{bound}$  are the MST responses of the free  
 110 and saturated target respectively,  $c_L$  is the ligand concentration (varied over 16 concentrations),  
 111  $c_T$  is the target concentration (constant for each experiment), and  $K_D$  is the dissociation  
 112 constant.

113

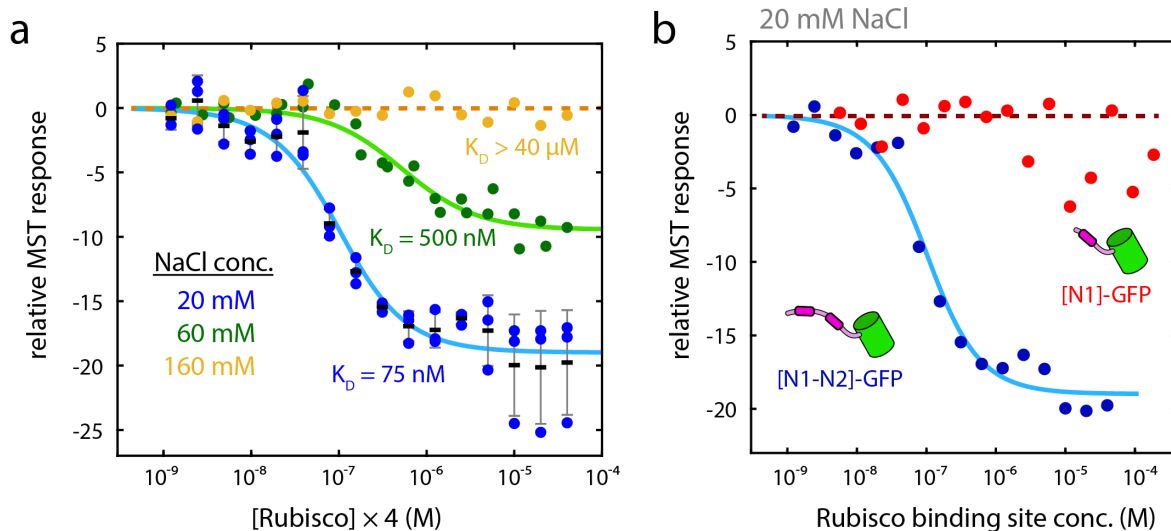
114 The ligand concentration is taken as the effective total concentration of binding sites  
 115 available to the target. For example, [N1-N2]-GFP can engage two of the eight sites on Rubisco,  
 116 therefore, the ligand concentration is four times the Rubisco holoenzyme ( $L_8S_8$ ) concentration.  
 117 Implicit in this simple treatment are a number of assumptions. One, the targets bind with all N-  
 118 peptides (e.g. 2 for [N1-N2]-GFP). Two, all possible binding configurations (that is, the  
 119 microscopic arrangements of the N-peptide binding locations allowable by the linking region)  
 120 have equivalent binding energies. And three, the thermophoretic response of the bound species  
 121 will be the same regardless of the number of targets bound to that particular Rubisco.  
 122 Deviations from these assumptions are expected to be small and do not justify the inclusion of  
 123 additional fit parameters.

123

124 Where applicable the fitting procedure was conducted by taking the median fit  
 125 parameters from bootstrap sampling. Specifically, a subset of the data was randomly selected  
 126 from among the experimental replicates. This subset was fit to Eq. S2 by least-squares fitting  
 127 and the parameters recorded. The process was repeated 10,000 times. The mean and  
 128 confidence intervals for  $K_D$  were determined from the resulting distribution.

128

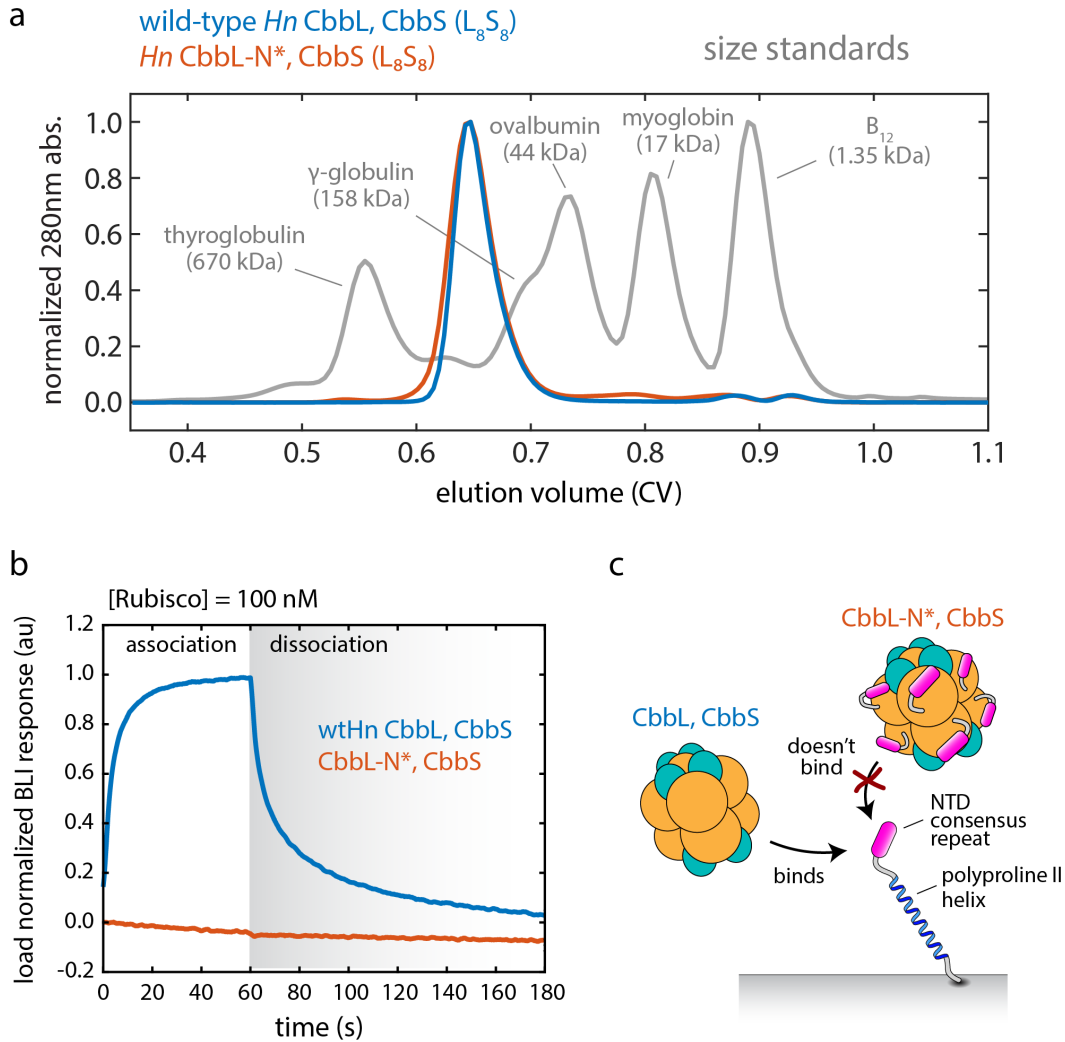
129 The single site (i.e. one N-peptide to one Rubisco site) binding energy proved to be too  
 130 weak to accurately determine under our experimental conditions (see Fig. S5b). Consequently,  
 131 we made rough estimates of the single site binding constants by assuming a linear free energy  
 132 relationship, that is, the free energy of binding for a single site was taken as half that of the  
 133 bivalent species ([N1-N2]-GFP; see Fig. 2e dashed lines). This should be construed as an  
 134 upper limit since multivalent ligands are empirically known to bind more weakly than the sum of  
 135 the individual site free energies. This shortfall is generally attributed to the entropy decreases in  
 the linker regions partially offsetting the favorable binding energies.<sup>1,2</sup>



