

1     **Evolution of Selective RNA Processing and Stabilization operons**  
2                     **in cellulosome-harboring *Clostridium* spp.**

3     **Yogendra Bhaskar<sup>1,3,\*</sup>, Mohammadhadi Heidari B.<sup>1,3</sup>, Chenggang Xu<sup>2</sup>, Jian Xu<sup>1,3,\*</sup>**

4  
5     <sup>1</sup>Single-Cell Center and CAS Key Laboratory of Biofuels and Shandong Key Laboratory of  
6     Energy Genetics, Qingdao Institute of Bioenergy and Bioprocess Technology, Chinese  
7     Academy of Sciences, Qingdao, Shandong, 266101, China

8     <sup>2</sup>College of Life Science, Shanxi University, Taiyuan, Shanxi, 030006, China

9     <sup>3</sup>University of Chinese Academy of Sciences, Beijing, 100049, China

10  
11     \***Correspondence:** Tel: +86 532 8066 2651; Fax: +86 532 8066 2654  
12     E-mail address: 2014in-yogendra@qibebt.ac.cn and xujian@qibebt.ac.cn

13  
14     **Key words:** stem-loop structure; stoichiometry of protein complexes; cellulosome; operon  
15     evolution; bacterial adaptation

16 **Abstract**

17 In selective RNA processing and stabilization (SRPS) operons, the stoichiometry of  
18 encoded proteins is determined by their respective 3'-end stem-loops (SLs), yet the evolution  
19 of this mechanism remains elusive. In cellulosomal operons of *Clostridium* spp., we show  
20 that the SLs and their associated genes form a monogamy companionship during the operon  
21 evolution. Based on  $\Delta G$  of such SLs, we propose CoSLOE (Composite SL-based Operon  
22 Evolution) model with evolutionary ratio (ER)  $>1$  or  $<1$  for positive or negative selection of  
23 SRPS operons. In the composite SL- $\Delta G$ -based tree (CoSL-tree) of cellulosomal operons,  
24 when traversing from leafs to the root nodes, ERs reveal diversifying/positive selection  
25 towards a less efficient cellulosomal system, consistent with glycoside-hydrolase gene  
26 variation both in-operon and genome-wide. A similar pattern is followed by the ATPase  
27 operon and the majority of orthologous SRPS operons genome-wide, suggesting conservation  
28 among operons in such selection. Thus SRPS operons via their transcript-stabilizing non-  
29 coding elements are highlighting a link between operon stoichiometry and operon evolution.

## 30 1 Introduction

31 In bacterial genomes, ~50% of the genes are organized and regulated in the form of  
32 operon (Osbourn & Field *et al.*, 2009). Within an operon, to ensure proper absolute and  
33 relative abundance of the component genes, one strategy adopted by certain bacteria is  
34 selective RNA processing and stabilization (SRPS), where the RNA molecule is cleaved by  
35 ribonuclease into fragments, and then with the involvement of the specific *cis*-elements  
36 (Stem-loops), mature mRNA transcripts stabilize to differential gene expression and  
37 eventually to the protein complex (Rochat *et al.*, 2013). The SRPS mechanism controls  
38 operons that encode a variety of key protein complexes and regulatory pathways such as the  
39 glycolysis pathway, maltose transport system, cellulosome complex and photosynthetic  
40 apparatus (Newbury *et al.*, 1987, Klug *et al.*, 1993, Ludwig *et al.*, 2001, Xu *et al.*, 2015).

41 Using the cellulosome-encoding *cip-cel* operon of *Clostridium cellulolyticum* (*Ccel*) as a  
42 model, we showed that the stem-loops generally located at the 3'-end of regulated genes  
43 precisely regulate structure and relative abundance of the subunit-encoding transcripts  
44 processed from a primary polycistronic RNA (Xu *et al.*, 2015). Importantly, the “ratio” of  
45 subunit-encoding transcripts for the *cip-cel* operon, which quantitatively specifies  
46 cellulosome stoichiometry, appears to be encoded by the genome (i.e., organism-specific) and  
47 insensitive to alteration of culture conditions, since change among glucose, cellobiose and  
48 cellulose did not result in ratio change (Xu *et al.*, 2015). These findings revealed a key role of  
49 such stem-loops (SLs; i.e., all such SLs present in a SRPS operon) in specifying proper  
50 function of SRPS-operon-encoded protein complexes (or metabolic pathways). Moreover,  
51 they strongly suggest potential links between the structure and function of these stem loops to  
52 organismal evolution.

53 However, key questions remain unanswered: (i) how do these SLs evolve? How  
54 conserved are these SLs among orthologous operons? What is the nature of such conservation?  
55 (ii) What is the link in evolution between these SLs and their companion genes in SRPS  
56 operons? (iii) What roles do these SLs play in the evolution of SRPS operons? Are these roles  
57 conserved for SPRS operons at a genome-wide scale? How similar or divergent are these  
58 roles across different genomes? Can evolution of SRPS operons be quantitatively modeled  
59 via these SLs? Here in cellulosomal operons of *Clostridium* spp., based on  $\Delta G$  of such SLs,  
60 we propose CoSLOE (Composite SL-based Operon Evolution) with evolutionary ratio  
61 (ER)  $>1$  or  $<1$  for positive or negative selection of SRPS operons. In the CoSL-tree of  
62 cellulosomal operons, when traversing from leafs to the root nodes, ERs reveal  
63 diversifying/positive selection towards a less efficient cellulosomal system, consistent with  
64 glycoside-hydrolase gene variation both in-operon and genome-wide. A similar pattern is  
65 followed by the ATPase operon and the majority of orthologous SRPS operons genome-wide,  
66 suggesting conservation among operons in such selection. Thus SRPS operons via their  
67 transcript-stabilizing non-coding elements are highlighting a link between operon  
68 stoichiometry and operon evolution.

## 69 2 Materials and Methods

## 70 2.1 Prediction of the stable stem-loops

71 SLs in *C. cellulolyticum* were predicted via the following steps (Bhaskar *et al.*, 2021). (i)  
72 Prediction of motifs using RNAmotif (Macke *et al.*, 2001); (ii) Estimation of free-energy and  
73 RNA secondary structure using RNAfold (Hofacker, 2003); (iii) Genome mapping of the  
74 predicted SLs; (iv) Screening of the SLs based on their stability for highly stable SLs. These  
75 stable SLs were then mapped to the operon map of the respective species, followed by  
76 functional classification based on the derived classification rules, whereby the SRPS operons  
77 were identified. Genome-encoded ratios were predicted for these SRPS operons using the  $\Delta G$   
78 of the harbored SLs (Bhaskar *et al.*, 2021).

## 79 2.2 Calculation of the $\Delta G$ -based ratio for an SRPS operon

80 Ratios were calculated in the SRPS operon using the  $\Delta G$  (free-energy) of the SLs present  
81 in and flanking the operon (**Fig. 1A**). For example, the ratio for a four-gene operon (with SLs  
82 found after first two genes and at the end of operon) “Gene-1 ( $\Delta G_1$ ) : Gene-2 ( $\Delta G_2$ ) : Gene-  
83 3 : Gene-4 ( $\Delta G_4$ )” would be “ $\Delta G_1$  :  $\Delta G_2$  :  $\Delta G_4$  :  $\Delta G_4$ ”. To normalize the ratio,  $\Delta G$  of all SLs  
84 in an operon were divided by the sum of all  $\Delta G$  (**Table S1**).

