

# A TetR-family protein activates transcription from a new promoter motif associated with essential genes for autotrophic growth in acetogens

## Authors

Renato de Souza Pinto Lemgruber<sup>1</sup>, Kaspar Valgepea<sup>1,2</sup>, Ricardo Axayacatl Gonzalez Garcia<sup>1</sup>, Christopher de Bakker<sup>1,□</sup>, Robin William Palfreyman<sup>1,3</sup>, Ryan Tappel<sup>4</sup>, Michael Köpke<sup>4</sup>, Séan Dennis Simpson<sup>4</sup>, Lars Keld Nielsen<sup>1</sup>, Esteban Marcellin<sup>1,3\*</sup>

## Affiliations

<sup>1</sup>Australian Institute for Bioengineering and Nanotechnology (AIBN), The University of Queensland, Brisbane, Queensland, Australia;

<sup>2</sup>ERA Chair in Gas Fermentation Technologies, Institute of Technology, University of Tartu, Tartu, Estonia;

<sup>3</sup>Queensland Node of Metabolomics Australia, The University of Queensland, Brisbane, Queensland, Australia;

<sup>4</sup>LanzaTech Inc., Skokie, Illinois, USA

\* **Correspondence:** Dr. Esteban Marcellin, e.marcellin@uq.edu.au

Present address: □Servatus Ltd., Innovation Centre, University of the Sunshine Coast, Sippy Downs, Australia

**Words: 7941**

**Figures: 4**

**Keywords:** *Clostridium autoethanogenum*; Wood-Ljungdahl pathway; Transcriptional regulation; Gas fermentation; Autotrophy

## Abstract

Acetogens can fix carbon (CO or CO<sub>2</sub>) into acetyl-CoA via the Wood-Ljungdahl pathway (WLP) that also makes them attractive cell factories for the production of fuels and chemicals from waste feedstocks. Although most biochemical details of the WLP are well understood and systems-level characterisation of acetogen metabolism has recently improved, key transcriptional features such as promoter motifs and transcriptional regulators are still unknown in acetogens. Here, we use differential RNA-sequencing to identify a previously undescribed promoter motif associated with essential genes for autotrophic growth of the model-acetogen *Clostridium autoethanogenum*. RNA polymerase was shown to bind to the new promoter motif using a DNA-binding protein assay and proteomics enabled the discovery of four candidates to potentially function directly in control of transcription of the WLP and other key genes of C<sub>1</sub> fixation metabolism. Next, *in vivo* experiments showed that a TetR-family transcriptional regulator (CAETHG\_0459) and the housekeeping sigma factor ( $\sigma^A$ ) activate expression of a reporter protein (GFP) in-frame with the new promoter motif from a fusion vector in *E. coli*. Lastly, a protein-protein interaction assay with the RNA polymerase (RNAP) shows that CAETHG\_0459 directly binds to the RNAP. Together, the data presented here advance the fundamental

46 understanding of transcriptional regulation of C<sub>1</sub> fixation in acetogens and provide a strategy for  
47 improving the performance of gas-fermenting bacteria by genetic engineering.

48  
49  
50

## 51 **1. Introduction**

52 The Wood-Ljungdahl pathway (WLP) of acetogens is speculated to be the first biochemical  
53 pathway on Earth that emerged when the atmosphere was still highly reduced and rich in CO,  
54 CO<sub>2</sub>, and H<sub>2</sub> (Fuchs, 2011; Russell and Martin, 2004; Weiss et al., 2016). These C<sub>1</sub> gases can be  
55 converted into acetyl-CoA through the WLP (Ragsdale and Pierce, 2008; Wood, 1991) and  
56 acetogens are the only known organisms using the WLP as a terminal electron-accepting,  
57 energy-conserving process to fix CO<sub>2</sub> into biomass (Drake et al., 2006; Fuchs, 2011). This  
58 pathway is responsible for the production of acetic acid in quantities surpassing the billion ton  
59 mark annually. It is estimated that the pathway contributes to fixing ~20% of the CO<sub>2</sub> on Earth  
60 (Drake et al., 2006; Ljungdahl, 2009). All this takes place with the WLP operating at the edge of  
61 thermodynamic feasibility (Schuchmann and Müller, 2014) and requires the use of the third  
62 mode of energy conservation, electron bifurcation, which likely contributed to the emergence of  
63 life on Earth (Herrmann et al., 2008; Li et al., 2008; Nitschke and Russell, 2011). Acetogens are  
64 also attractive cell factories for the sustainable production of fuels and chemicals from gaseous  
65 waste feedstocks (e.g. syngas from gasified municipal solid waste and industrial waste gases)  
66 (Claassens et al., 2016; Dürre and Eikmanns, 2015; Liew et al., 2016; Molitor et al., 2016).  
67 While the field has advanced enormously in the last decade (Liew et al., 2016; Molitor et al.,  
68 2016), better fundamental understanding of acetogen metabolism is needed to guide rationale  
69 metabolic engineering, for example, to increase their substrate uptake or product yields.

70

71 Recent quantitative studies of acetogen physiology have expanded understanding of their  
72 metabolism considerably (reviewed in (Molitor et al., 2017; Schuchmann and Müller, 2014)).  
73 Although most biochemical details of the WLP are well established (Ragsdale, 1991, 1997,  
74 2008) and systems-level understanding of acetogen metabolism has recently improved (Valgepea  
75 et al., 2017a, 2018), key transcriptional features such as promoter motifs and transcriptional  
76 regulators controlling the expression of genes needed for autotrophic growth are yet unknown.  
77 This information could benefit acetogen metabolic engineering and improve our understanding  
78 of their complex transcriptional regulation (Aklujkar et al., 2017; Marcellin et al., 2016;  
79 Nagarajan et al., 2013; Tan et al., 2013). Prediction of promoter motifs strictly based on  
80 computational analysis (based solely on the organism's genome sequence) has the drawback of  
81 detection of promoter-like sequences across the genome, which is particularly pronounced in  
82 non-conserved DNA motifs (Patrik, 2006). An instrumental step towards more accurate promoter  
83 motif identification was the development of the differential RNA-sequencing (dRNA-Seq)  
84 technology, first described in 2010 by Sharma and colleagues (Sharma et al., 2010) for the  
85 human pathogen *Helicobacter pylori*.

86

87 dRNA-Seq enables the experimental determination of transcription start sites (TSSs) and correct  
88 mapping of TSSs enables genome-wide identification of promoters and gene expression  
89 regulatory sequences, besides providing experimental data for a more accurate genome  
90 annotation. Once a TSS has been experimentally determined, promoter sequences can be mapped  
91 from there. Thus, characterisation of the transcriptional architecture (i.e. TSSs and promoter

92 motifs) and a more accurate annotation of acetogen genomes have the potential to yield valuable  
93 insights into the complex transcriptional regulation of acetogens. To date, only one study has  
94 determined TSSs in acetogens, using *Eubacterium limosum* (Song et al., 2017). Here, we used  
95 dRNA-Seq as a tool to identify the TSSs in the model-acetogen *Clostridium autoethanogenum*  
96 grown under autotrophic and heterotrophic conditions. The subsequent search for promoter  
97 motifs detected a previously undescribed motif associated with essential genes in acetogens. We  
98 then provide experimental evidence for the relevance of this new promoter motif (names  
99 hereafter  $P_{\text{cauto}}$ ) by identifying a TetR-family protein that activates gene expression from this  
100 motif by directly binding to the RNA polymerase.