136  
137 **Figure S5**  
138

139 **a**, MST responses for [N1-N2]-GFP association to Rubisco. The concentration of the target, [N1-N2]-GFP,  
140 was 50 nM. The abscissa represents the concentration of effective binding sites and is four times the  
141 Rubisco  $L_8S_8$  concentration since each target will engage two of the eight possible sites. Binding  
142 experiments were performed at 20, 60, and 160 mM NaCl. At 20 mM NaCl three replicates were  
143 performed across 16 Rubisco concentrations. Black lines indicate the means while the gray whiskers  
144 show  $\pm$  one standard deviation. At 60 mM NaCl the experiment was performed twice with slightly varying  
145 concentrations. At 160 mM NaCl data from one representative experiment is shown. The fits to the 20 mM  
146 and 60 mM NaCl data are according to Eq. S2 and represent the mean fit parameters from bootstrap  
147 sampling of the data. For 160 mM NaCl no binding could be determined over this concentration range  
148 and the dashed orange line is drawn at zero response as a visual guide. **b**, Comparison between a  
149 double N-peptide, [N1-N2]-GFP, and single N-peptide, [N1]-GFP, species by MST. Both had 50 nM  
150 target. The Rubisco binding site concentration is specific to the two different targets. For [N1-N2]-GFP it is  
151 the concentration of  $L_8S_8$  multiplied by 4 and for [N1]-GFP it is the concentration of  $L_8S_8$  multiplied by 8  
152 since the former has four potential binding sites on the Rubisco holoenzyme while the latter has eight.  
153 The [N1]-GFP data points are the mean values values from (a). The [N1]-GFP data points are from one  
154 representative experiment and indicates no conclusive binding over the concentration range. The dashed  
155 red line is at zero response as a visual guide.





157  
158 **Figure S6**

159  
160 **a**, Size exclusion analysis of wild-type Rubisco and the N\*-peptide fusion construct. Both elute at volumes  
161 commensurate with compact L<sub>8</sub>S<sub>8</sub> complexes. A run with the Bio-Rad Gel Filtration Standard is included  
162 for comparison. Standard masses are indicated. **b**, BLI responses of wtRubisco and the N\*-peptide fusion  
163 Rubisco at 100 nM with N\*-polyPro as the surface bait. The fusion showed no binding. **c**, Proposed  
164 cartoon model of differential BLI binding activities. N\*-peptide fusion Rubisco is apparently self-passivated  
165 by saturating the binding sites from stable association of the fused N\*-peptides.

166 *X-ray crystallography refinement statistics*

167

168 Table S1. Data collection and refinement statistics (molecular replacement)

169

<i>H. neapolitanus</i> CbbL-N*, CbbS	
<b>Data collection</b>	
Space group	C 2
Cell dimensions	
<i>a</i> , <i>b</i> , <i>c</i> (Å)	171.83 153.95 108.06
$\alpha$ , $\beta$ , $\gamma$ (°)	90 124.70 90
Resolution (Å)	104.1 - 2.4 (2.486 - 2.4) <sup>a</sup>
$R_{\text{sym}}$ or $R_{\text{merge}}$	0.1244 (0.5876)
$I / \sigma I$	12.84 (2.97)
Completeness (%)	99.90 (99.92)
Redundancy	6.9 (6.4)
<b>Refinement</b>	
Resolution (Å)	104.1 - 2.4 (2.486 - 2.4)
No. reflections	89958 (8958)
$R_{\text{work}} / R_{\text{free}}$	0.184 (0.235) / 0.248 (0.301)
No. atoms	18568
Protein	17618
Ligand/ion	0
Water	950
<i>B</i> -factors	41.56
Protein	41.65
Ligand/ion	n/a
Water	40.04
R.m.s. deviations	
Bond lengths (Å)	0.008
Bond angles (°)	1.19

170 Values in parentheses are for highest-resolution shell.

171

172 a) data from a single crystal

173 *Enrichment of binding motif from peptide array data*

174

175 Cai et al. performed an experiment testing the binding of Rubisco to a peptide array  
176 chip.<sup>3</sup> The peptide library was composed of every 8-mer CsoS2 peptide (from *Prochlorococcus*  
177 *marinus* MIT9313) tiled residue-by-residue across the entire protein. The chip was incubated  
178 with Rubisco, washed, and then assayed with fluorescently labeled anti-RbcL antibody. The  
179 relative fluorescence at each site then provides some indication of the Rubisco binding activity.  
180 Cai et al. observed a large number of potential hits scattered throughout the CsoS2 sequence  
181 making a determination of the interaction motif challenging. The raw data was generously  
182 provided by the authors in the Supplementary Material. This data set was re-analyzed in light of  
183 the biochemical and structural evidence of the binding motif presented herein.

184 Since the original data did not have a clear indication of a specific binding site, we chose  
185 to look at it in a statistical manner. R3, R9, and R10 are key conserved N\*-peptide residues  
186 containable within an 8-mer (G17 and K18 are too far separated). We examined those peptides  
187 containing at least two basic residues (i.e. K or R) as generically positively charged species and  
188 ones matching the particular R spacings consistent with any pair of arginines among R3, R9,  
189 and R10 as a test of the specific motif (i.e. RxxxxxR, RxxxxxxR, or RR).

190 The results showed that the doubly basic peptides demonstrate enriched fluorescence  
191 signal relative to the full peptide library. Peptides matching the motif regular expressions above,  
192 however, had significantly greater enrichment even than the doubly basic subset (Fig. S7a). A  
193 bootstrap analysis was performed to assess the likelihood of obtaining an equivalently high  
194 median fluorescence via random peptides from either the whole population or the doubly basic  
195 subset (Fig. S7b). Out of 10,000 trials none were found to exceed the motif median  
196 fluorescence, implying a p-value less than  $10^{-4}$ .

197 This outcome indicates that the peptide array binding propensities are indeed consistent  
198 with the binding motif identified in this work. The relatively high abundance of 8-mers in the  
199 peptide library containing portions of the binding motif resulted in multiple “hotspots” scattered  
200 throughout the CsoS2 sequence and made prospective unique identification of the specific  
201 binding site(s) impossible. The statistical strength of the retrospective motif analysis attests the  
202 genuine signal in the peptide array experiment.