## 85 2.3 Phylogenetic analysis of the cellulosomal and the ATP synthase operons

86 The genomes and associated annotations of 13 cellulosome operon-harboring Clostridial  
87 species including *Ruminiclostridium cellulolyticum* H10 (*Ccel*; NC\_011898.1),  
88 *Ruminiclostridium papyrosolvans* DSM 2782 (*Cpap*; GCF\_000175795.2), *Clostridium*  
89 *saccharoperbutylacetonicum* (*Csac*; NC\_020291.1), *Clostridium* sp. BNL1100 (*Cbnl*;  
90 GCF\_000244875.1), *Clostridium felsineum* DSM 794 (*Cfel*; GCF\_002006355.1),  
91 *Ruminiclostridium josui* JCM 17888 (*Cjos*; GCF\_000526495.1), *Ruminiclostridium*  
92 *cellobioparum* subsp. *termitidis* CT1112 (*Cter*; GCF\_000350485.1), *Clostridium*  
93 *acetobutylicum* ATCC 824 (*Cace*; NC\_015687.1), *Clostridium cellulovorans* 743B (*Cloc*;  
94 NC\_014393.1), *Ruminiclostridium hungatei* DSM 14427 (*Chun*; GCF\_002051585.1),  
95 *Clostridium puniceum* DSM 2619 (*Cpun*; GCF\_002006345.1), *Clostridium roseum* DSM  
96 7320 (*Cros*; GCF\_002006215.1) and *Ruminiclostridium cellobioparum* DSM 1351=ATCC  
97 15832 (*Ccell*; GCF\_000621505.1) (**Table S1**), were downloaded from NCBI. Cellulosome  
98 operons from *Cpap*, *Csac*, *Cbnl*, *Cfel*, *Cjos*, *Cter*, *Cace*, *Cloc*, *Chun*, *Cpun*, *Cros* and *Ccell*  
99 were identified by the available annotation and BLAST (Altschul *et al.*, 1990), where the *cip-*  
100 *cel* operon (encoding the cellulosome) from *Ccel* was used as a query with the e-value cutoff  
101 of  $1e-5$ . Organismal phylogeny (16S-tree) of these species was derived using the 16S rRNA  
102 sequence, where all positions containing gaps and missing data were eliminated, which  
103 resulted in a total of 1,326 positions in the final multiple-sequence alignment. Phylogenetic  
104 analyses for the cellulosome and the ATPase operons were conducted in MEGA7 (Kumar *et*  
105 *al.*, 2016) via the Maximum Likelihood method.

106 The  $\Delta G$ -based dendrogram of SLs was performed using the *pvclust* (Suzuki &  
107 Shimodaira, 2006) package in R (CRAN <http://cran.r-project.org/>) (**Fig. S1A**). To calculate  
108 the SLs'  $\Delta G$ -based dendrogram (CoSL-tree),  $\Delta G$  of all SLs in an operon were divided by the  
109 sum of all  $\Delta G$ , which generated a normalized proportion for an operon, and empty cells (i.e.,

110 values are “non-applicable”) were replaced by the average value of that proportion while  
111 clustering (**Table S1**). The normalized  $\Delta G$  proportions of 13 clostridia species were supplied  
112 to *pvclust* with the Euclidean distance method and the *ward.D2* hierarchical clustering, with  
113 bootstrapping for 1000 times. The Ka/Ks values for genes were calculated using the Codeml  
114 tool of *PAML* package (Yang *et al.*, 2007).

## 115 **2.4 Structural alignment analysis of orthologous SLs**

116 The orthologous SLs from the cellulosomal operons and the ATP synthase operons were  
117 aligned structurally and sequence-wise using the *LocARNA* alignment and folding tool (Smith  
118 *et al.*, 2010, Will *et al.*, 2012). Evolution of the orthologous SLs was shown using the  
119 multiple-alignment of SL sequences and dot-bracket notations.

## 120 **2.5 Derivation of Composite Stem-Loop-based Operon Evolution (CoSLOE) model**

121 The CoSLOE model was described via an equation that calculates the evolutionary ratio  
122 (ER):

$$123 \quad \frac{G1}{G2} \times \frac{S1}{S2} \times \frac{CV1}{CV2} = ER$$

124 where G1 and G2, S1 and S2, CV1 and CV2 are the number of genes, number of SLs and  
125 coefficient of variations (CV) respectively, in the two operons. CV is the ratio of standard  
126 deviation ( $\text{Ratios}_D$ ) and mean ( $\overline{\text{Ratio}}$ ) of the  $\Delta G$ -based ratio of an operon.

## 127 **3 Results**

### 128 **3.1 Phylogenetic analysis of the SLs in cellulosome-encoding SRPS operons from 13** 129 **Clostridial genomes**

130 To probe the role of SLs in the function and evolution of SRPS operons at the whole-  
131 genome scale, we developed an approach to identify the SRPS operons based on the genome-  
132 wide predicted stable stem-loops (SLs) and then use the free-energy ( $\Delta G$ ) of these stable SLs  
133 to calculate ratios of SRPS operons (Bhaskar *et al.*, 2021) (**Fig. 1A**). The  $\Delta G$ -based ratios  
134 were calculated for the cellulosome complex operon (*cip-cel*) in *Ccel*, which can model  
135 stoichiometry of the encoded complex. To probe how this mechanism has evolved, we  
136 extended the analysis to twelve additional mesophilic Clostridial spp.: *C. papyrosolvans*  
137 (*Cpap*), *C. saccharoperbutylacetonicum* (*Csac*), *C. sp. BNL1100* (*Cbnl*), *C. felsineum* (*Cfel*),  
138 *C. josui* (*Cjos*), *C. termitidis* (*Cter*), *C. acetobutylicum* (*Cace*), *C. cellulovorans* (*Cloc*), *C.*  
139 *hungatei* (*Chun*), *C. puniceum* (*Cpun*), *C. roseum* (*Cros*), and *C. cellobioparum* (*Ccell*; **Table**  
140 **S1**). These operons are orthologous, as indicated by orthology of genes, functional  
141 conservation of encoded proteins and the global similarity in operon structure. Our  $\Delta G$ -based  
142 method predicted 7, 7, 5, 5, 5, 6, 5, 4, 3, 3, 3 and 3 SLs in *Cbnl*, *Cpap*, *Cjos*, *Cter*, *Ccell*,  
143 *Chun*, *Cloc*, *Cace*, *Cpun*, *Cros*, and *Cfel* respectively (**Table S1**). The  $\Delta G$ -based ratio for  
144 these Clostridial species were also highly skewed, e.g., the ratios of *Cbnl* and *Cpap* are “-  
145 24.4:-26.3:-25.9:-25.9:-25.9:-15.3:-15.3:-21.2:-21.2:-18.3:-21.5:-21.5” and “-23.6:-26.3: -

146 25.3:-25.3:-25.3:-16.7:-16.7:-16.8:-16.8:-23.5:-23.9:-23.9” respectively. Similarly, *Ccell* and  
147 *Cter* exhibit identical ratios, so do *Cace*, *Cfel* and *Cros* (**Fig. 1B**; **Table S1**).

148 To probe how such operon properties have evolved, the  $\Delta G$ -based proportions of all  
149 harbored SLs in an operon (which we termed “composite SLs” or CoSL) were used to  
150 generate a dendrogram (CoSL-tree; **Fig. S1A**). CoSL-tree was then compared to the 16S  
151 rRNA-based tree (16S-tree; i.e., the organismal phylogeny; **Fig. S1B**). Predicted ratios from  
152 the 13 cellulosomal operons were combined to form a data matrix, which was then used for  
153 the hierarchical clustering with 1000 iterations to generate the ratio-based tree (**Methods**).  
154 Intriguingly, the species were classified differently in the two clades derived from CoSL-tree  
155 (**Fig. S1A**) and 16S-tree (**Fig. S1B**). For example, (i) *Cace* and *Cros* are in Clade 1 of CoSL-  
156 tree, yet found in Clade 2 of 16S-tree; (ii) *Cpun* is an out-group in CoSL-tree, whilst *Cfel* is  
157 an out-group in 16S-tree; (iii) *Cros* and *Cfel* are clustered in CoSL-tree yet distantly apart in  
158 16S-tree. Such difference between CoSL-tree and 16S-tree indicates the deviation of SRPS  
159 operon evolution from organismal taxonomy.

### 160 3.2 Gene-SL relationship during evolution of Clostridial cellulosomal operons

161 To probe the roles of SLs in cellulosomal operon evolution, seven orthologous SLs were  
162 first identified in the intergenic regions of the 13 orthologous cellulosomal operons, via  
163 comparison of their sequences, structures and organization in the operons (**Fig. 1C**; **Fig. 2**;  
164 **Fig. S2**). However, not all the Clostridial species harbor similar numbers of orthologous SLs  
165 and at identical positions (**Fig. 1C**; **Table S1**): 7 SLs in *Cbnl* and *Cpap*, 6 in *Ccel* and *Chun*,  
166 5 in *Cjos*, *Cter*, *Ccell* and *Cloc*, 4 in *Cpun* and *Cace* and 3 in *Cfel*, *Cros* and *Csac*. The  
167 presence of these SLs suggests SPRS mechanisms in these 13 cellulosomal operons (for *Cloc*,  
168 the role of multiple promoters is also involved (Doi *et al.*, 1998)).