101

## 102 **2. Materials and Methods**

### 103 **2.1 Bacterial strains and growth conditions**

104 *Clostridium autoethanogenum* strain DSM 10061 was obtained from The German  
105 Collection of Microorganisms and Cell Cultures (DSMZ). Cells were grown as described before  
106 (Marcellin et al., 2016). Briefly, heterotrophic and autotrophic growth were investigated in serum  
107 bottles on fructose (5 g/L) and on steel mill off-gas (35% CO, 10% CO<sub>2</sub>, 2% H<sub>2</sub> and 53% N<sub>2</sub>),  
108 respectively. Cells were grown at 37 °C on a shaker (100 RPM, rounds per minute) and sampled  
109 for dRNA-Seq analysis from the exponential growth phase ( $OD_{600\text{nm}} = 0.5-0.6$ ).

110

### 111 **2.2 Differential RNA-sequencing (dRNA-Seq)**

112 Extraction and preparation of RNA for cDNA library construction were performed as described  
113 elsewhere (Marcellin et al., 2016). Briefly, RNA was extracted using TRIzol followed by  
114 column purification with RNAeasy (Qiagen). The resulting total RNA pools were sent to Vertis  
115 Biotechnologie AG (Freisig, Germany) for sequencing. The cDNA libraries were prepared using  
116 the 5'tagRACE method (Fouquier D'Hérouel et al., 2011). Firstly, the 5' Illumina TruSeq  
117 sequencing adapter carrying sequence tag TCGACA was ligated to the 5'-monophosphate groups  
118 (5'P) of processed transcripts (TAP- on Figure 1A). Samples were then treated with Tobacco  
119 Acid Pyrophosphatase (TAP) to convert 5'-triphosphate (5'PPP) structures of primary transcripts  
120 into 5'P ends to which the 5' Illumina TruSeq sequencing adapter carrying sequence tag  
121 GATCGA was ligated (TAP+ on Figure 1A). Next, first-strand cDNA was synthesised using an  
122 N6 randomised primer to which the 3' Illumina TruSeq sequencing adapter was ligated after  
123 fragmentation.

124 The 5' cDNA fragments were amplified with PCR using a proof reading enzyme and  
125 primers designed for TruSeq sequencing according to the manufacturer's instructions. The main  
126 advantage of using the 5'tagRACE method (Fouquier D'Hérouel et al., 2011) for dRNA-Seq  
127 comes from amplifying the 5' ends of processed and primary transcripts in a single PCR reaction,  
128 which preserves their quantitative representation in an RNA pool. Finally, 5' cDNAs were  
129 purified using the Agencourt AMPure XP Kit (Beckman Coulter Genomics) and analysed by  
130 capillary electrophoresis before sequencing the single-end libraries using the Illumina NextSeq  
131 500 system and a MID 150 Kit with 75 bp read length.

132

### 133 **2.3 Determination of transcription start sites (TSSs)**

134 Sequencing reads were aligned and mapped to the genome of *C. autoethanogenum* DSM  
135 10061 (CP006763.1) using the software TopHat2 (Kim et al., 2013) without trimming or  
136 removal of any reads. Reads were processed with the TSSAR (TSS Annotation Regime) software  
137 (Amman et al., 2014) for automated *de novo* determination of TSSs from dRNA-Seq data using

138 the following parameters: p-Value 1e-3, Noise threshold 10, Merge range 5. The identified TSSs  
139 were classified as primary (within 250 nt upstream of an annotated gene), internal (within an  
140 annotated gene), antisense (on the opposite strand of an annotated gene), or orphan (not assigned  
141 to any of the previous classes) (Figure 1B). Since our main aim was the identification of the  
142 TSSs of essential genes for autotrophic growth in acetogens (e.g WLP), we focused on the  
143 primary TSSs.

144

#### 145 **2.4 Search for promoter motifs and the Shine-Dalgarno sequence**

146 To determine promoter motifs, we searched for consensus sequence motifs 50 nt upstream of  
147 primary TSSs using the MEME software (Bailey et al., 2009) with the following parameters: -  
148 dna, -max size 10000000, -mod zoops, -nmotifs 50, -minw 4, -maxw 50, -revcomp, -oc. Only  
149 motifs with E-value  $\leq 0.05$  and at least 13 TSSs associated to it (i.e. at least two genes associated  
150 to it, Figure 1C) were considered and ranked based on the number of assigned TSSs  
151 (Supplementary file 1).

152

153 To search for the Shine-Dalgarno sequence, 30 nt upstream of annotated genes (CP006763.1 and  
154 NC\_022592.1) were searched with the MEME software (Bailey et al., 2009) using the same  
155 parameters as in the promoter motif search, except for -nmotifs 10, -maxw 30.

156

#### 157 **2.5 Search for the new promoter motif in acetogens**

158 Occurrence of the new promoter motif (see results) in *C. autoethanogenum*, *C. ljungdahlii*, *C.*  
159 *ragdalei*, *C. coskatii*, *M. thermoacetica*, and *E. limosum* was determined using the FIMO tool  
160 (Grant et al., 2011) within the MEME software by searching for the sequence up to 300 nt  
161 upstream of annotated genes (since no TSS data is available) with default FIMO parameters.  
162 Occurrence in each acetogen relative to *C. autoethanogenum* was normalised with the number of  
163 annotated genes.

164

#### 165 **2.6 DNA-binding protein assay**

166 Firstly, *C. autoethanogenum*—DSM 19630—cells were acquired from autotrophic bioreactor  
167 chemostat cultures (CO or CO+H<sub>2</sub>) described in a separate work (Valgepea et al., 2018). Briefly,  
168 cells were grown in bioreactor chemostat cultures in the chemically defined medium on either  
169 CO or CO+H<sub>2</sub> at 37 °C, pH = 5, dilution rate of  $\sim 1 \text{ day}^{-1}$  ( $\mu \sim 0.04 \text{ h}^{-1}$ ), and at a biomass  
170 concentration  $\sim 1.4 \text{ gDCW/L}$ . Cells were pelleted by immediate centrifugation ( $20,000 \times g$  for 2  
171 min at 4 °C), and stored at -80 °C until analysis.

172 Frozen pellets were thawed, resuspended in BS/THES buffer described in (Jutras et al.,  
173 2012) with pH adjusted to 7.0, and passed five times through the EmulsiFlex-C5 High Pressure  
174 Homogenizer (Avestin Inc.) according to the manufacturer's instructions, with the final sample  
175 volume adjusted to 35 mL with the BS/THES buffer. Samples were then centrifuged ( $35,000 \times g$   
176 for 15 min at 4 °C) and the supernatant filtered using a 0.22  $\mu\text{M}$  filter (Merck).

177

178 The DNA-binding protein assay was based on a pull-down/DNA affinity chromatography  
179 method described by Jutras and co-workers (Jutras et al., 2012) with the following modifications.  
180 The DNA sequences were of 125 bp length containing the respective promoter sequence in the  
181 middle with flanking regions downstream and upstream. pH of the buffers was adjusted to 7. The  
182 bait-target/ligand binding step was performed with 1 mL of cell extract without the addition of  
183 non-specific competitor DNA.