203

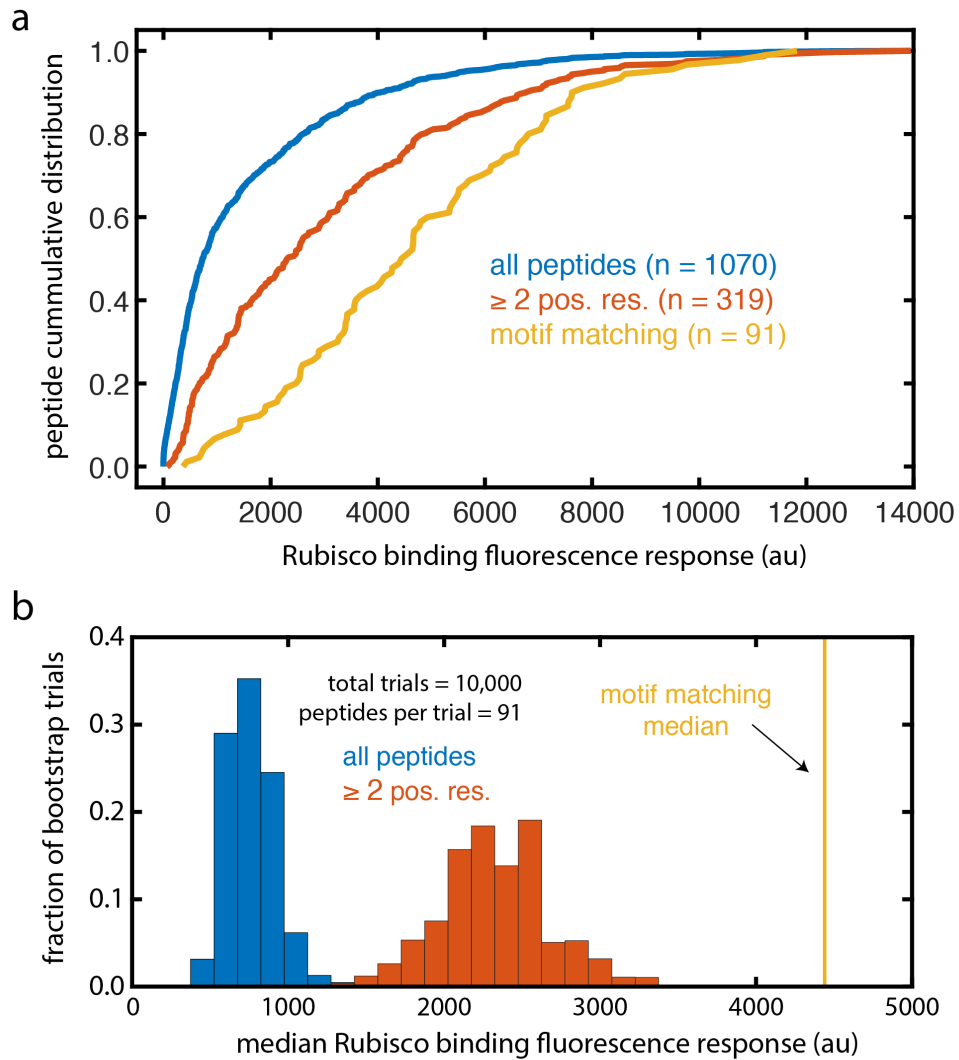
204

205

206

207

208



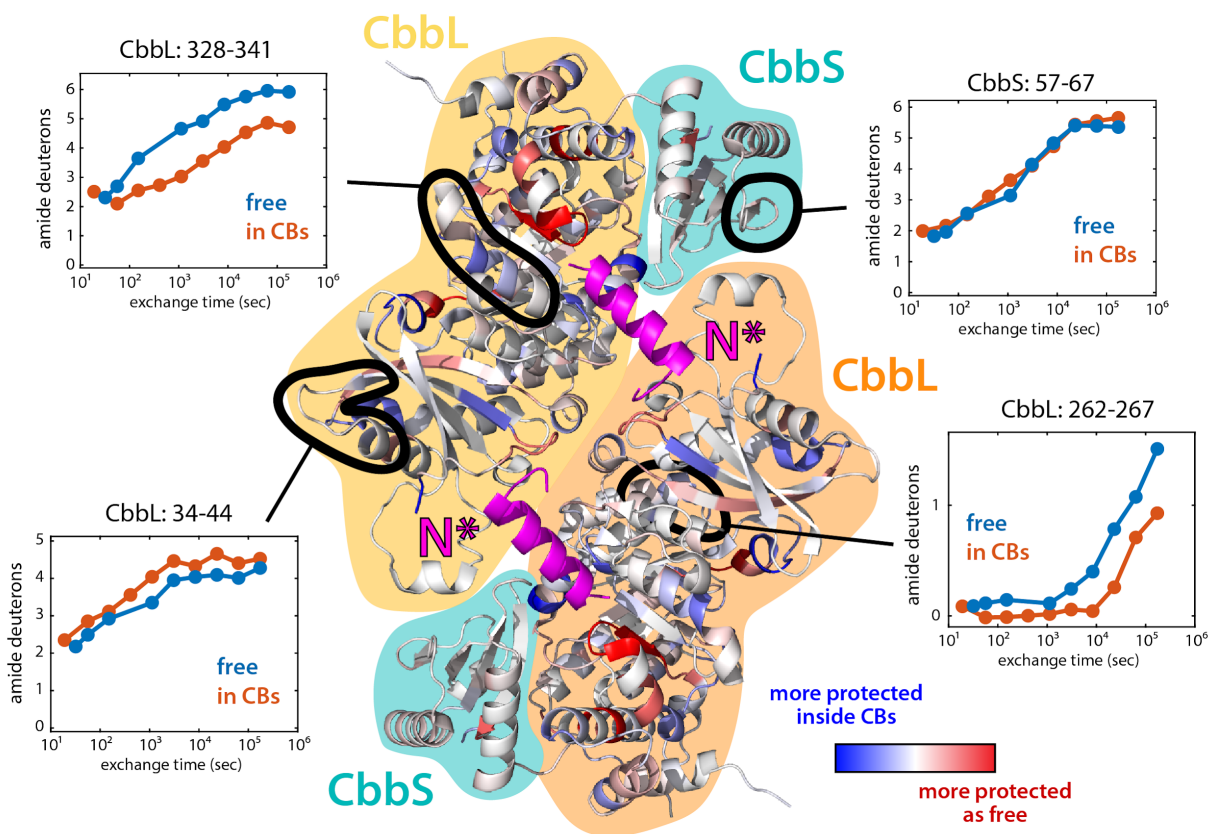
209  
 210  
 211

**Figure S7**

212 **a**, Cumulative distributions of Rubisco binding fluorescence response for CsoS2 array peptides including  
 213 the full dataset, those with more than two basic residues, and those matching the N\*-peptide arginine  
 214 motif. **b**, Distributions of bootstrap results. 91 peptides were taken at random (with replacement) from  
 215 either the full dataset or those with two or more basic residues and the median fluorescence response  
 216 calculated. 10,000 trials were conducted with each set and none exceeded the motif matching median.

217 *Hydrogen / deuterium exchange of Rubisco inside and outside carboxysomes*

218 Hydrogen / deuterium exchange (HDX) mass spectrometry experiments were performed  
219 on purified Rubisco and Rubisco encapsulated within purified carboxysomes. Overall, the  
220 differences between these two states was very minor with no particular regions possessing  
221 systematic differential protection (Fig. S8).