169 Interestingly, although the region between a SL and its associated genes can be inserted  
170 by another gene, the SLs are always positioned with their associated genes in a sequential  
171 fashion that is conserved among a set of orthologous operons. Thus, to probe the gene–SL  
172 relationship, orthologous SLs were aligned via sequence and structural similarity (**Fig. S2**).  
173 Compatible base pairs (in the stem sequences) were found in SL-1, 2, 3, 5 and 7,  
174 underscoring the structural similarity among the orthologous SLs (**Fig. 2**). Specifically, (i)  
175 SL-1 is present in all the Clostridial species (except *Cloc*), and SL-1 and 2, in their  
176 respective clades, are of similar length and identical  $\Delta G$  to other orthologous SLs, yet show  
177 higher variation in their loop sequences (**Fig. 2A, B**; **Fig. S2A, B**); (ii) SL-3 shows less  
178 sequence variation in the two clades than SL-1 and 2, possibly due to its role as terminator  
179 SLs (**Fig. 2C**; **Fig. S2C**); (iii) SL-4, 5 and 6 are clade-specific, as they are absent in Clade 2  
180 except the SL-4 in *Cloc* (**Fig. 2D**; **Fig. S2D, E, F**); (iv) similar to SL-3, SL-7 carries a low  
181 level of sequence variation (**Fig. 2E**; **Fig. S2G**). Such variation in SL sequence and structure  
182 depicts their evolutionary distance.

183 Intriguingly, a dockerin-encoding gene, located at the 7<sup>th</sup> position of operon in *Cloc*, the  
184 9<sup>th</sup> in subclade of *Ccel-Cjos-Cbnl-Cpap* (except *Cjos*) and the 12<sup>th</sup> in Subclade 1.1 species  
185 (except *Chun*) is always controlled by the orthologous SL-5 (**Fig. 1C, 2D**; **Fig. S2E**).  
186 Similarly, a cellulase-encoding gene, situated at distinct positions among cellulosomal

187 operons, is controlled by an orthologous SL-3. In addition, clade-specific derivative  
188 homologous SLs in the cellulosomal operons also show such loyalty with their respective  
189 companion genes, e.g., (i) *Cloc* harbors an extra cellulase-encoding gene carrying SL-3A  
190 (homologous to SL-3; **Fig. S2H**); (ii) SL-2A is found in *Ccel* and *Chun*, similar to SL-2 (**Fig.**  
191 **S2I**); (iii) SL-7A is found in *Csac*, which is similar to SL-7 and works as terminator to the  
192 operon (**Fig. S2J**). These observations suggest monogamy as one feature of the gene-SL  
193 relationship during evolution of SRPS operons.

### 194 **3.3 The Composite Stem-Loop based Operon Evolution (CoSLOE) model for SPRS** 195 **operons**

196 Taking advantage of the link between SLs and evolution of operon, we propose  
197 Composite Stem-Loop based Operon Evolution (CoSLOE) for the SRPS operons (**Fig. 3**).  
198 The model consists of (**Equation I**): (i) the number of genes in the operon (G), where the  
199 addition of one gene shows the positive selection, while an equal number of genes suggests  
200 neutral operons; (ii) the number of SLs (S), which plays crucial roles in regulation,  
201 stabilization and termination of genes; (iii) variance of  $\Delta G$  of the SLs in operons (CV), where  
202 multiple SLs with distinct free-energy together specify and control the stoichiometry of gene  
203 expression. Therefore, the evolutionary ratio (ER) of an operon with respect to the other  
204 operons is,

$$205 \quad \frac{G1}{G2} \times \frac{S1}{S2} \times \frac{CV1}{CV2} = ER \text{ (for ideal condition, ER =1)} \quad \text{(I)}$$

206 where G1 and G2, S1 and S2, CV1 and CV2 are the number of genes, number of SLs and  
207 coefficient of variations (CV) respectively, in the two operons. CV is ratio of standard  
208 deviation and mean of the ratio of  $\Delta G$  of SLs for an operon. Positive or purifying selection of  
209 the operon is indicated by  $ER > 1$  and  $ER < 1$  respectively, while ER of 1 corresponds to  
210 neutral selection (i.e., ideal condition).

### 211 **3.4 CoSLOE reveals selection pressure on the cellulosomal operons**

212 To probe their evolution, pairwise ERs for the 13 cellulosomal operons in CoSL-tree  
213 were derived via CoSLOE (**Equation I**; **Fig. 3**). In Subclade 1.2 (**Fig. 1C**), (i) *Cbnl* and  
214 *Cpap* show ER of 1.01 and 0.99 with each other respectively, suggesting that the selection  
215 pressure is almost neutral and *Cbnl* is positively selected towards the root; (ii) the next  
216 nearest species is *Cjos*, which lacks one gene and two SLs possibly due to the deletion or  
217 horizontal transfer of genes, shows the ER of 0.57, 0.57 and 0.52 with *Cpap*, *Cbnl* and *Ccel*  
218 respectively (**Table 1A**), i.e. equally separated from all the three clostridia; (iii) however, the  
219 Ka/Ks values, at the gene level selection, for the first gene of *Cjos* are 1.42, 1.45, and 1.5  
220 with *Ccel*, *Cbnl*, *Cpap* respectively (**Table S2**), suggesting the first gene of these operons is  
221 under positive selection towards *Cjos*; (iv) the operon ER for *Ccel* is 1.93 1.10 and 1.10 with  
222 *Cjos*, *Cbnl* and *Cpap* respectively (**Table 1A**), which depict the positive selection with the  
223 addition of a new gene at 11<sup>th</sup> place in operon (**Fig. S5**). Taken together, in Subclade 1.2 of  
224 CoSL-tree, species are under positive selection while going from *Cpap* to *Ccel* (**Fig. 1C**), and

225 also while going from *Cros* to *Cace* (due to the much higher *Cace-Cros* ER of 7.76 than *Cfel-*  
226 *Cros* ER of 1.05; **Table 1B**).

227 In Subclade 1.1, the SLs in *Ccell* operon are more similar to *Cter* than to *Chun*, *i.e.*, the  
228 operon ERs for *Cter-Ccell* and *Chun-Ccell* are 1.10 and 1.83 respectively (**Table 1A**), while  
229 those for *Cter-Chun* and *Ccell-Chun* are 0.60 and 0.55 respectively. Thus *Ccell* and *Cter* are  
230 in purifying selection, while Subclade 1.1 is under positive selection towards *Chun* (similar to  
231 as Subclade 1.2; **Fig. 1C**; **Fig. S5**).

232 The Clade 2 species in CoSL-tree are more dynamic in evolution than Clade 1, in that  
233 they show more variable number of genes and SLs. *Cloc*, *Csac* and the out-grouped *Cpun*  
234 exhibit a certain degree of similarity to the Clade 1 species, but feature the addition of new  
235 SLs such as SL-3A and SL-7A (homologous to SL3 and SL7 respectively; **Fig. 1C**).  
236 Moreover, their cellulosomal operons are distinct, *e.g.*, *Cpun* and *Csac* operons harbor no  
237 cohesin, glycoside hydrolase (GH) or dockerin genes. In fact, ERs for *Cpun* and *Csac* versus  
238 *Cloc* are 1.37 and 2.11 respectively (**Table 1B**), consistent with positive selection.

239 Notably, if the ERs are calculated without considering SLs (and the CV) in **Equation 1**,  
240 then the number of genes by itself is not sufficient to detect the selection. For example, in  
241 Subclade 1.2 (*Cpap-Cbnl-Cjos-Ccel*; *Cros-Cfel-Cace*), the equal number of genes would  
242 suggest an ER of 1, which however is misleading. Therefore, in computing CoSL-based ER,  
243 the SLs are essential for deriving ERs in CoSLOE.

### 244 **3.5 The CoSLOE model of cellulosomal operons is supported by variation in enzyme** 245 **genes**

246 In CoSLOE, purifying/negative selection occurs when the tree is traversed from the root  
247 to the leaf nodes, and diversifying/positive selection takes place when traversing from leafs to  
248 the root nodes (**Fig. 4**). In the cellulosomal operon (**Fig. 1C**), positive evolution takes place in  
249 the *Ccel-Cjos-Cbnl-Cpap* direction (root to leaf), in the *Chun-Cter-Ccell* direction and in the  
250 *Cace-Cros-Cfel* direction respectively, with the root node being the most positively selected  
251 and the leaf nodes the most negatively selected.