184 Next, either salmon sperm (Thermo) or Poly dI-dC (Sigma) were used as non-specific  
185 competitor DNA in the subsequent washing steps. Briefly, Dynabeads<sup>TM</sup> M-280 Streptavidin  
186 (Thermo Fisher Scientific) were mixed with DNA containing either the promoter sequence of  
187 CAETHG\_1615, 1617 (WLP genes assigned with the new promoter motif), or 3224 (a glycolytic  
188 gene as a control for our assay since it was assigned the well-known TATAAT motif, which  
189 should yield binding of the RNAP and the housekeeping  $\sigma$  factor  $\sigma^A$ ). Next, the cell extract was  
190 added and samples were incubated for 30 min at room temperature. This was followed by two  
191 washing steps with the BS/THES buffer (Jutras et al., 2012) to remove proteins not bound to the  
192 target DNA. Finally, protein elution was performed in Tris-HCl (pH 7) with a successively  
193 increasing concentration of NaCl (200, 300, 500, 750 mM, 1M, and 2M). The eluted protein  
194 solutions were analysed by gel electrophoresis NuPAGE® Novex®Bis-Tris (Invitrogen) and  
195 visualized using Sypro® Ruby (Molecular Probes) according to the manufacturer's instructions.  
196 The 500 mM NaCl eluate yielded the most prominent bands and therefore this eluate was used  
197 for further analysis. No bands were observed in the negative control when water was used  
198 instead of DNA (data not shown), confirming that the identified proteins were pulled down by  
199 the DNA sequences.

200

## 201 **2.7 Protein digestion for mass spectrometry-based proteomics**

202 For the digestion of proteins from gel band excision, the gel bands of interest were cut and de-  
203 stained for 1 h with a buffer of 50 mM ammonium bicarbonate (ABC) in 50% acetonitrile  
204 (ACN). Following buffer removal, 50  $\mu$ L of 10 mM DTT was added and samples were incubated  
205 for 30 min at 60 °C to reduce disulphide bonds. Next, the DTT solution was removed, and 50  $\mu$ L  
206 of 55 mM iodoacetamide (IAA) was added and samples were incubated for 30 min in the dark at  
207 room temperature to alkylate sulfhydryl groups. After removal of the IAA solution, gel pieces  
208 were washed twice with 100  $\mu$ L of 50 mM ABC, and dehydrated with 100% ACN. Protein  
209 digestion was performed overnight at 37 °C by rehydrating gel pieces with 50  $\mu$ L of  
210 Trypsin/Lys-C mix (10 ng/ $\mu$ L in 25 mM ABC) and 100  $\mu$ L of ABC.

211 Extraction of peptides from gel pieces was performed by repeating the following steps  
212 five times: addition of 100  $\mu$ L of 0.1% formic acid (FA) in 50% ACN and sonication of samples  
213 in a water bath for 10 min. Samples were then concentrated to near dryness using a centrifugal  
214 vacuum concentrator (Eppendorf) and resuspended in 50  $\mu$ L of 0.1% FA. Finally, samples were  
215 desalted using C<sub>18</sub> ZipTips (Merck Millipore) as follows: the column was wetted using 0.1% FA  
216 in 100% ACN, equilibrated with 0.1% FA in 70% ACN, and washed with 0.1% FA before  
217 loading the sample and washing again with 0.1% FA. Finally, peptides were eluted with 0.1%  
218 FA in 70% ACN, and then diluted 10-fold with 0.1% FA for mass spectrometry analysis.

219

220 For the digestion of proteins from the whole purified DNA bound material, the whole purified  
221 DNA-bound material from the DNA-protein binding assay was incubated for 30 min at 95 °C.  
222 Next, 30  $\mu$ L of 10 mM DTT was added and samples were incubated for 45 min at 55 °C to  
223 reduce disulphide bonds. Then, 40  $\mu$ L of 55 mM IAA was added and samples were incubated for  
224 30 min in the dark at room temperature to alkylate sulfhydryl groups. Protein digestion was  
225 performed overnight at 37 °C using 50  $\mu$ L of Trypsin/Lys-C mix (10 ng/ $\mu$ L in 25 mM ABC) and  
226 stopped by lowering the pH to 3 using FA. Finally, the samples were desalted and prepared for  
227 mass spectrometry analysis as described above.

228

## 229 **2.8 Protein identification using mass spectrometry**

230 Detection of proteins in both the digestion products of gel band excision and the whole captured  
231 material was performed using a QTOF Sciex 5600 or a Thermo Orbitrap Elite mass spectrometer  
232 (depending on instrument availability) with details described elsewhere (Kappler and Nouwens,  
233 2013) (Yang et al 2016) with a modified liquid chromatographic (LC) gradient. Protein  
234 identification was performed using the software ProteinPilot v5.0 (ABSciex) with the Paragon  
235 Algorithm against the NC\_022592.1 and CP006763 genome annotations with the following  
236 search parameters: Trypsin+LysC digestion; IAA as cysteine alkylation; Thorough search effort;  
237 FDR analysis. Only proteins below 1% false discovery rate (FDR; estimated global) and with at  
238 least two peptides with more than 95% confidence were considered as identified.

239

## 240 **2.9 Molecular Biology Techniques**

241 The full list of bacterial strains, plasmids, and primers used in this work for the *in vivo*  
242 transcription assay and protein overexpression step are shown in Supplementary file 2. Luria-  
243 Bertani (LB) broth or agar with antibiotics were used for growth .

244 *E. coli* DH5 $\alpha$  was used as the cloning strain and performed transformations according to  
245 the manufacturer's instructions (BIOLINE). *E. coli* BL21 was used in the *in vivo* transcription  
246 assay and protein overexpression step. *E. coli* BL21 chemically competent cells were prepared  
247 using the RuCl<sub>2</sub> method (Green and Rogers, 2013).

248 PCR amplification of targeted sequences was performed using the Phusion polymerase  
249 (NEB) and the OneTaq polymerase (NEB). Plasmid were assembled using standard ligation with  
250 the T4 DNA ligase or using Gibson assembly (Gibson et al., 2009).

251

### 252 **2.9.1 Construction of a $\sigma$ -factor candidate expression system in *E. coli***

253 Candidates for potential  $\sigma$  factors were selected based on protein identification using mass  
254 spectrometry (see above) from proteins annotated as transcriptional regulatos (Table 1).  
255 Additionally, we also built a plasmid for the L-seryl-tRNA(Sec) selenium transferase  
256 (CAETHG\_2839) (identified as a stronger band in the pull-down assay (Figure 3B)), and the  
257 housekeeping  $\sigma$  in *Clostridia* ( $\sigma^A$ ) (CAETHG\_2917) (Figure 3B).

258 The potential  $\sigma$  factor candidates were cloned into plasmid pET28a+ to be expressed  
259 under the control of a T7 promoter. DNA sequences were PCR amplified using the primers  
260 shown in Supplementary file 2 and purified using a QIAGEN kit. Next, the plasmid pET28a+  
261 was linearised using restriction enzymes NdeI and HindIII and purified using a QIAGEN kit.  
262 Codon optimisation was required to express the  $\sigma$  factor candidates of TetR-family protein  
263 (CAETHG\_0459) and  $\sigma^A$  (CAETHG\_2917) before DNA sequences were synthesised as gene  
264 block (gBlock®) fragments.