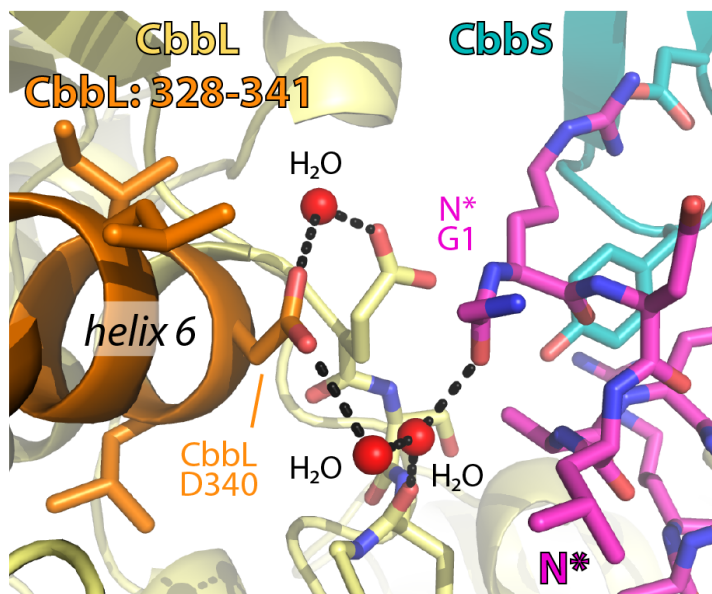


222  
223 **Figure S8**

224  
225 The structure displayed contains two CbbLs and two CbbS and shows the L<sub>2</sub> dimer interface across  
226 which the N\*-peptide (in magenta) binds. The Rubisco cartoon is colored according to the differential  
227 protection to amide hydrogen exchange. Those residues in blue experience greater protection within  
228 purified carboxysomes and those in red experience greater protection as free Rubisco. The comparison  
229 between these states was carried out with HDExaminer (Sierra Analytics) using moderate smoothing.  
230 Four specific peptides outlined in black highlight some of the diversity of HDX behavior. Most peptides  
231 that were observed from both states had essentially identical exchange kinetics as exemplified by the top  
232 right subpanel for CbbS: 57-67. Less common were peptides with different exchange profiles between  
233 encapsulated and unencapsulated Rubisco. CbbL: 34-44 (lower left subpanel) had slightly more  
234 protection in free Rubisco. CbbL: 328-341 (upper left subpanel) and CbbL: 262-267 (lower right subpanel)  
235 both had greater protection inside carboxysomes. Since it has the most dramatic results and comes  
236 closest to the N\*-peptide, the interactions of CbbL: 328-341 are examined in greater detail in Fig. S9  
237 below.  
238

239 Several specific peptides present in both samples had distinctive exchange profiles. Of  
240 particular note is CbbL: 328-341 which showed greater protection inside carboxysomes and is  
241 found in relatively close proximity to the N\*-peptide (Fig. S8). N\* contacts CbbL D340 through a  
242 water-mediated hydrogen bond network extending to N\* G1 (Fig. S9). The Rubisco-N\* fusion  
243 protein used for the crystal structure has a two serine linker joining N\* to the CbbL C-terminus  
244 which is not observed to be ordered. It is possible that residues upstream of the N-peptide  
245 binding motifs of CsoS2 interact more extensively with CbbL: 328-341 in the carboxysome.  
246 Finally, CbbL: 328-341 covers much of helix 6 and part of loop 6 which plays a role in Rubisco's  
247 activity and CO<sub>2</sub>/O<sub>2</sub> specificity.<sup>4</sup> Consequently, it is conceivable that CsoS2 binding leads to  
248 changes in Rubisco's catalytic properties tailored to the unique chemical environment of  
249 carboxysome lumen.

250  
251



252  
253 **Figure S9**

254  
255 Structure of the interactions bridging the N\*-peptide (magenta) and CbbL: 328-341 (orange), a peptide  
256 demonstrating significant HDX differences between carboxysomal and free Rubisco. A water-mediated  
257 hydrogen bond network extends between CbbL D340 and N\* G1 perhaps accounting in part for the  
258 peptide's greater HDX protection in carboxysomes.

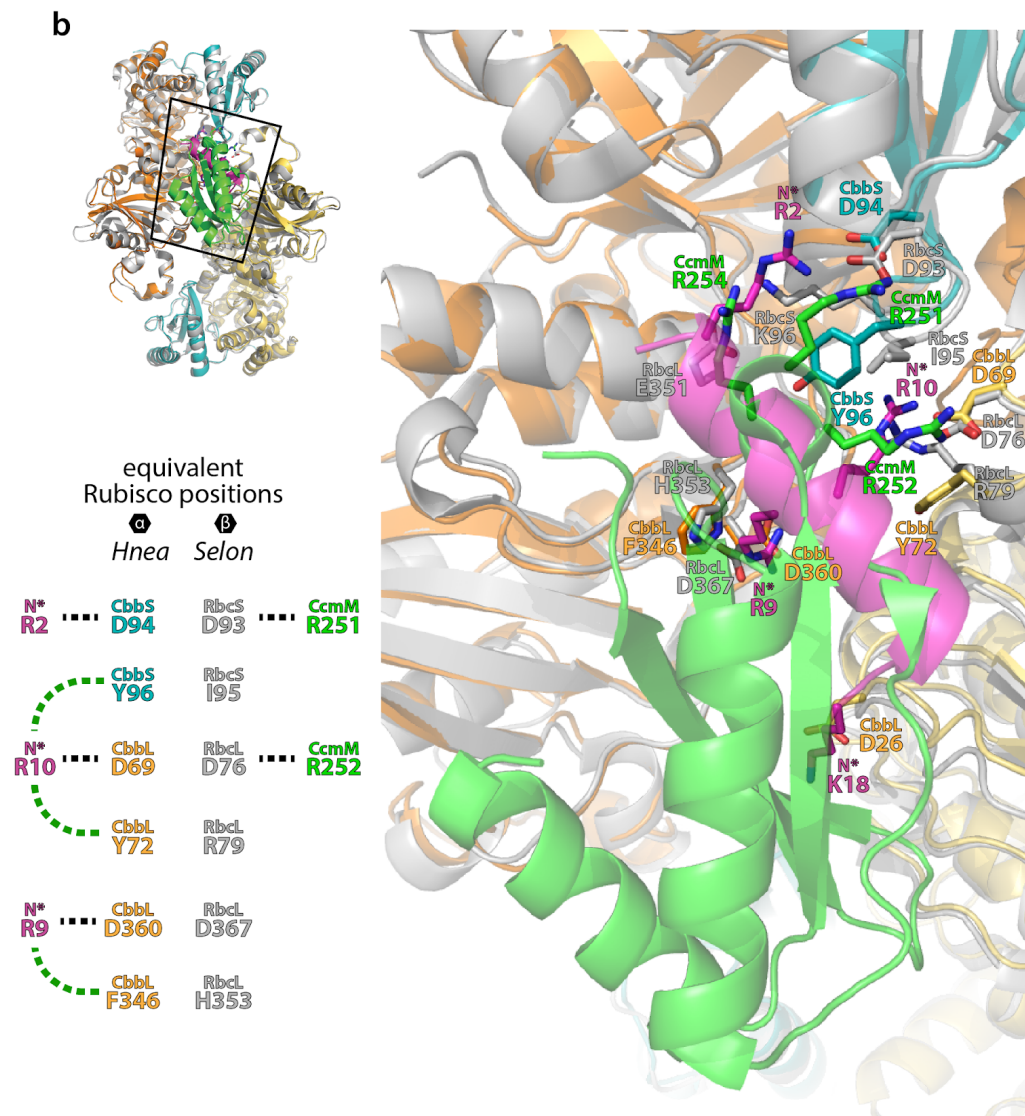
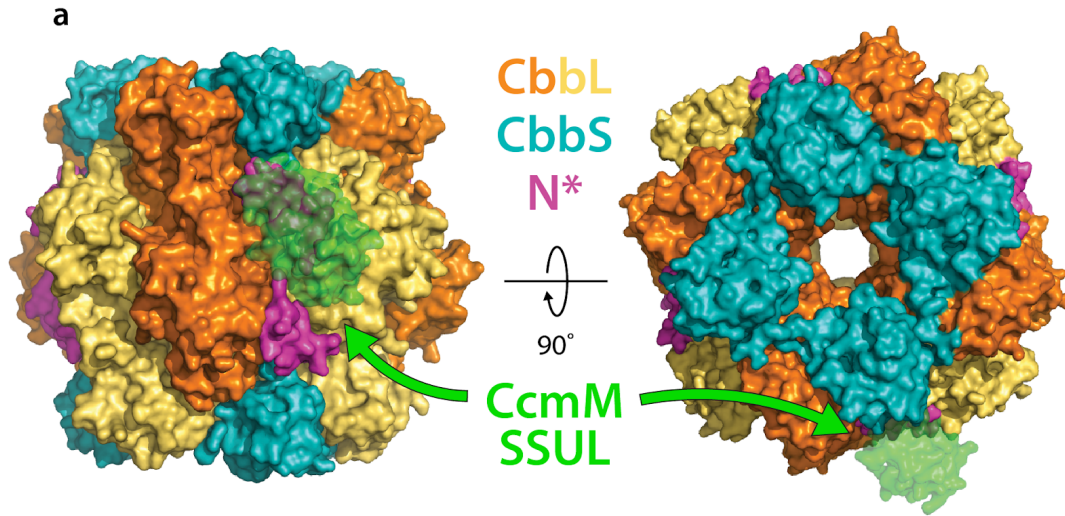
259  
260