252 To probe the biological significance of these findings, the genome-wide numbers of  
253 carbohydrate-active enzymes (CAZymes) and carbohydrate-binding module (CBM) were  
254 compared, since these enzymes are major parts of the cellulosomal system (Busch *et al.*,  
255 2017). For example, in Subclade 1.2, for the *Ccel*, *Cjos*, *Cbnl* and *Cpap* genome (which  
256 exhibit > 95% similarity in 16S rRNA sequences; **Fig. 1C, 6A**), (i) the CAZymes (including  
257 glycoside hydrolases or GHs, carbohydrate esterases or CEs, and polysaccharide lyases or  
258 PLs) harbored is 111 (94 GHs, 13 CEs, 4 PLs), 116 (92 GHs, 19 CEs, 5 PLs), 127 (103 GHs,  
259 19 CEs, 5 PLs) and 122 (103 GHs, 16 CEs, 3 PLs) respectively (Dassa *et al.*, 2017),  
260 exhibiting an overall pattern of increase; (ii) for GH5 (*Ccel*: 7; *Cjos*: 7; *Cbnl*: 8; *Cpap*: 7),  
261 GH9 (13, 14, 14, 14) and GH43 (9, 13, 13, 13), an increase in number is apparent when  
262 traversing from root to the leaf nodes (*Ccel-Cjos-Cbnl-Cpap*); (iii) a similar pattern (*i.e.*,  
263 increase in number) is observed in CBMs (54, 59, 67, 71) and to a less degree, dockerins (69,  
264 72, 88, 68) (Dassa *et al.*, 2017) (**Fig. 4A**). Thus *Ccel* is an outlier in terms of the genome-



265 wide CAZyme number. Moreover, of 26Kb in size, the cellulosomal operon of *Ccel* is the  
266 largest (**Fig. S5**; *Cjos*: 22.5Kb; *Cbnl*: 25Kb; *Cpap*: 25Kb), and harbors unique genes such as  
267 pectin degrading enzymes (Pagès *et al.*, 2003, McDonald *et al.*, 2008) (*RglIIY*) and longer  
268 hybrid linkers (Pinheiro *et al.*, 2008) that join cohesins to scaffoldins. Similarly, in the *Cace*-  
269 *Cfel*-*Cros* cluster, the *Cace* cellulosomal operon harbors one additional enzyme (Sialidase;  
270 **Fig. 4B**) yet lacks cellulosomal complex activity (Sabathé *et al.*, 2002), in opposite to *Cros*  
271 and *Cfel* which are used for the retting process (Angelini *et al.*, 2013). These observations are  
272 consistent with CoSLOE-derived positive selection of the *Ccel* cellulosomal operon.

273 Similarly, in Subclade 1.1, the cellulosomal operon of *Chun* uniquely harbors a  
274 xyloglucanase gene. Thus the near-root cellulosome operons are positively selected towards  
275 less efficient cellulosic activity or addition of auxiliary functionality, supporting CoSLOE-  
276 derived operon selection.

### 277 **3.6 Evolution of the ATP synthase operons via CoSLOE is similar to the cellulosomal** 278 **operons**

279 CoSL-tree of the ATP synthase operons is similar to that of the cellulosome operons (**Fig.**  
280 **5A**; **Fig. 1C**), except that *Cpun* is clustered with *Chun* in the former. Notably, within each of  
281 the *Chun*-*Cell*-*Cter*, *Ccel*-*Cpap*-*Cbnl*-*Cjos* and *Cros*-*Cace*-*Cfel* subclades, operon sequences  
282 are nearly 100% similar, and the gene sequences of subunit alpha and beta are conserved  
283 across 13 species (**Fig. 5A**). The less variation in gene sequences (than cellulosome operon)  
284 among 13 *Clostridium* species is probably due to the strict functional conservation of the  
285 ATP synthase complex.

286 As in cellulosome operons, gene-SL relationship was probed in the ATPase operons.  
287 Three orthologous SLs were predicted in ATPase operon, where (i) SL-1 is preserved in  
288 *Chun*, *Cjos*, *Cbnl*, *Cpap* and *Ccel* (always flanking at 3' UTR of subunit C); (ii) SL-2 is  
289 present at 3' UTR of subunit alpha in *Cpun*, *Cjos*, *Cbnl*, *Cpap*, *Ccel*, and *Csac*; (iii) SL-3 is  
290 conserved throughout the 13 species at the 3' UTR of epsilon chain and terminating the  
291 operon; (iv) Only one SL is predicted in *Cace*, *Cfel*, *Cros*, and *Cloc*, suggesting that their  
292 ATP synthase operons seem not regulated by the SRPS mechanism (**Fig. 5B**). Taken together,  
293 these associations between genes and SLs show their relationship, which is consistent with  
294 the observation in cellulosome operons.

295 As for ER, in Clade 1, the ERs for *Cter*-*Cpun*, *Ccell*-*Cpun* and *Chun*-*Cpun* are 0.34, 0.41  
296 and 0.80 respectively (**Table S3**), revealing negative selection towards leaf nodes in Clade 1  
297 (**Fig. 5A**), *i.e.* *Ccell* and *Cter* appear to undergo purifying selection, while *Chun* is positively  
298 selected towards *Cpun* (ER for *Cpun*-*Chun*: 1.24; **Table S3**). In Subclade 2.1, similarly, ERs  
299 of *Cjos*-*Ccel*, *Cbnl*-*Ccel* and *Cpap*-*Ccel* are 0.94, 0.93, 1.20, respectively (**Table S3**), *i.e.* the  
300 overall flow of *Cjos*-*Cbnl*-*Ccel* is consistent with positive selection, except for *Cpap* (**Fig.**  
301 **5A**).

302 Interestingly, the ATP synthase operons exhibit an evolution pattern similar to the  
303 cellulosome operons, by positive selection in the *Cter*-*Ccell*-*Chun* and *Cjos*-*Cbnl*-*Cpap*-*Ccel*,  
304 direction (purifying selection in the reverse direction; **Fig. 5A**). Since ATP synthase operon is

305 functionally conserved in most species (Neupane *et al.*, 2019), less variability was present in  
306 the genes and SLs. However, the root species of Subclade 1.1 (*Chun*), 1.2 (*Ccel*) and 2.1  
307 (*Cloc*) harbors smaller operons, longer operons and an additional enzyme at the 3' UTR  
308 region, respectively (**Fig. 5A**), consistent with the evolutionary pattern suggested by the  
309 observations in cellulosome operons.

### 310 **3.7 Genome-wide application of CoSLOE reveals the direction of organismal selection**

311 The evolutionary flow in a tree represents the different directions that species follow due  
312 to the selection-pressure on them, during evolution. To probe SL-driven evolutionary  
313 selection-pressure, orthologous SRPS operons were probed using CoSLOE. However, due to  
314 the lack of orthology among SRPS operons in 13 species, operon evolution was probed clade-  
315 wise in CoSL-tree (**Fig. 1C**), *i.e.* Subclade 1.2 (*Ccel-Cjos-Cbnl-Cpap*). In Subclade 1.2,  
316 orthologous SRPS operons are scattered across the genomes of *Ccel*, *Cjos*, *Cbnl* and *Cpap*  
317 which are in the form of chromosome (*Cbnl* and *Ccel*) or contigs (*Cpap*-31 and *Cjos*-2; **Fig.**  
318 **6A**; **Table S4**).

319 Here, five out of the 25 SRPS operons, *i.e.*, Op617, Op622, Op716, Op863 and Op1745,  
320 follow the *Ccel-Cjos-Cbnl-Cpap* direction (black arrows; **Fig. 6A**). Interestingly, the other 80%  
321 SRPS operons show the positive selection flow in the *Cpap-Cbnl-Cjos-Ccel* direction, *i.e.*  
322 from the leaf nodes to the root nodes, which is similar to the cellulosomal cluster  
323 evolutionary flow, e.g., for Op142, Op376 and Op898 (red arrows; **Fig. 6A**). These  
324 observations suggest that the SRPS mechanism, although evolutionarily conserved, can  
325 reveal selection-pressure that is distinct from organismal phylogeny.

## 326 **4 Discussion and conclusion**

327 In existing frameworks of operon evolution, coding sequences (*i.e.*, subunits of protein  
328 complex or components of metabolic pathway encoded by the operon) have been thought to  
329 play a major role. They can drive structural variation and functional adaptation of operons  
330 towards a specific niche (Gogarten *et al.*, 2002, Francino *et al.*, 2012), for example, by  
331 deletion or insertion of the whole genes or via synonymous/non-synonymous mutations of  
332 their sequences. However, it remains elusive whether non-coding elements play a role in such  
333 adaptation of operons.