265 Plasmids with the  $\sigma$  factor candidates were then assembled by Gibson assembly using  
266 equimolar concentrations of the linearised backbone plasmid and the PCR fragment in a 20  $\mu$ L  
267 reaction. After incubation at 50 °C, 5  $\mu$ L of the Gibson mix was then used to transform *E. coli*  
268 DH5 $\alpha$  by heat shock. After recovery on SOC media at 37° C for 60 min, 100  $\mu$ L of cells were  
269 spread on LB agar plates containing kanamycin (50  $\mu$ g/mL). Plates were then incubated at 37 °C  
270 for 16 h and kanamycin resistant colonies were tested by colony PCR for proper assembly using  
271 pET\_conf(FWD)/pET\_conf(REV) primers (Supplementary file 2). A colony that tested positive  
272 for assembly was then picked and grown overnight on LB media containing kanamycin.  
273 Plasmids were recovered from 5 mL of overnight culture using a QIAGEN miniprep kit and the  
274 digestion profile was verified with the assembly. Plasmids were then used to transform *E. coli*  
275 BL21 chemically competent cells (described above).

276 *E. coli* BL21 strains harbouring  $\sigma$  factor candidate-expressing plasmids were then grown  
277 overnight on LB media containing kanamycin and 1 mM IPTG (Isopropyl  $\beta$ -D-1-  
278 thiogalactopyranoside). Next, 2 mL of overnight culture were spun down and the supernatant  
279 was removed. Next, the cell pellet was resuspended in the BugBuster master mix solution  
280 (Novagen) for protein extraction following the manufacturer's instructions. The insoluble and  
281 soluble fractions were loaded into an SDS-PAGE gel to confirm the overexpression of the  $\sigma$   
282 factor candidates (data not shown).

283

### 284 **2.9.2 Construction of a $P_{\text{cauto}}$ -GFP-UV reporter fusion system in *E. coli***

285 To determine whether the  $\sigma$  factor candidates could activate transcription, we assembled a GFP-  
286 based reporter expression system under the control of the  $P_{\text{cauto}}$ . Firstly, plasmid pBR322 was  
287 digested with HindIII, purified, and used as the backbone followed by PCR amplification of the  
288 DNA sequence containing  $P_{\text{cauto}}$  from the *C. autoethanogenum* genome (500 bp upstream of the  
289 start codon of the gene CAETHG\_1617) and purification using a QIAGEN kit.

290 Next, the GFP-UV gene was PCR amplified from plasmid pBR\_PprpR-GFPuv and  
291 purified after which the three DNA fragments were added at an equimolar concentration to a  
292 Gibson assembly mix subsequently incubated at 50 °C. 5  $\mu$ L of the Gibson mix was used to  
293 transform chemically competent *E. coli* DH5 $\alpha$  cells by heat shock and after incubation at 37 °C,  
294 100  $\mu$ L of cells were spread on LB agar plates containing ampicillin (100  $\mu$ g/mL) and incubated  
295 at 37 °C for 16 h. Ampicillin resistant colonies were then tested by colony PCR using the primer  
296 sets of  $P_{\text{cauto}}$ -GFP\_conf(FWD-1)/  $P_{\text{cauto}}$ -GFP\_conf(REV-1) and  $P_{\text{cauto}}$ -GFP\_conf(FWD-2)/  $P_{\text{cauto}}$ -  
297 GFP\_conf(REV-2) (supplementary file 2). Confirmed colonies were picked and grown overnight  
298 on LB containing antibiotic for plasmid recovery. The digestion profile confirmed the assembly  
299 of plasmid pBR\_  $P_{\text{cauto}}$ -GFP.

300 The  $P_{\text{cauto}}$ -GFP-UV was excised from pBR\_  $P_{\text{cauto}}$ -GFP using restriction enzyme HindIII.  
301 Digestion mix was loaded on a 1% agarose gel and the  $P_{\text{cauto}}$ -GFP-UV region recovered using a  
302 QIAGEN gel extraction kit. Then, the recovered DNA sequence was cloned into plasmid  
303 pACYC184, which was previously digested with HindIII and purified using a QIAGEN kit.

304 Ligation was performed according to the manufacturer's instruction and 5  $\mu$ L of the mix  
305 was used to transform *E. coli* DH5 $\alpha$  competent cells. After heat shock and incubation, 100  $\mu$ L of  
306 cells were spread on LB agar containing chloramphenicol (30  $\mu$ g/mL) and incubated at 37 °C for  
307 16 h. Chloramphenicol-resistant colonies were tested by colony PCR for proper assembly.  
308 Positive colonies were then grown overnight on LB media and the plasmid was recovered.  
309 Assembly of plasmid pACYC\_  $P_{\text{cauto}}$ -GFP was confirmed by digestion profile and Sanger  
310 sequencing (AGRF, Australia) (data not shown).

311

### 312 **2.9.3 Construction of a variants for the $P_{\text{cauto}}$ promoter motif region.**

313 Later a new reporter system including the  $P_{\text{cauto}}$  and the GFPuv sequences was built to  
314 remove the 500 bp upstream region in pBR\_  $P_{\text{cauto}}$ -gfp. The idea was to keep only the sequence  
315 used for the pull-down assay plus including the ribosomal binding site (Shine-Dalgarno  
316 sequence) to be tested *in vivo* with TetR-family protein (CAETHG\_0459) and  $\sigma^A$   
317 (CAETHG\_2917) (see net section), the two proteins that responded positively in the *in vivo*  
318 assay (see results). This new plasmid, pBR\_  $P_{\text{cauto}}$ 130\_gfp, was built by cloning the PCR product  
319 of primers WLP130F and WLP130R using pBR\_  $P_{\text{cauto}}$ -gfp as template, at the HindIII site of  
320 pBR322 by Gibson assembly (supplementary file 2). Then, the  $P_{\text{cauto}}$ 130\_gfp region was excised  
321 from pBR\_  $P_{\text{cauto}}$ 130\_gfp using HindIII and ClaI, and cloned by ligation in pACYC184 to build

322 plasmid pAC\_P<sub>cauto</sub>130\_gfp. A variation of the promoter region (pAC\_P<sub>cauto</sub>30C\_gfp) was also  
323 built to introduce single nucleotide changes in the WLP promoter motif. Changes were as follow:  
324 **ctggagcaggttttgtagttgcagtaactggttcaata**, changed to **ccatcaaggtcttaagttgcagtaactggttcaata**. This  
325 promoter was again tested with the TetR-family protein (CAETHG\_0459) and  $\sigma^A$   
326 (CAETHG\_2917). All plasmids maps used can be found in supplemental information.