261 *Structural comparison to CcmM/Rubisco interaction*

262 The recent work by Wang et al.<sup>5</sup> on the structure of the CcmM/Rubisco complex and  
263 attendant liquid-liquid phase separation affords a direct comparison between that scaffold  
264 interaction underlying the  $\beta$ -carboxysome assembly and the CsoS2/Rubisco interaction,  
265 described herein, central to  $\alpha$ -carboxysome assembly. Striking parallels are evident but the  
266 molecular details are distinct and bear no obvious evolutionary connection. That both systems  
267 converged upon multivalent binding to nearly identical Rubisco sites and have propensities  
268 toward phase separation, is a fascinating coincidence and perhaps a hint at some optimality of  
269 this assembly strategy.

270 In both cases the scaffolding element binds at the union of two L<sub>2</sub> dimers and a small  
271 subunit (Fig. S10a). Consequently, the binding site only exists in the fully assembled L<sub>8</sub>S<sub>8</sub>  
272 Rubisco holoenzyme. Unlike the N\*-peptide which can apparently simultaneously bind at eight  
273 possible sites, CcmM-SSUL occludes the immediately adjacent site and therefore has only four  
274 possible sites per Rubisco.

275 Wang et al. identified two primary regions of interaction between CcmM-SSUL and  
276 Rubisco which they called "Interface I" and "Interface II". Interface I closely overlaps with the  
277 N\*/Rubisco binding and is the area of focus below in Fig. S10b. (Interface II is located near the  
278 loop at the bottom right and is not shown in detail.) Interface I is largely electrostatic in nature  
279 with a series of three arginines (CcmM R251, R252, and R254) making important contacts.  
280 R251 and R252 form salt bridges to two aspartates (RbcS D93 and RbcL D76) which are  
281 positionally equivalent to those engaged in salt bridges to N\* R2 and N\* R10, respectively (see  
282 Fig. S10b table). Despite utilizing some of the same residues for salt bridges, the scaffold  
283 geometries are remarkably different. In CcmM, residues R251, R252, and R254 fan out from a  
284 short helix, called  $\alpha$ 2, whose axis is directed down inward to Rubisco. The N\*-peptide's helix  
285 axis, in contrast, runs perpendicular to the Rubisco surface. Finally, CcmM/Rubisco contains no  
286 apparent cation- $\pi$  interactions which feature prominently in the N\*/Rubisco binding interface.