334 The stoichiometry of SRPS operons, found genome-wide, can be modeled based on the  
335 genome sequence of SLs alone (Bhaskar *et al.*, 2021), suggesting a quantitative model of  
336 evolution at the whole operon level, in parallel to the evolution at the individual coding gene  
337 level (e.g.,  $K_a/K_s$  (Kimura *et al.*, 1968)). Based on  $\Delta G$  of such SLs, we proposed CoSLOE,  
338 with evolutionary ratio (ER)  $>1$  or  $<1$  for positive or negative selection of SRPS operons. In  
339 the CoSL-tree of cellulosomal operons, when traversing from leafs to the root nodes, ERs  
340 reveal diversifying/positive selection towards a less efficient cellulosomal system, consistent  
341 with glycoside-hydrolase gene variation both in-operon and genome-wide. A consistent  
342 pattern is followed by the ATPase operon and the majority of orthologous SRPS operons  
343 genome-wide, suggesting conservation among operons in such selection. Therefore, CoSLOE  
344 provides a new layer of insights into operon evolution that is distinct from existing models

345 **(Fig. 6B)**. Specifically, (i) Driving forces: for individual genes, mutation, recombination,  
346 genetic drift and selection are known evolution drivers; for SRPS operons, in addition to  
347 addition/deletion/mutation of genes, CoSLOE introduces SLs a previously unrecognized  
348 driver. (ii) Theoretical models: in addition to the known models of gene evolution (neutral  
349 theory (Nei *et al.*, 2005)) and operon evolution (selfish operon theory (Lawrence & Roth *et*  
350 *al.*, 1996), co-regulation model (Price *et al.*, 2005) and piece-wise model (Fani *et al.*, 2005)),  
351 CoSLOE provides a new framework for quantitatively modeling SRPS operon evolution. (iii)  
352 Rate of selection: CoSLOE compares rate of selection between two orthologous operons,  
353 which is conceptually similar to Ka/Ks or dN/dS which compares between orthologous genes.  
354 (iv) Outcome of selection: gene evolution generally results in changed protein sequence, yet  
355 the operon evolution depicted by CoSLOE results in altered ingredient or stoichiometry of  
356 the whole protein complex or metabolic pathway. (v) Direction of selection: just like Ka/Ks  
357 for orthologous genes, CoSLOE offers strategy to model the direction for orthologous  
358 operons. (vi) Phylogeny: a tree based on  $\Delta G$  of SLs of SRPS operons can model the selection  
359 of operon and organism, in contrast to 16S-rRNA gene trees that model organismal  
360 phylogeny. (vii) Underlying sequence: instead of relying on coding sequences, CoSLOE  
361 takes advantage of the non-coding sequences to model operon evolution, and highlights the  
362 role of *cis*-elements in shaping operon evolution and organism adaptation. (viii) Origin of  
363 operons: CoSLOE suggests the SLs (and their relationship with associated genes) as a key  
364 player in the original formation of operon structure, in addition to horizontal gene transfer,  
365 deletion of intervening genes and addition of ORFan genes (Price *et al.*, 2006).

366 Notably, we have tested CoSLOE on just 13 Clostridial species, and expansion of the  
367 model to a broader range of species is limited by the paucity of experimental data and lack of  
368 computational approaches to identify SPRS operons. Therefore, to what degree the model is  
369 applicable across microorganisms is not yet clear, and answer to this question is perhaps  
370 ultimately dependent on the breath and boundary of SPRS mechanism. Despite these  
371 limitations, for SRPS operons, our findings here reveal the link between operon stoichiometry  
372 and operon evolution, and propose a new *cis*-element-based framework to model the  
373 direction and rate of SRPS operon evolution.

## 374 **5 Acknowledgements**

375 This work was supported by University of Chinese Academy of Sciences Scholarship for  
376 International PhD students.

## 377 **6 Author contribution**

378 YB and JX designed the study; YB performed the computational analysis; YB and JX  
379 analyzed the data; MHB and CX provided critical suggestions; YB and JX wrote the paper.

## 380 **7 Competing interests**

381 The authors declare no conflicts of interest.

## 382 **8 Data availability**

383       The data underlying this article are available in the article and in its online  
384   supplementary material.

385 **9 References**

- 386 Altschul SF, Gish W, Miller W, Myers EW & Lipman DJ (1990) Basic local alignment  
387 search tool. *J Mol Biol* **215**: 403-410.
- 388 Angelini LG, Tavarini S & Foschi L (2013) Spanish broom (*Spartium junceum* L.) as new  
389 fiber for biocomposites: the effect of crop age and microbial retting on fiber quality. Vol.  
390 2013 p.^pp. Hindawi Publishing Corporation.
- 391 Bhaskar Y, Su X, Xu C & Xu J (2021) Predicting Selective RNA Processing and  
392 Stabilization Operons in *Clostridium* spp. *Front Microbiol* **12**: 1281.
- 393 Busch A, Kunert G, Heckel DG & Pauchet Y (2017) Evolution and functional  
394 characterization of CAZymes belonging to subfamily 10 of glycoside hydrolase family 5  
395 (GH5\_10) in two species of phytophagous beetles. *PloS one* **12**: e0184305.
- 396 Dassa B, Borovok I, Lombard V, Henrissat B, Lamed R, Bayer EA & Moraïs S (2017) Pan-  
397 Cellulosomics of Mesophilic Clostridia: Variations on a Theme. *Microorganisms* **5**: 74.
- 398 Doi RH, Park J-S, Liu C-c, Malburg Jr LM, Tamaru Y, Ichiishi A & Ibrahim A (1998)  
399 Cellulosome and noncellulosomal cellulases of *Clostridium cellulovorans*. *Extremophiles* **2**:  
400 53-60.
- 401 Fani R, Brillì M & Lio P (2005) The origin and evolution of operons: the piecwise building  
402 of the proteobacterial histidine operon. *J Mol Evol* **60**: 378-390.
- 403 Francino MP (2012) The ecology of bacterial genes and the survival of the new. *Int J Evol*  
404 *Biol* **2012**.
- 405 Gogarten JP, Doolittle WF & Lawrence JG (2002) Prokaryotic evolution in light of gene  
406 transfer. *Mol Biol Evol* **19**: 2226-2238.
- 407
- 408 Hofacker IL (2003) Vienna RNA secondary structure server. *Nucleic Acids Res* **31**: 3429-  
409 3431.
- 410 Kimura M (1968) Evolutionary rate at the molecular level. *Nature* **217**: 624-626.
- 411 Klug G (1993) The role of mRNA degradation in the regulated expression of bacterial  
412 photosynthesis genes. *Mol Microbiol* **9**: 1-7.
- 413 Kumar S, Stecher G & Tamura K (2016) MEGA7: Molecular Evolutionary Genetics  
414 Analysis version 7.0 for bigger datasets. *Mol Biol Evol* **33**: 1870-1874.

- 415 Lawrence JG & Roth JR (1996) Selfish operons: horizontal transfer may drive the evolution  
416 of gene clusters. *Genetics* **143**: 1843-1860.
- 417 Ludwig H, Homuth G, Schmalisch M, Dyka FM, Hecker M & Stülke J (2001) Transcription  
418 of glycolytic genes and operons in *Bacillus subtilis*: evidence for the presence of multiple  
419 levels of control of the gapA operon. *Mol Microbiol* **41**: 409-422.
- 420 Macke TJ, Ecker DJ, Gutell RR, Gautheret D, Case DA & Sampath R (2001) RNAMotif, an  
421 RNA secondary structure definition and search algorithm. *Nucleic Acids Res* **29**: 4724-4735.
- 422 McDonald AG, Boyce S & Tipton KF (2008) ExplorEnz: the primary source of the IUBMB  
423 enzyme list. *Nucleic Acids Res* **37**: D593-D597.
- 424 Nei M (2005) Selectionism and neutralism in molecular evolution. *Mol Biol Evol* **22**: 2318-  
425 2342.
- 426 Neupane P, Bhujra S, Thapa N & Bhattarai HK (2019) ATP Synthase: Structure, Function and  
427 Inhibition. Vol. 10 p.^pp. 1-10.
- 428 Newbury SF, Smith NH & Higgins CF (1987) Differential mRNA stability controls relative  
429 gene expression within a polycistronic operon. *Cell* **51**: 1131–1143.
- 430 Osbourn AE & Field B (2009) Operons. *Cell Mol Life Sci* **66**: 3755-3775.  
431
- 432 Pagès S, Valette O, Abdou L, Bélaïch A & Bélaïch J-P (2003) A rhamnogalacturonan lyase  
433 in the *Clostridium cellulolyticum* cellulosome. *J Bacteriol* **185**: 4727-4733.
- 434 Pinheiro BA, Proctor MR, Martinez-Fleites C, Prates JA, Money VA, Davies GJ, Bayer EA,  
435 FontesM CM, Fierobe H-P & Gilbert HJ (2008) The *Clostridium cellulolyticum* dockerin  
436 displays a dual binding mode for its cohesin partner. *J Biol Chem* **283**: 18422-18430.
- 437 Price MN, Arkin AP & Alm EJ (2006) The life-cycle of operons. *PLoS Genet* **2**: e96.
- 438 Price MN, Huang KH, Arkin AP & Alm EJ (2005) Operon formation is driven by co-  
439 regulation and not by horizontal gene transfer. *Genome Res* **15**: 809-819.
- 440 Rochat T, Bouloc P & Repoila F (2013) Gene expression control by selective RNA  
441 processing and stabilization in bacteria. *FEMS Microbiol Lett* **344**: 104–113.
- 442 Sabathé F, Bélaïch A & Soucaille P (2002) Characterization of the cellulolytic complex  
443 (cellulosome) of *Clostridium acetobutylicum*. *FEMS Microbiol Lett* **217**: 15-22.