327

328

#### 329 **2.9.4 *In vivo* transcription activation of P<sub>cauto</sub>-GFP(UV) fusion by the candidate genes in** 330 ***E. coli***

331 *E. coli* BL21 was used for the *in vivo* assay. Firstly, six biological replicate cultures of  
332 cells were grown in a 96-well plate (Corning Costar catalogue number #3799) carrying the  
333 pACYC plasmid with or without (to correct for the autofluorescence of the cells) the promoter-  
334 GFPuV fusion reporter in *trans* with a pET plasmid carrying each of the  $\sigma$  factor candidates.  
335 Additionally, a system with cells carrying either the pACYC promoter-GFPuV fusion reporter or  
336 its backbone plasmid plus the pET plasmid with no candidate was used as the control.

337 Cells were grown in 150  $\mu$ L of LB media containing kanamycin and chloramphenicol at  
338 30 °C and agitation of 200 RPM. At mid-exponential phase, cells were sub-cultured to a black  
339 96-well plate (Greiner #655090) to an initial OD of 0.05-0.1 in LB media containing kanamycin  
340 and chloramphenicol supplemented with either 0.0 mM IPTG (No IPTG) or 1.0 mM IPTG. The  
341 *in vivo* experiment was performed at 30 °C and agitation of 200 RPM.

342 Growth was followed by measuring the optical density (OD) at 600 nm while  
343 fluorescence intensity (FI; for GFP expression) was measured using the excitation filter of 355  
344 nm and an emission filter of 520 nm. The experiment was conducted using the FLUOstar Omega  
345 microplate reader and the Omega software v.1.20 (BMG LabTech). Fluorescence intensity was  
346 normalized per OD (FI/OD) and the signal resulting from the cells harbouring the backbone  
347 plasmid only was subtracted from the cells carrying the promoter-GFP fusion reporter  
348 (Normalized FI/OD).

349 For the WLP promoter motif variants (described in the previous sentence) four biological  
350 replicates were used.

351 Student's t-test (two-tailed) was performed between each of the candidate's normalized  
352 FI/OD value without and with IPTG and between the control system. A candidate gene was  
353 considered to activate gene expression from P<sub>cauto</sub> if it met both of the following two conditions:  
354 1) there was a statistically significant difference (p-value<0.01) in FI/OD between the candidate  
355 without and with IPTG; 2. there was a statistically significant difference (p-value<0.01) between  
356 the FI/OD signal of the candidate and the control vector (PET\_) with IPTG.

357

#### 358 **2.9.5 Overproduction and purification of TetR-family protein (CAETHG\_0459)**

359 To enable the test whether the TetR-family protein CAETHG\_0459 activates  
360 transcription from P<sub>cauto</sub> by interacting directly with the RNAP, the target protein had to be  
361 heterologously expressed and purified for the protein-protein interaction assay (see 2.9.5).

362 The *E. coli* strain harbouring the plasmid pET\_TetR1 (CAETHG\_0459) was grown at 30  
363 ° C and 200 RPM until mid-exponential phase in LB media containing kanamycin. Cells were  
364 sub-cultured to 1 L LB media containing kanamycin to an initial OD of 0.05-0.1 and  
365 subsequently grown until OD ~1 at 30 °C and 200 RPM. Then, 1.0 mM IPTG was added and  
366 cells were left growing until OD ~3. Cells were pelleted from 1 L culture by centrifugation at  
367 5,000  $\times g$  for 20 min at 4 °C, the pellet was resuspended in 5 mL of the BugBuster Master Mix



368 (Merck Millipore #71456) per gram of wet cell weight with EDTA-free protease inhibitor  
369 cocktail (Sigma #11836170001), and then incubated in a rotating mixer for 20 min at room  
370 temperature. Next, cells debris were removed by centrifugation at  $16,000 \times g$  for 20 min at 4 °C  
371 and the supernatant (supplemented with 20 mM Imidazole) was loaded on a 1 mL Ni<sup>+</sup>-  
372 HisTrapHP column (GE Healthcare #71-5027-68 AK) and washed with a buffer containing 100  
373 mM Tris-HCl (pH 7), 100 mM NaCl, 20 mM Imidazole.

374 The TetR-family protein protein CAETHG\_0459 was eluted in the same wash buffer  
375 containing a stepwise imidazole gradient (50-500 mM) following a buffer exchange performed  
376 using a HiTrap Desalting column (GE Healthcare #17-1408-01). Finally, the purified protein was  
377 stored in 50 mM Na<sub>2</sub>HPO<sub>4</sub>, 300 mM NaCl, pH7, 50% glycerol. Protein purity was analysed by  
378 gel electrophoresis using NuPAGE® Novex®Bis-Tris (Invitrogen) and stained with  
379 SimplyBlue™ SafeStain (Novex). Protein concentration was measured by the Direct Detect  
380 Spectrometer (Merck Millipore).

### 381 382 **2.9.6 TetR-family protein (CAETHG\_0459)-RNA polymerase Core enzyme interaction** 383 **experiment**

384 The protein-protein interaction (PPI) experiment was performed as described previously  
385 (Raffestin et al., 2005) with some modifications. The purified TetR-family protein  
386 (CAETHG\_0459) with 6-His-tag (2 µg) was coupled to Ni<sup>+</sup>-NTA agarose beads (Thermo  
387 #88831) in 800 µL of buffer A (50 mM Na<sub>2</sub>HPO<sub>4</sub>, 300 mM NaCl, 50 mM imidazole, pH 7). The  
388 beads coupled with the target protein were then washed three times in buffer B (50 mM Na<sub>2</sub>HPO<sub>4</sub>,  
389 300 mM NaCl, 0.1% Tween 20, 50 mM imidazole, pH 7). Next, the beads-protein complex was  
390 incubated with *E. coli* RNA polymerase Core enzyme (2.5 µg) (BioLabs #M0550S) at 37 °C for  
391 2 h. After two washes in buffer A, the beads-protein complex was suspended in 15 µL of  
392 Laemmli Buffer (32.9 mM Tris-HCl, pH6.8, 13.15 % (w/v) glycerol, 1.05 % SDS, 0.005%  
393 bromophenol blue, 355 mM 2-mercaptoethanol), heated at 100 °C for 5 min, and analysed by gel  
394 electrophoresis using NuPAGE® Novex®Bis-Tris (Invitrogen) and stained with SimplyBlue™  
395 SafeStain (Novex). The negative control was performed by incubating the RNA polymerase  
396 Core enzyme with Ni<sup>+</sup>-NTA agarose beads following the same procedure.

### 397 398 **2.9.7 Visualization of cells harbouring the Pcauto-GFP(UV) fusion and the TetR-family** 399 **protein (CAETHG\_0459) plasmids by microscopy**

400 Cells carrying the Pcauto-GFP(UV) fusion reporter and the TetR-family protein  
401 (CAETHG\_0459) plasmids were analysed by microscopy to visualize the expression of GFP.  
402 For this, cells were plated in an LB agar plate (LB media containing 6 g/L of agar) containing 1.0  
403 mM IPTG, kanamycin, and chloramphenicol. After overnight incubation at 37° C, colonies were  
404 visualized using the ZOE™ Fluorescent Cell Imager (Bio-Rad) using the manufacturer's  
405 instructions and following parameters: Gain: 40; Exposure (ms): 340; LED intensity: 22;  
406 Contrast: 59.