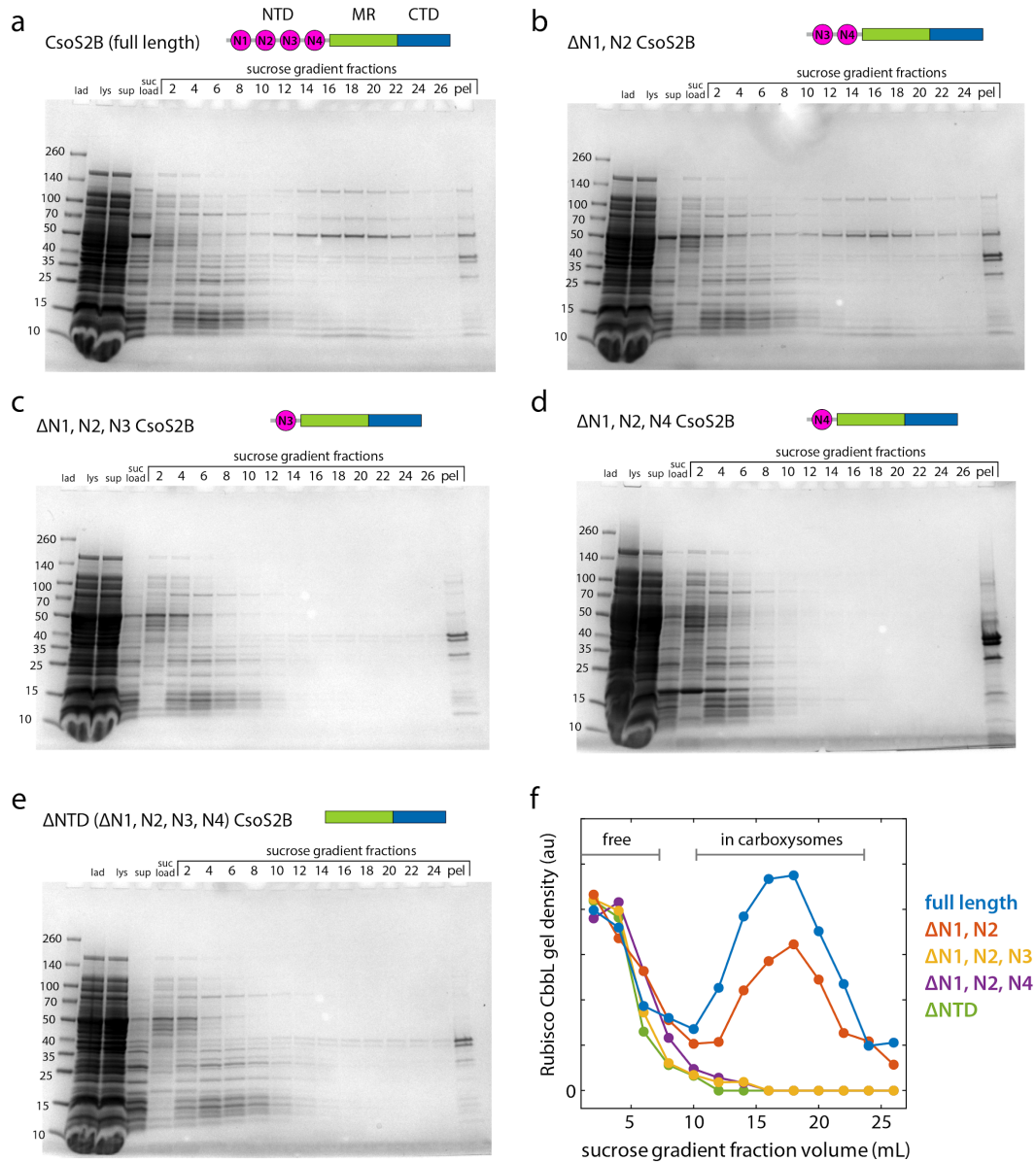
288 **Figure S10**

289

290 **a**, Surface representation of the N\*/Rubisco complex with aligned CcmM-SSUL from the model of Wang  
291 et al.<sup>5</sup> in semi-transparent green. **b**, Detailed comparative view of the scaffold/Rubisco interaction  
292 interface. The inset table pairs equivalent Rubisco positions from alignment and the dashed lines indicate  
293 select specific interactions to the corresponding scaffold element shown with salt bridges in black and  
294 cation- $\pi$  interactions in green. “*Hnea*” is the  $\alpha$ -carboxysomal Form IA Rubisco from *Halothiobacillus*  
295 *neapolitanus* with CbbL (in orange/yellow) and CbbS (in cyan). The N\*-peptide-bound structure is from  
296 the current study with PDB ID: XXXX. “*Selon*” is the  $\beta$ -carboxysomal Form IB Rubisco from  
297 *Synechococcus elongatus* PCC 7942 with large subunit, RbcL, and small subunit, RbcS, both in grey.  
298 The bound small subunit-like repeat, CcmM-SSUL1, is shown in green. The atomic model was  
299 determined from cryo-electron microscopy single particle analysis and has PDB ID: 6HBC.

300

301



303  
304 **Figure S11**

305  
306 **a-e**, 4-20% SDS-PAGE gels of carboxysome purifications for each of the NTD truncation constructs,  
307 shown schematically. The carboxysomes were purified according to established protocols culminating  
308 with ultracentrifugation on a sucrose step gradient having 5-mL layers of 10, 20, 30, 40, and 50% (w/v)  
309 sucrose. Each fraction was 1 mL. The “pel” fraction is the resuspended pellet from bottom of the gradient.  
310 Normal carboxysomes typically occur as a broad band peaked around 18 mL. **f**, Rubisco large subunit gel  
311 density as a function of fraction volume. As the major carboxysome component, Rubisco is a sensitive  
312 proxy for the intact particles. Only full length CsoS2 and one retaining two of the N-peptide repeats (ΔN1,  
313 N2) resulted in purifiable carboxysomes.

314 *Protein sequences*

315 Select features and mutation sites are indicated by highlights and described after each  
316 sequence.

317

318 ***Halothiobacillus neapolitanus* Rubisco:**

319 CbbL (large subunit):

320 MSAVKKYSAGVKEYRQTYWMPEYTPLDSDILACFKITPQPGVDREEAAAAVAESSTGTWTTV  
321 WTDLLTDMDYKGRAYRIEDVPGDDAAFYAFIAYPIDLFEEGSVNVFTSLVGNVFGFKAVRGL  
322 RLEDVRFPLAYVKTCGGPPHGIQVERDKMNKYGRPLLGTIKPKLGLSAKNYGRAVYECLRGG  
323 LDFTKDDENINSQPFMRWRDRFLVQDATETAEAQTGERKGHYLNVTAPTPEEMYKRAEFAK  
324 EIGAPIIMHDYITGGFTANTGLAKWCQDNGVLLHIHRAMHAVIDRNPNGIHFVRLTKILRLSGG  
325 DHLHTGTVVGKLEGDRASTLGWIDLLRESFIPEDRSRGIFFDQDWGSMPGVFAVASGGIHWVH  
326 MPALVNIFGDDSVLQFGGGTLGHPWGNAAGAAANRVALEACVEARNQGRDIEKEGKEILTAAA  
327 QHSPELKIAMETWKEIKFEFDTVDKLDTQNRWSHPQFEK

328

329 CbbS (small subunit):

330 MAEMQDYKQSLKYETFSYLPPMNAERIRAQIKYAIQGWSPGIEHVEVKNSMNQYWYMWKLP  
331 FFGEQNVNDNLAEIEACRSAYPTHQVKLVAYDNYAQLGLAFVYVRGN

332

333 StrepII affinity tag

334 CbbL Y72; Y72A or Y72R

335 CbbL F346; F346A

336 CbbS Y96; Y96A

337

338 ***Halothiobacillus neapolitanus* CsoS2:**

339 MHHHHHHPSQSGMNPADLSGLSGKELARARRAALSKQGKAAVSNKTASVNRSTKQAASSINT  
340 NQVRSSVNEVPTDYQMADQLCSTIDHADFGTESNRVRDLCRQRREALSTIGKKAVKTNGKPS  
341 GRVRPQQSVVHNDAMIENAGDTNQSSSTSLNNEELSEICSIADDMPERFGSQAKTVRDICRARR  
342 QALSERGTRA VPPKPQSQGGPGRNGYQIDGYLDTALHGRDAAKRHREMLCQYGRGTAP SCK  
343 PTGRVKNSVQSGNAAPK K VETGHTLSGGSVTGTQVDRKSHVTGNEPGTCRAVTGTEYVGTE  
344 QFTSFCNTSPKPNATKVNVT T TARGRPVSGTEVSRTEKVTGNESGVC RNV TGTEYMSNEAHF  
345 SLCGTAAPKPSQADKVMFGATARTHQV VSGSDEF RPSSVTGNESGAKRTITGSQY ADEGLARL  
346 TINGAPAKVARTHTFAGSDVTGTEIGRSTRVTGDESGSCR SISGTEYLSNEQFQSFC DTKPQR  
347 SPFKVGGDR TNKGQSVTGNL VDRSELVTGNEPGSCSRVTGSQY GQSKICGGGVG KVRSMRT  
348 LRGT SVSGQLDHAPKMSGDERGGC MPVTGNEY YGREHFEPFC TSTPEPEAQ STEQSLTCE  
349 GQIISGTSVDASDLVTGNEIGEQLISGDAYVGAQQTGCLPTSPRFNQTGNVQSMGFKNTNQP  
350 EQNFAPGEVMPTDFSIQTPARSAQN RITGNDIAPSGRITGPGMLATGLITGTPEFRHAARELVG  
351 SPQPMAMAMANRNKAAQAPVVQPEVVATQEKPELVCAPRSDQMDRVSGEGKERCHITGDD  
352 WSVNKHITGTAGQWASGRNPSMRGNARVVETS AFANRNV PKPEKPGSKITGSSGNDTQGSLI  
353 TYSGGARG