- 444 Smith C, Heyne S, Richter AS, Will S & Backofen R (2010) Freiburg RNA Tools: a web  
445 server integrating INTARNA, EXPARNA and LOCARNA. *Nucleic Acids Res* **38**: W373-  
446 W377.
- 447 Suzuki R & Shimodaira H (2006) Pvclust: an R package for assessing the uncertainty in  
448 hierarchical clustering. *Bioinformatics* **22**: 1540-1542.
- 449 Will S, Joshi T, Hofacker IL, Stadler PF & Backofen R (2012) LocARNA-P: accurate  
450 boundary prediction and improved detection of structural RNAs. *RNA* **18**: 900-914.
- 451 Xu C, Huang R, Teng L, Jing X, Hu J, Cui G, Wang Y, Cui Q & Xu J (2015) Cellulosome  
452 stoichiometry in *Clostridium cellulolyticum* is regulated by selective RNA processing and  
453 stabilization. *Nat Commun* **6**: 6900.
- 454 Yang Z (2007) PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**:  
455 1586-1591.  
456

## 457 10 Figure Legends

458 **Figure 1. Composite SLs (CoSLs) in the cellulosomal operons from 13 Clostridial**  
459 **species. (A)** Schematic representation of the SRPS operon via  $\Delta G$  of the harbored composite  
460 SLs. Upstream Controlled Unit (UCU) represents the region (which is upstream to a SL and  
461 can include multiple genes) that is regulated by a SL via the SRPS mechanism. **(B)**  $\Delta G$  of the  
462 harbored SLs in the cellulosomal operons from 13 *Clostridium* spp., showing skewness of the  
463  $\Delta G$  within an operon and divergence of pattern among orthologous operons. **(C)** Composite  
464 SLs (CoSLs)-based tree of the cellulosomal operons using the orthologous SLs. Genes are  
465 colored based on the encoded protein.

466

467 **Figure 2. Structural alignment of SLs harbored in the cellulosome operons from 13**  
468 **Clostridial species reveal the relationship between SLs and their associated genes.** The  
469 multiple sequence alignment is shown for the five orthologous SLs: SL-1 **(A)**, SL-2 **(B)**, SL-3  
470 **(C)**, SL-5 **(D)** and SL-7 **(E)**. SLs were aligned via sequence or structure using *LocARNA* and  
471 shown with the consensus structure in the Dot bracket form (middle). Compatible base pairs  
472 are colored based on the standard format, where the hue shows sequence conservation among  
473 the number of different types of compatible base pairs (C-G, G-C, A-U, U-A, G-U or U-G) in  
474 the corresponding columns. Color saturation decreases with the number of incompatible base  
475 pairs. The bar plot represents conservation of compatible base pairs (higher bar for higher  
476 conservation, and vice versa).

477

478 **Figure 3. Proposing CoSLOE to quantitatively model the evolution of SRPS operons.** In  
479 CoSLOE (Composite Stem-Loop based Operon Evolution), operon ER is calculated based on  
480 number of genes, number of SLs and coefficient of the variation (CV) of the  $\Delta G$ -based (i.e.,  
481  $\Delta G$  of the CoSLs) ratio.  $\Delta G_{SD}$  and  $\overline{\Delta G}$  represent the standard deviation and the mean,  
482 respectively. The ratio determines the selection pressure between two operons, i.e., ratio of 1  
483 represents the neutral selection, while positive (or negative) selection occurs when ratio is  $>1$   
484 (or  $<1$ ). These ratios in a clade of a tree determine the direction of the evolution, i.e. positive  
485 selection for species with ratio  $>1$  and negative/purifying selection for species with ratio  $<1$ .

486

487 **Figure 4. Evolution of the cellulosome-encoding SRPS operons based on CoSLOE.** The  
488 two clades in the tree represent two different operon evolution scenarios yet with an identical  
489 evolutionary directional flow, i.e., the movement from root to leaf nodes defines purifying  
490 (i.e., negative) selection and the reverse movement (i.e., leaf to root) depicts diversifying  
491 selection (i.e., positive selection). Upstream Controlled Unit (UCU) represents the region  
492 (which is upstream to a SL and can include multiple genes) that is regulated by a SL via the  
493 SRPS mechanism. **(A)** In Clade 1 (with four species), positive selection resulted in distinct  
494  $\Delta G$  of SLs (each corresponding to a UCU) in the outermost operon of the clade, while those  
495 operons with similar  $\Delta G$  are conserved in the leaf nodes. Variation in a UCU can be caused  
496 by gain/loss of the SLs along with the corresponding genes (i.e., depicting appearance and  
497 disappearance of new genes). **(B)** In Clade 2 (with three species), positive selection also led  
498 to distinct  $\Delta G$  of SLs (each corresponding to a UCU) in the outermost species, suggesting  
499 identical evolutionary flow in both clades. However, the positively selected operons carry  
500 characteristics distinct from the leaf-node operons, despite an identical number of SLs and



501 genes. Thus the change in the  $\Delta G$  of SLs can lead to operons with discrete function. Color  
502 gradient of genes represents positive (darker color) or purifying (light color) selection.

503

504 **Figure 5. Evolution of the ATP-synthase-encoding SRPS operons in the 13 Clostridial**  
505 **species based on the CoSLOE model.** (A) ATP synthase operons sequences were mapped  
506 sequence-wise according to the CoSL-based phylogeny, where the sequence similarity is  
507 shown by dark black (100% similarity) and light black (64% similarity) color gradient and  
508 the genes are colored based on their encoding protein. The ERs determine the positive and  
509 negative selection pattern (shown via green and red arrows, respectively). (B) The multiple  
510 sequence alignment of SL-3 from the ATP synthase operons in the 13 Clostridial species.

511

512 **Figure 6. The genome-wide CoSLOE model defines the direction of organismal selection.**  
513 (A) Genome-wide evolution of SRPS operons based on the CoSLOE model (Table S4).  
514 Orthologous SRPS operons are scattered across the genomes of *C. cellulolyticum* (*Ccel*), *C.*  
515 *josui* (*Cjos*), *C. sp. BNL1100* (*Cbnl*), and *C. papyrosolvans* (*Cpap*), represented here in the  
516 form of chromosome (*Cbnl* and *Ccel*) and contigs (*Cpap-31* and *Cjos-2*). Totally, five out of  
517 25 SRPS operons follow the *Ccel-Cjos-Cbnl-Cpap* direction and the other 80% SRPS  
518 operons show a positive selection flow in the *Cpap-Cbnl-Cjos-Ccel* direction, *i.e.*, from the  
519 leaf nodes to the root nodes. Dark and grey region/band represents operon (Op), and their  
520 thickness shows the length. Direction of arrows represents direction of selection pressure,  
521 either positive (red) or negative (black). (B) Link and distinction between CoSLOE and the  
522 existing models for operon evolution. Information derived from the CoSLOE model is  
523 highlighted in red.