407  
408

## 409 **3. Results**

### 410 **3.1 Differential RNA-sequencing (dRNA-Seq)**

411  
412 In this work, we aimed to determine the TSSs of essential genes for autotrophic growth of the  
413 model-acetogen *C. autoethanogenum* (e.g. genes in the WLP and of hydrogenases). We thus

414 performed dRNA-Seq analysis (Sharma et al., 2010) of autotrophic (CO, CO<sub>2</sub>, and H<sub>2</sub>; referred  
415 to as ‘syngas’) and heterotrophic (fructose) cultures of *C. autoethanogenum* to experimentally  
416 determine TSSs and promoter motif(s) associated with essential genes for autotrophic growth in  
417 acetogens.

418  
419 Previously described batch cultures (Marcellin et al., 2016) were sampled during exponential  
420 growth and subjected to dRNA-Seq cDNA library preparation and sequencing. The cDNA  
421 libraries were prepared using the 5'tagRACE method (Fouquier D'Hérouel et al., 2011), an  
422 improved library preparation method compared to TEX (5'-phosphate-dependent Terminator  
423 RNA exonuclease) that has the advantage of preserving the quantitative representation of 5' ends  
424 between processed (5'-P end) and primary (5'-PPP end) transcripts (see Methods). TSSs were  
425 determined by comparing the libraries enriched for processed (TAP-) and primary (TAP+)  
426 transcripts (Figure 1A) using the TSSAR tool (Amman et al., 2014).

427

### 428 **3.2 Overall dRNA-Seq features of *C. autoethanogenum***

429

430 We classified TSSs as primary, internal, antisense, and orphan (Figure 1B, Table S1) and found  
431 primary TSSs only for around half of the annotated genes (3,983) in *C. autoethanogenum*  
432 (Brown et al., 2014) (Table S1). More than 60% of the genes contain only one primary TSS,  
433 while the rest show up to 12 TSSs (Figure 1C, Table S1). Focusing on the 14 main metabolic  
434 groups of *C. autoethanogenum* genes as described in (Brown et al., 2014), we detected primary  
435 TSSs for all genes except for the Nfn transhydrogenase complex (CAETHG\_1580) (Table S2).  
436 While primary TSSs were detected for seven of the 11 genes of the WLP biosynthetic gene  
437 cluster (CAETHG\_1606-21), only half of the WLP TSSs were shared between syngas and  
438 fructose. For example, genes of the WLP methyl branch (CAETHG\_1614-17) contained 20  
439 primary TSSs on syngas compared to only nine on fructose. On the other hand, the TSSs  
440 associated with Hydrogenases and ATPase genes were found in similar numbers between syngas  
441 and fructose.

442

443 Determination of nucleotide base preferences for transcription initiation within five nucleotides  
444 downstream and upstream of the primary TSSs showed a clear enrichment of adenine (A) and  
445 guanine (G) at +1 (~90%) and thymine (T) at -1 for both syngas (Figure 1D) and fructose (data  
446 not shown). Overall, adenine and cytosine were the most and least preferred nucleotide bases,  
447 respectively.

448 Analysis of 5'untranslated regions (5'UTRs)—the sequence between the TSS and the  
449 annotated start codon—indicates transcripts potentially associated with post-transcriptional  
450 regulation and thus of mRNA stability and translational efficiency (Cho et al., 2014). Calculation  
451 of 5'UTR lengths for primary TSSs showed a median length of 63 nt with 65% of TSSs <100 nt  
452 for both growth conditions (Figure 1E and Table S1). Genes with longer UTR lengths tend to be  
453 regulated more at the post-transcriptional level (Cho et al., 2009; David et al., 2006). On the  
454 other hand, leaderless mRNAs—mRNAs with no or <10 nt 5'UTR—are translated in the absence  
455 of upstream signals (typically the Shine-Dalgarno sequence) (Shine and Dalgarno, 1974; Zheng  
456 et al., 2011) used for regulating translational efficiency through ribosome binding. We found ~70  
457 (~2%) leaderless mRNAs with <10 nt 5'UTRs, none of which were in the WLP, Hydrogenases,  
458 Acetate or Ethanol groups (Figure 1E and Table S3).

459 In addition to the ability to determine TSSs, dRNA-Seq analysis also facilitates a more  
460 accurate annotation of the genome. Based on the TSSs and the Shine-Dalgarno (AGGAGG)  
461 position that was found to be highly conserved within 9-14 nt upstream of the first start codon  
462 (ATG/CTG/GTG/TTG) (Figure 1F), we re-annotated the start codon for 38 genes and confirmed  
463 the changes in one gene by peptide identification using mass spectrometry (Table S4). Moreover,  
464 either the start or stop codon of an additional 99 genes, which had previously been annotated in  
465 different frames, were manually corrected. The corrections have been deposited into NCBI under  
466 the accession number BK010482 and the complete manually corrected genbank file of *C.*  
467 *autoethanogenum* is available in Table S5.

468

### 469 **3.3 Discovery of a new promoter motif**

470 The RNA polymerase (RNAP) needs to form a holoenzyme with a  $\sigma$  factor in bacteria to  
471 recognise a specific promoter motif (sequence) and initiate transcription (Feklistov et al., 2014;  
472 Gruber and Gross, 2003). Experimentally determined TSS data from dRNA-Seq analysis is ideal  
473 for *in silico* determination of promoter motifs, which is important for understanding  
474 transcriptional regulation, especially in less-studied bacteria such as acetogens.

475 We searched for consensus sequence motifs 50 nt upstream of primary TSSs using the  
476 MEME software (Bailey et al., 2009) and were able to determine seven promoter motifs in *C.*  
477 *autoethanogenum* (E-value  $\leq 0.05$ ) (Tables S6 and 7 for syngas and fructose growth,  
478 respectively). Of those identified, only three motifs were assigned with more than 100 TSSs and  
479 shared between the two datasets, likely representing the most conserved motifs in *C.*  
480 *autoethanogenum* (Figure 2A).

481 The top motif was found 10 nt upstream of primary TSSs (447 and 543 TSSs for syngas  
482 and fructose, respectively; E-value  $< 10^{-111}$ ) and resembles the Pribnow box (TATGnTATAAT),  
483 which is associated with the housekeeping  $\sigma$  factors of *Escherichia coli* ( $\sigma^{70}$ ; (Walker and Osuna,  
484 2002)), *Helicobacter pylori* ( $\sigma^{80}$ ; (Sharma et al., 2010)) and *Clostridium acetobutylicum* ( $\sigma^A$ ;  
485 (Sauer et al., 1994, 1995)). Expectedly, the well-known -35 TTGACA and -10 TATAAT motifs  
486 (TATA box in eukaryotes and archaea) for housekeeping  $\sigma$  factors (Burgess and Anthony, 2001)  
487 was also among the top-3 promoter consensus sequences (392 and 262 TSSs for syngas and  
488 fructose, respectively; E-value  $< 10^{-46}$ ). These two motifs were assigned for most of the genes of  
489 glycolysis/gluconeogenesis and the TCA cycle (Table S2).