354

355

356 Hexahistidine affinity tag

357 NTD peptide repeats

358 Middle region peptide repeats

359 CTD peptide repeats

360 CTP

361

362

363 **CsoS2 N-terminal domain (NTD):**

364 MSHHHHHHPSQSGMNPADLSGLSGKELARARRAALSKQGKAAVSNKTASVNRSTKQAASSIN  
365 TNQVRSSVNEVPTDYQMADQLCSTIDHADFGTESNRVRDLCRQRREALSTIGKKAVKTNGKPS  
366 GRVRPQQSVVHNDAMIENAGDTNQSSSTSLNNEELSEICSIADDMPERFGSQAKTVRDICRARR  
367 QALSERGTRA VPPKPQSQGGPGRNGYQIDGYLDTALHGRDAAKRHREMLCQYGRGTAP SCK  
368 PTGRVKNSVQSGNAAPK K V

369

370 Hexahistidine affinity tag

371 NTD peptide repeats

372 Basic residues in N-peptides all mutated to alanines

373

374

375 **CsoS2 Middle Region (MR):**  
376 MSHHHHHHAPK K V E T G H T L S G G S V T G T Q V D R K S H V T G N E P G T C R A V T G T E Y V G T E Q F T S F C  
377 N T S P K P N A T K V N V T T T A R G R P V S G T E V S R T E K V T G N E S G V C R N V T G T E Y M S N E A H F S L C G T A  
378 A K P S Q A D K V M F G A T A R T H Q V V S G S D E F R P S S V T G N E S G A K R T I T G S Q Y A D E G L A R L T I N G A P A  
379 K V A R T H T F A G S D V T G T E I G R S T R V T G D E S G S C R S I S G T E Y L S N E Q F Q S F C D T K P Q R S P F K V G Q  
380 D R T N K G Q S V T G N L V D R S E L V T G N E P G S C S R V T G S Q Y G Q S K I C G G G V G K V R S M R T L R G T S V S  
381 G Q Q L D H A P K M S G D E R G G C M P V T G N E Y Y G R E H F E P F C T S T P E P E A Q

382  
383 Hexahistidine affinity tag  
384 Middle region peptide repeats

385  
386  
387  
388 **CsoS2 C-terminal domain (CTD):**  
389 MSHHHHHH T S T P E P E A Q S T E Q S L T C E G Q I I S G T S V D A S D L V T G N E I G E Q Q L I S G D A Y V G A Q Q T  
390 G C L P T S P R F N Q T G N V Q S M G F K N T N Q P E Q N F A P G E V M P T D F S I Q T P A R S A Q N R I T G N D I A P S G  
391 R I T G P G M L A T G L I T G T P E F R H A A R E L V G S P Q P M A M A M A N R N K A A Q A P V V Q P E V V A T Q E K P E L V  
392 C A P R S D Q M D R V S G E G K E R C H I T G D D W S V N K H I T G T A G Q W A S G R N P S M R G N A R V V E T S A F A N  
393 R N V P K P E K P G S K I T G S S G N D T Q G S L I T Y S G G A R G

394  
395 Hexahistidine affinity tag  
396 CTD peptide repeats  
397 CTP

398  
399  
400 **N\*-polyProline:**  
401 M S W K H H H H H H E N L Y F Q S A A V G G G S G G G S G G P P P P P P P P P A P A P A P A P P P P P P P P P A P A P A P  
402 A P P P P P P P P G G G S G G G S G G G R D L A R A R R E A L S Q Q G K A A V G G G S G G G S G G S G

403  
404 Hexahistidine affinity tag  
405 polyproline II helix  
406 NTD peptide consensus repeat (N\*)

407  
408  
409  
410 **random N\*-polyProline:**  
411 M S W K H H H H H H E N L Y F Q S A A V G G G S G G G S G G P P P P P P P P P A P A P A P A P P P P P P P P P A P A P A P  
412 A P P P P P P P P G G G S G G G S G G G R R K L R A A G R A L Q V E Q A D S R A G G G S G G G S G G S G

413  
414 Hexahistidine affinity tag  
415 polyproline II helix  
416 random NTD peptide

417 **polyProline flexible polar AA control:**  
418 MSWKHHHHHHENLYFQSAAVGGGSGGGSGGPPPPPPPPAPAPAPAPPPPPPPPPAPAPAP  
419 APPPPPPPPGGGSGGGSGGGGTGSTGSGSSSSSGSGTSGTGGGGSGGGSGGGSG  
420  
421 Hexahistidine affinity tag  
422 polyproline II helix  
423  
424  
425  
426 **MR R1-polyProline:**  
427 MSWKHHHHHHENLYFQSAAVGGGSGGGSGGPPPPPPPPAPAPAPAPPPPPPPPPAPAPAP  
428 APPPPPPPPGGGSGGGSGGGKVVETGHTLSGGSVTGTQVDRKSHVTGNEPGTCRAVTGTEYV  
429 GTEQFTSFCGGGSGGGSGGGSG  
430  
431 Hexahistidine affinity tag  
432 polyproline II helix  
433 MR repeat 1  
434  
435  
436  
437 **CTD R8-polyProline:**  
438 MSWKHHHHHHENLYFQSAAVGGGSGGGSGGPPPPPPPPAPAPAPAPPPPPPPPPAPAPAP  
439 APPPPPPPPGGGSGGGSGGGSIQTPARSAQNRITGNDIAPSGRITGPGMLATGLITGTPEFGG  
440 GSGGGSGGGSG  
441  
442 Hexahistidine affinity tag  
443 polyproline II helix  
444 CTD repeat 8  
445  
446  
447 **CTP-polyProline:**  
448 MSWKHHHHHHENLYFQSAAVGGGSGGGSGGPPPPPPPPAPAPAPAPPPPPPPPPAPAPAP  
449 APPPPPPPPGGGSGGGSGGGVVKPEKPGSKITGSSGNDTQGSLITYSGGARGGGGSGGGSG  
450 GSG  
451  
452 Hexahistidine affinity tag  
453 polyproline II helix  
454 CTP  
455

456 **[N1-N2]-GFP:**  
457 MHHHHHHENLYFQSPSQSGMNPADLSGLSGKELARARRAALSKQGKAAVSNKTASVNRSTK  
458 QAASSINTNQVRSSVNEVPTDYQMADQLCSTIDHADFGTESNRVRDLCRQRREALSTIGKKAV  
459 KTNGKPSGRVRPQQSGSGSGGSKGEELFTGVVILVELDGDVNGHKFSVRGEGEGDATNGK  
460 LTLKFICTTGKLPVPWPTLVTTLYGVQCFSRYPDHMKQHDFFKSAMPEGYVQERTISFKDDGT  
461 YKTRAEVKFEGDTLVNRIELKGIDFKEDGNILGHKLEYNFNHNVYITADKQKNGIKANFKIRHNV  
462 EDGSVQLADHYQQNTPIGDGPVLLPDNHYLSTQSKLSKDPNEKRDHMLLEFVTAAGITHGMD  
463 ELYKWSHPQFEK

464  
465 Hexahistidine affinity tag  
466 NTD peptide repeats (N1 and N2)  
467 NTD flexible “interstitial” sequence  
468 Superfolder GFP<sup>6</sup> (with dimer abolishing K206 variant)<sup>7</sup>