524 **11 Table Legends**

525 **Table 1. Evolutionary ratio (ER) matrix for the cellulosome complex operons from the**  
 526 **13 Clostridial species. (A)** ERs for Clade-1 of CoSL-tree, where the subclades of *Cpap-*  
 527 *Cbnl-Cjos-Ccel*, *Ccell-Cter-Chun* and *Cros-Cfel-Cace* are colored with orange, yellow and  
 528 green respectively. **(B)** ERs for Clade-2 of CoSL-tree.

529 **(A)**

	<i>Cpap</i>	<i>Cbnl</i>	<i>Cjos</i>	<i>Ccel</i>	<i>Ccell</i>	<i>Cter</i>	<i>Chun</i>	<i>Cros</i>	<i>Cfel</i>	<i>Cace</i>
<i>Cpap</i>	1	0.99	1.75	0.91	2.50	2.28	1.37	10.21	9.72	1.32
<i>Cbnl</i>	1.01	1	1.76	0.91	2.52	2.30	1.38	10.29	9.79	1.33
<i>Cjos</i>	0.57	0.57	1	0.52	1.43	1.31	0.78	5.84	5.56	0.75
<i>Ccel</i>	1.10	1.10	1.93	1	2.77	2.52	1.51	11.27	10.73	1.45
<i>Ccell</i>	0.40	0.40	0.70	0.36	1	0.91	0.55	4.08	3.88	0.53
<i>Cter</i>	0.44	0.43	0.77	0.40	1.10	1	0.60	4.47	4.26	0.58
<i>Chun</i>	0.73	0.73	1.28	0.66	1.83	1.67	1	7.46	7.10	0.96
<i>Cros</i>	0.10	0.10	0.17	0.09	0.25	0.22	0.13	1	0.95	0.13
<i>Cfel</i>	0.10	0.10	0.18	0.09	0.26	0.23	0.14	1.05	1	0.14
<i>Cace</i>	0.76	0.75	1.33	0.69	1.90	1.73	1.04	7.76	7.39	1

530

531 **(B)**

	<i>Cloc</i>	<i>Csac</i>	<i>Cpun</i>
<i>Cloc</i>	1	2.12	1.38
<i>Csac</i>	0.47	1	0.65
<i>Cpun</i>	0.72	1.54	1

## 532 12 Supplementary Tables and Figures

533 **Table S1. Free energy ( $\Delta G$ ) of harbored SLs in the SRPS operons that encode cellulosome from**  
534 **13 Clostridial species.**

535

536 **Table S2.  $K_a/K_s$  values for the first gene in the cellulosomal operon from *Cpap*, *Cbnl*, *Cjos* and**  
537 ***Ccel*.**

538

539 **Table S3. The evolutionary ratio (ER) matrix for the ATP synthase operons from the 13**  
540 **Clostridial species.**

541

542 **Table S4. The number of genes and SLs in all SRPS operons in the *Ccel*, *Cjos*, *Cbnl* and *Cpap***  
543 **genomes.**

544

545 **Figure S1. Phylogenetic tree of 13 Clostridial species based on the predicted  $\Delta G$  of the SLs**  
546 **(CoSL-tree) (A) or the 16S rRNA sequences (16S-tree) (B).**

547

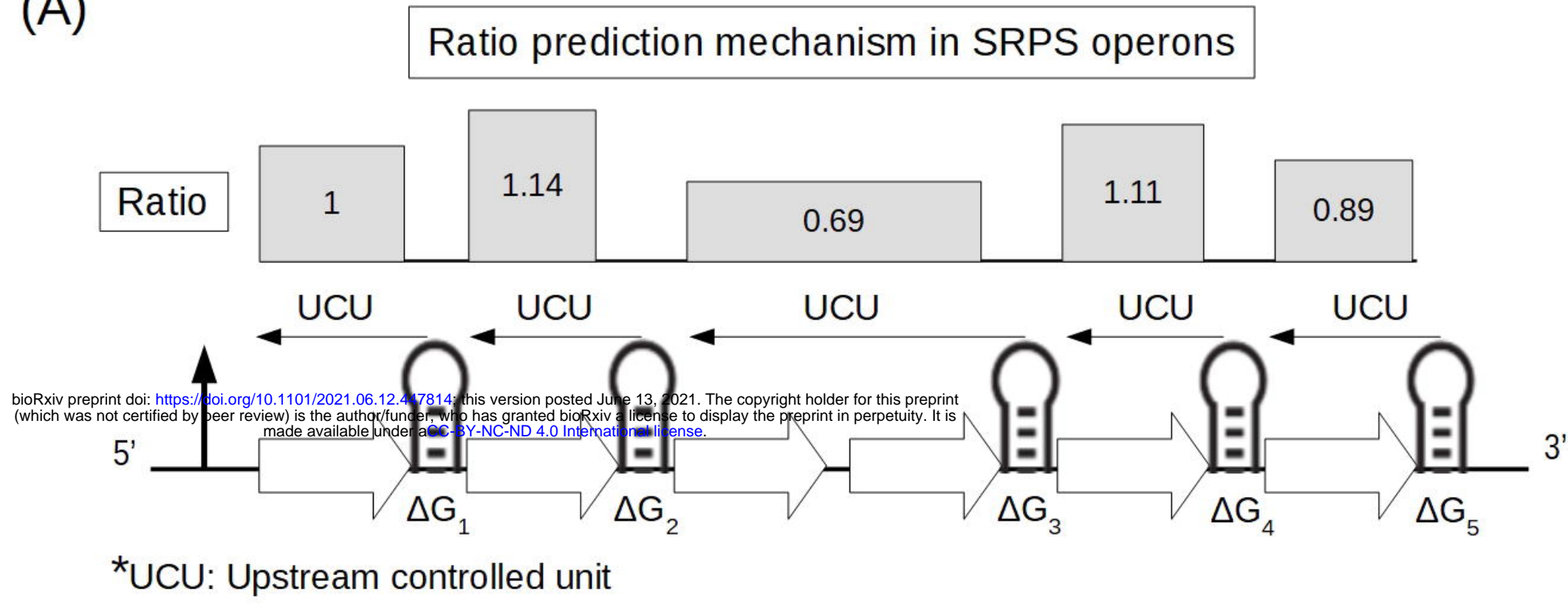
548 **Figure S2. Clade-wise representation of the SL-1 (A), SL-2 (B), SL-3 (C), SL-4 (D), SL-5 (E),**  
549 **SL-6 (F), SL-7 (G), SL-3A (H), SL-2A (I) and SL-7A (J), based on their sequence and structure**  
550 **similarity.**

551

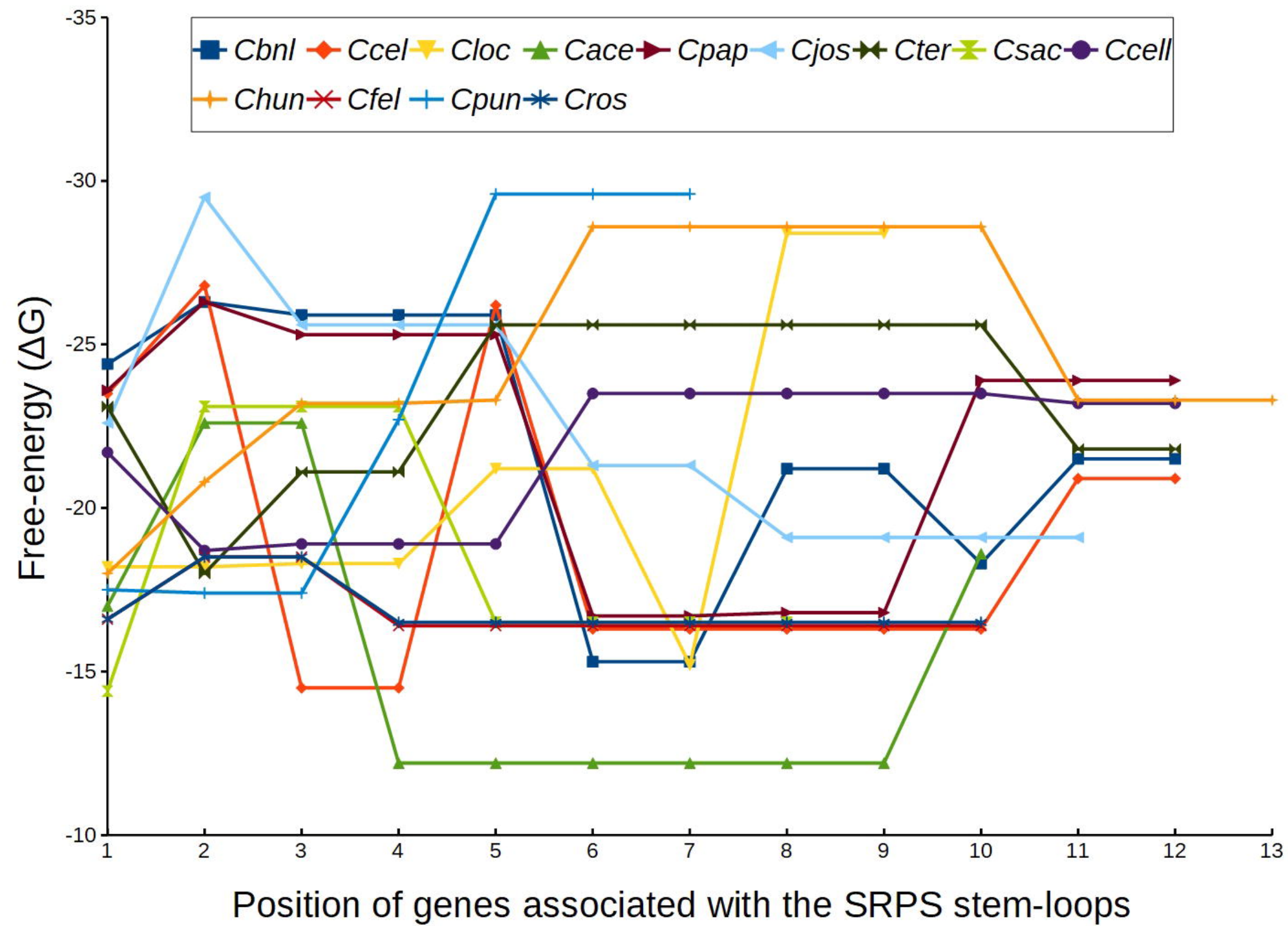
552 **Figure S3. Linear comparison of the nucleic acid sequences of the cellulosome operons from**  
553 **the 13 Clostridial species.**