490 The third most abundant promoter motif has, to the best of our knowledge, not previously  
491 been reported in the literature (Figure 2A).  $P_{\text{cauto}}$  is highly conserved both during growth on  
492 syngas (Motif 02 in Table S5; 392 TSSs; E-value  $< 10^{-174}$ ) and fructose (Motif 03 in Table S6; 224  
493 TSSs; E-value  $< 10^{-77}$ ). Importantly,  $P_{\text{cauto}}$  seems to be involved in the transcriptional regulation of  
494 essential genes for acetogens and was assigned to genes of the WLP cluster (CAETHG\_1606-21)  
495 and the metabolic groups, as described in (Brown et al., 2014), of Hydrogenases, Acetate,  
496 ATPase, and Pyruvate (Figure 2B; Tables S2, S6, and S7). We confirmed the unique presence of  
497 the “new promoter motif” upstream of the TSSs. Investigation of its upstream regions up to 100  
498 or 150 nt showed no other motif apart from the one conserved within 50 nt upstream of TSSs.  
499 This new promoter is well characterised by an evenly interspaced (A/T)G repetition with an  
500 almost central A/T position (Figure S1). These observations potentially indicate the presence of a  
501 new  $\sigma$  factor or transcriptional regulator of critical importance in acetogens.

502

### 503 **3.4 RNA polymerase and proteins annotated as transcriptional regulators specifically bind**

504  $P_{\text{cauto}}$

505 We performed DNA-protein binding assays to determine if the RNAP and/or other  
506 protein(s) bind to  $P_{\text{cauto}}$ . The promoter sequences of two WLP genes (CAETHG\_1615 and 1617,  
507 Methylene-tetrahydrofolate reductase domain-containing protein and Methenyltetrahydrofolate  
508 cyclohydrolase, respectively) annotated with  $P_{\text{cauto}}$  were used for the DNA-protein binding assay  
509 using the promoter pull down/DNA affinity chromatography method (Figure 3A; (Jutras et al.,  
510 2012)). The promoter sequence of a glycolytic gene (CAETHG\_3424, glyceraldehyde-3-  
511 phosphate dehydrogenase, type I) was included as a control for the assay since it was assigned  
512 the well-known TATAAT motif, which should yield binding of the RNAP and the housekeeping  
513  $\sigma$  factor,  $\sigma^A$ . DNA-bound proteins captured using streptavidin-coupled magnetic Dynabeads<sup>TM</sup>  
514 were identified using mass spectrometry of the digestion products of the whole captured material  
515 and of gel band excisions. Since this DNA-protein binding assay requires significant amounts of  
516 cellular protein material, especially for efforts to identify low abundance proteins such as  $\sigma$   
517 factors or transcriptional regulators, autotrophic bioreactor chemostat cultures (CO or CO+H<sub>2</sub>) of  
518 *C. autoethanogenum* described in a separate work (Valgepea et al., 2018) were sampled for this  
519 analysis.

520 The promoter pull down/DNA affinity chromatography method (Figure 3A; (Jutras et al.,  
521 2012)) was fine-tuned for *C. autoethanogenum*. Eluting the proteins with 500 mM NaCl yielded  
522 the most prominent bands while no bands were observed in the negative control when water was  
523 used instead of DNA (data not shown), which confirms that the identified proteins were pulled  
524 down by the DNA sequences (see Methods). The alpha and beta subunits of the RNAP  
525 (CAETHG\_1920 and 1954-55) were successfully identified for both  $P_{\text{cauto}}$  (CAETHG\_1615 and  
526 CAETHG\_1617) and the TATAAT motif control (Figure 3B). Additionally, the RNAP omega  
527 subunit was identified in the whole purified DNA-bound material for both motifs (Table S8).  
528 The housekeeping  $\sigma^A$  (CAETHG\_2917) was detected for the TATAAT motif control as expected  
529 (Figure 3B). A stronger band was identified in the  $P_{\text{cauto}}$  gels around 50 kDa and identified as a  
530 protein annotated as L-seryl-tRNA(Sec) selenium transferase (CAETHG\_2839; 51.5 kDa)  
531 (Figure 3B). Finally, mass spectrometry analysis of the whole purified DNA-bound material  
532 identified three proteins annotated as transcriptional regulators (based on NC\_022592.1) that  
533 were unique for the  $P_{\text{cauto}}$  (Table 1) and found for both CO and CO+H<sub>2</sub> cultures across technical  
534 replicates of the DNA-protein binding assay (Table S8).

### 535 536 **3.5 TetR-family transcriptional regulator (CAETHG\_0459) activates transcription from** 537 **$P_{\text{cauto}}$ *in vivo***

538 To determine whether any of the three identified protein candidates annotated as transcriptional  
539 regulators that uniquely bind to  $P_{\text{cauto}}$  (Table 1) could activate transcription from this promoter,  
540 we created a transcriptional fusion reporter vector harbouring the sequence of  $P_{\text{cauto}}$  in-frame  
541 with a green fluorescence protein (GFPuV). We also tested transcriptional activation using the L-  
542 seryl-tRNA(Sec) selenium transferase (CAETHG\_2839) (identified as a stronger band in the  
543 pull-down assay (Figure 3B)), and using the housekeeping  $\sigma$  factor in clostridia ( $\sigma^A$ )  
544 (CAETHG\_2917), since it has been reported that promoter binding sites of different  $\sigma$  factors  
545 can overlap (Cho et al., 2014). Transcriptional activation of  $P_{\text{cauto}}$  with concomitant GFP  
546 production was investigated in *E. coli* by inducing the expression of the candidate activator  
547 proteins from a second T7 protein over-expression vector cloned into plasmid pET28e+ by the  
548 addition of IPTG (see Methods). Fluorescence was measured at early-exponential growth (OD  
549 ~0.26) as FI/OD.

550 After subtracting the signal from cells harbouring the two plasmids but lacking the fusion  
551 reporter (promoter + GFP, see Methods), only induction of the TetR-family transcriptional  
552 regulator protein (CAETHG\_0459) (out of the three transcriptional regulator candidates) led to  
553 statistically higher levels of GFP expression ( $p < 0.01$ ) compared to the control vector with no  
554 candidate (Figure 4A and Table S9). Interestingly, induction of  $\sigma^A$  also led to transcription  
555 activation ( $p < 0.01$ ). We then confirmed expression of GFP in the strain expressing  
556 CAETHG\_0459 grown on a plate with IPTG using fluorescence microscopy (Figure 4B). This  
557 shows that both CAETHG\_0459 and  $\sigma^A$  independently activate transcription from  $P_{\text{cauto}}$ .  
558 Importantly, the motif is associated with the expression of essential genes in gas-fermenting  
559 acetogens including genes in the WLP and hydrogenases (Table S2, S6-7).

560 The 130bp variant (which includes the sequence used for the pull-down assay plus the  
561 ribosomal binding site) also showed statistical significance ( $p$ -value  $< 0.01$ ) of fluorescence  
562 increase when TetR-family transcriptional regulator protein (CAETHG\_0459) was present.  
563 Similarly,  $\sigma^A$  could also activate transcription, however only at the level of  $p$ -value  $< 0.05$ .  
564 Interestingly when mutations were included in the promoter motif, TetR- (CAETHG\_0459)  
565 could no longer activate expression of GFP, as expected (Figure 4A).