469  
470  
471 **[N1]-GFP:**  
472 MSHHHHHHPSQSGMNPADLSGLSGKELARARRAALSKQGKAAVSNKTASVNRSTKQAASSIN  
473 TNQVRSSVNEVPTDYQMADQLCSTIDHADFGTESNRVENLYFQSSGSGSGGSKGEELFTGVV  
474 PILVELDGDVNGHKFSVRGEGEGDATNGKLTTLKFICTTGKLPVPWPTLVTTLYGVQCFSRYPD  
475 HMKQHDFFKSAMPEGYVQERTISFKDDGTYKTRAEVKFEGDTLVNRIELKGIDFKEDGNILGHK  
476 LEYNFNHNVYITADKQKNGIKANFKIRHNVEDGSVQLADHYQQNTPIGDGPVLLPDNHYLSTQ  
477 SKLSKDPNEKRDHMLLEFVTAAGITHGMD ELYKWSHPQFEK

478  
479 Hexahistidine affinity tag  
480 NTD peptide repeat (N1)  
481 NTD flexible “interstitial” sequence  
482 Superfolder GFP (with dimer abolishing K206 variant)

483  
484  
485  
486

487 **Rubisco-N\* fusion (CbbL-N\*, CbbS):**

488 CbbL-N\*:

489 MSAVKKYSAGVKEYRQTYWMPEYTPLDSDILACFKITPQPGVDREEAAAAVAESSTGTWTTV  
490 WTDLLTDMDYKGRAYRIEDVPGDDAAFYAFIAYPIDLFEEGSVVNVFTSLVGNVFGFKAVRGL  
491 RLEDVRFPLAYVKTCGGPPHGIQVERDKMNKYGRPLLGCTIKPKLGLSAKNYGRAVYECLRGG  
492 LDFTKDDENINSQPFMRWRDRFLFVQDATETAEAQTGERKGHYLNVTAPTPEEMYKRAEFAK  
493 EIGAPIIMHDYITGGFTANTGLAKWCQDNGVLLHIHRAMHAVIDRNPNGIHFRLTKILRLSGG  
494 DHLHTGTVVGKLEGDRASTLGWIDLLRESFIPEDRSRGIFFDQDWGSMPGVFAVASGGIHWVH  
495 MPALVNI FGDDSVLQFGGGT LGHPWGNAAGAAANRVALEACVEARNQGRDIEKEGKEILTAAA  
496 QHSPELKIAMETWKEIKFEFDTVDKLDTQNRSS **GRDLARARREALSQQGKAAV**GSWSHPQFE

497 **K**

498

499 CbbS:

500 MAEMQDYKQSLKYETFSYLPPMNAERIRAQIKYAI AQGWSPGIEHVEVKNSMNQYWYMWKLP  
501 FFGEQNV DNVLAIEIACRSAYPTHQVKLVAYDNYA QSLGLAFVVYRGN

502

503 **StrepII affinity tag**

504 **NTD consensus repeat**

505



506 *Plasmids used*

507

508 All plasmids in the table below were made for this study.

509 Table S2

<b>Name</b>	<b>Relevant genotype</b>	<b>Resistance</b>
pBz15	wild-type <i>Halothiobacillus neapolitanus</i> Form I Rubisco with StrepII tag (CbbL-StrepII, CbbS), pET-14b backbone	Amp
pLz74	Same as pBz15 but with N*-peptide fusion on CbbL, construct used for protein crystallization	Amp
pBz118	His-tagged <i>Hn</i> CsoS2. This protein contains a programmed ribosomal frameshift and expresses as a short and long form (CsoS2A and CsoS2B, respectively), <sup>8</sup> pET-14b backbone	Amp
pBz109	His-tagged CsoS2 N-terminal domain (NTD), pET-14b backbone	Amp
pBz106	His-tagged CsoS2 Middle Region (MR), pET-14b backbone	Amp
pBz110	His-tagged CsoS2 C-terminal domain (CTD), pET-14b backbone	Amp
pLz47	Consensus N*-peptide after polyproline II helix, N-terminal His-tag, pET-14b backbone	Amp
pLz55	Same as pLz47 but with randomized N*-peptide sequence	Amp
pLz26	His-tagged first two N-peptides of NTD fused to Superfolder GFP ([N1-N2]-GFP) used for MST, pET-14b backbone	Amp
pBz172	His-tagged first N-peptide of NTD fused to Superfolder GFP ([N1]-GFP) used for MST, pET-14b backbone	Amp
pLz37	Derived from pHnCB10 for heterologous expression of carboxysome operon proteins. <sup>9</sup> Contains full-length CsoS2.	Cm
pLz75	Same as pLz37 but with CsoS2 truncated such that NTD repeats N1 and N2 are removed	Cm
pLz76	Same as pLz37 but with CsoS2 truncated such that NTD repeats N1, N2, and N3 are removed	Cm
pLz77	Same as pLz37 but with CsoS2 discontinuously truncated such that repeats N1, N2, and N4 are removed	Cm
pLz78	Same as pLz37 but with CsoS2 truncated to remove the entire NTD (i.e. N1, N2, N3, and N4 are removed)	Cm
pLz54	Same as pLz47 but with 20 polar amino acids instead of N*-peptide	Amp

pLz56	Same as pLz47 but with MR repeat 1 instead of N*-peptide	Amp
pLz57	Same as pLz47 but with CTD repeat 8 instead of N*-peptide	Amp
pLz58	Same as pLz47 but with C-terminal peptide (CTP) instead of N*-peptide	Amp

510

511

512

513 **References**

514

- 515 1. Kane, R. S. Thermodynamics of multivalent interactions: influence of the linker. *Langmuir* **26**,  
516 8636–8640 (2010).
- 517 2. Kitov, P. I. & Bundle, D. R. On the nature of the multivalency effect: a thermodynamic model.  
518 *J. Am. Chem. Soc.* **125**, 16271–16284 (2003).
- 519 3. Cai, F. *et al.* Advances in understanding carboxysome assembly in prochlorococcus and  
520 synechococcus implicate csos2 as a critical component. *Life (Basel)* **5**, 1141–1171 (2015).
- 521 4. Karkehabadi, S., Satagopan, S., Taylor, T. C., Spreitzer, R. J. & Andersson, I. Structural  
522 analysis of altered large-subunit loop-6/carboxy-terminus interactions that influence catalytic  
523 efficiency and CO<sub>2</sub>/O<sub>2</sub> specificity of ribulose-1,5-bisphosphate carboxylase/oxygenase.  
524 *Biochemistry* **46**, 11080–11089 (2007).
- 525 5. Wang, H. *et al.* Rubisco condensate formation by CcmM in  $\beta$ -carboxysome biogenesis.  
526 *Nature* **566**, 131–135 (2019).
- 527 6. Pédelacq, J.-D., Cabantous, S., Tran, T., Terwilliger, T. C. & Waldo, G. S. Engineering and  
528 characterization of a superfolder green fluorescent protein. *Nat. Biotechnol.* **24**, 79–88  
529 (2006).
- 530 7. Zacharias, D. A., Violin, J. D., Newton, A. C. & Tsien, R. Y. Partitioning of lipid-modified  
531 monomeric GFPs into membrane microdomains of live cells. *Science* **296**, 913–916 (2002).
- 532 8. Chaijarasphong, T. *et al.* Programmed Ribosomal Frameshifting Mediates Expression of the  
533  $\alpha$ -Carboxysome. *J. Mol. Biol.* **428**, 153–164 (2016).
- 534 9. Bonacci, W. *et al.* Modularity of a carbon-fixing protein organelle. *Proc Natl Acad Sci USA*  
535 **109**, 478–483 (2012).