Figure 1

(A)



(B)



(C)

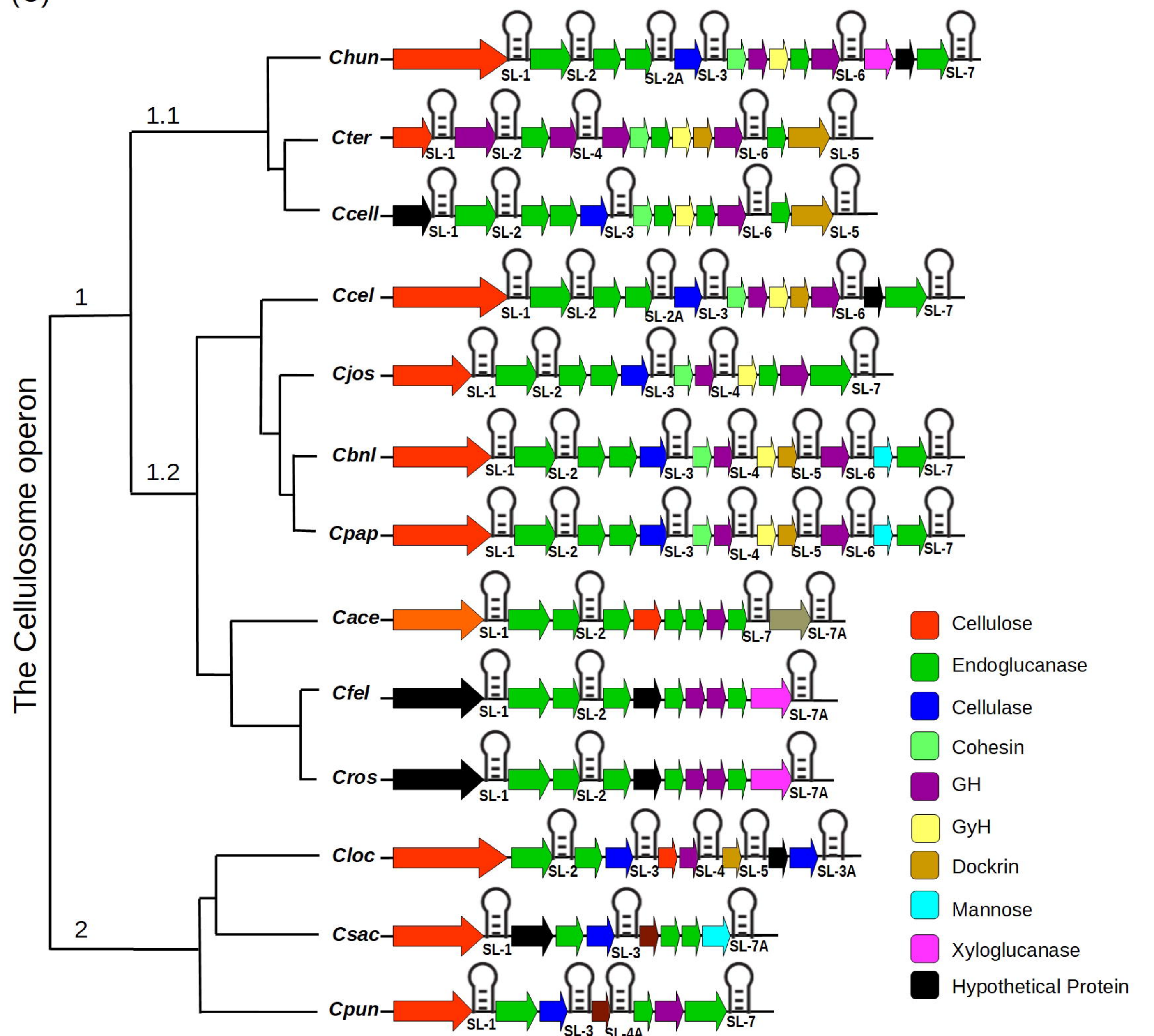




Figure 3

Composite Stem-Loop based Operon Evolution (CoSLOE) model

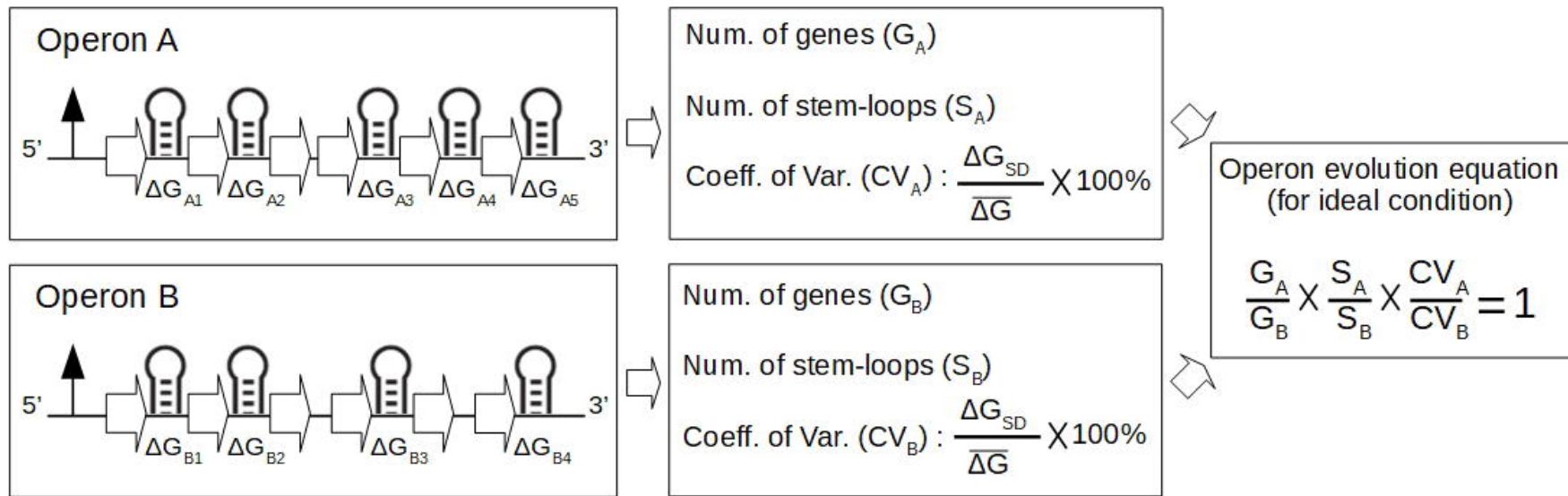


Figure 4