566

### 567 **3.6 CAETHG\_0459 directly binds to the RNA polymerase core enzyme**

568 As TetR-family proteins often act as transcriptional regulators (Cuthbertson and Nodwell,  
569 2013), we next investigated whether TetR-family protein CAETHG\_0459 activates transcription  
570 from  $P_{\text{cauto}}$  by interacting directly with the RNAP. Transcriptional regulators can reversibly  
571 interact with the RNAP Core enzyme independently of a DNA sequence to help activate  
572 transcription from a range of promoters (Burgess and Anthony, 2001; Feklistov et al., 2014). We  
573 thus performed an *in vitro* protein-protein interaction assay to test whether protein  
574 CAETHG\_0459 directly interacts with RNA polymerase Core in the absence of DNA. The  
575 purified His-tagged CAETHG\_0459 protein linked to  $\text{Ni}^{2+}$ -beads was incubated with the RNAP  
576 Core enzyme (see Methods). SDS-PAGE analysis clearly demonstrated an interaction between  
577 the core RNA polymerase and CAETHG\_0459 (Figure 4C lane 6) and shows that  
578 CAETHG\_0459 acts as a positive transcriptional regulator that activates transcription from  $P_{\text{cauto}}$   
579 by directly binding to the RNAP.

580

### 581 **3.7 $P_{\text{cauto}}$ is represented in other acetogens**

582 We next investigated if  $P_{\text{cauto}}$  was represented in other industrially relevant acetogens  
583 with available genomes: *Clostridium ljungdahlii*, *C. ragsdalei*, *C. coskatii*, *Moorella*  
584 *thermoacetica*, and *Eubacterium limosum* (Bengelsdorf et al., 2016; Redl et al., 2017; Shin et al.,  
585 2016; Song et al., 2017). We performed the reverse of the methodology previously used to search  
586 for consensus sequence motifs by looking for the occurrence of  $P_{\text{cauto}}$  300 nt upstream of  
587 annotated genes (since no TSS data was available) using the FIMO tool (Grant et al., 2011)  
588 within MEME. As expected based on their phylogenetic proximity (Bengelsdorf et al., 2013;  
589 Brown et al., 2014; Shin et al., 2016), *C. ljungdahlii*, *C. ragsdalei*, and *C. coskatii* showed  
590 similar occurrences of  $P_{\text{cauto}}$  (Figure 2C). Interestingly, while the representation in *M.*  
591 *thermoacetica* was very low,  $P_{\text{cauto}}$  seems to be present also in *E. limosum*. This result highlights  
592 the need for experimental determination of TSSs in more acetogens.

593

## 594 **4. Discussion**

595 Acetogens offer an enormous potential for the production of fuels and chemicals from  
596 gaseous waste feedstocks (Claassens et al., 2016; Dürre and Eikmanns, 2015; Liew et al., 2016;  
597 Molitor et al., 2016), with ethanol already being produced at industrial scale by LanzaTech.  
598 Acetogens have two major carbon fixation pathways: the WLP for autotrophic growth and  
599 glycolysis for heterotrophic growth. Although both the WLP and glycolysis/gluconeogenesis  
600 pathways operate during autotrophic and heterotrophic growth, the WLP carries a substantially  
601 higher metabolic flux during autotrophy (Valgepea et al., 2017a, 2018) and *vice versa* (Valgepea  
602 et al., 2017b). We and others have shown that transcriptional regulation between autotrophic and  
603 heterotrophic growth in acetogens is not trivial (Aklujkar et al., 2017; Marcellin et al., 2016;  
604 Nagarajan et al., 2013; Tan et al., 2013). We thus aimed to determine TSSs and transcriptional  
605 features of promoter motifs and transcriptional regulators associated with essential genes  
606 (including genes of the WLP) in the model-acetogen *C. autoethanogenum*.

607 Our study revealed a new promoter motif and the identification of two proteins  
608 activating gene expression from the new motif (the TetR-family protein (CAETHG\_0459) and  
609 the housekeeping  $\sigma^A$  (CAETHG\_2917)). An alternative TetR transcriptional regulator has been  
610 previously found to be a  $\sigma$  factor in *Clostridium tetani*, and its homologues, TcdR in *C. difficile*,  
611 BotR in *C. botulinum*, and UviA in *C. perfringens* have also been found to regulate toxin  
612 production (Dupuy et al., 2006; Dupuy and Matamouros, 2006; Raffestin et al., 2005). In  
613 combination, these results suggest that TetR proteins can play an important role in transcriptional  
614 regulation in *clostridia*. These studies support our PPI assay potentially suggesting that the TetR-  
615 family protein might function as a  $\sigma$  factor in *C. autoethanogenum*, but further studies (*in vitro*  
616 transcription assay) are needed to confirm this. In fact, unequivocal demonstration of  $\sigma$  factor  
617 activity requires that a protein is necessary and sufficient for activation of promoter recognition  
618 and transcription initiation by RNAP, independent of any other  $\sigma$  factor subunit. Thus our results  
619 do not exclude the possibility that a native  $\sigma$  factor of the *in vivo* expression host (*E. coli*), e.g.  
620  $\sigma^{70}$ , could have induced the TetR-family protein to drive transcription from  $P_{\text{cauto}}$ . Additional  
621 studies should also be performed to study whether both the  $\sigma^A$  and the TetR-family protein show  
622 an overlap in the promoter motif for transcriptional activation (Cho et al., 2014).

623 Notably, there are several TetR-family proteins, commonly regarded as transcriptional  
624 regulators (Cuthbertson and Nodwell, 2013), annotated in the *C. autoethanogenum* genome. In  
625 pathogenic *clostridia* these TetR-family proteins are often described as alternative  $\sigma$  factors,  
626 belonging to a class of  $\sigma$  factors called extracytoplasmic function (ECF)  $\sigma$  factors (Feklistov et  
627 al., 2014; Sineva et al., 2017). Their discovery led to a novel class of  $\sigma$  factors (group 5), which  
628 show a -35 and/or -10 conserved region in their target promoters (Dupuy et al., 2005, 2006;  
629 Dupuy and Matamouros, 2006; Staroń et al., 2009). It will be interesting to see whether  
630 transcription from  $P_{\text{cauto}}$  described here with an interspaced repetition of (A/T)G notably distinct  
631 from the canonical -35/-10 conserved regions is also activated by a novel  $\sigma$  factor.

632 Our work also shows that the housekeeping  $\sigma$  factor ( $\sigma^A$ ) in *clostridia* can activate  
633 transcription from  $P_{\text{cauto}}$  associated with essential genes for autotrophic growth in acetogens.  
634 Interestingly, in another acetogen *E. limosum*, the promoter regions of genes of the WLP,  
635 hydrogenases, and ATPase contain the well-known -35 TTGACA and -10 TATAAT motifs for  
636 the housekeeping  $\sigma$  factor ( $\sigma^A$ ) (Burgess and Anthony, 2001; Song et al., 2017). This potentially  
637 indicates that the housekeeping  $\sigma^A$  in acetogens can initiate transcription from different promoter  
638 motifs and illustrates well the great extent of genetic diversity among the non-taxonomic group  
639 of acetogens. While the WLP itself is highly conserved, it is not surprising that transcriptional  
640 regulation is diverse (Drake et al., 2006; Shin et al., 2016). The work presented here also

641 highlights the importance of  $P_{\text{cauto}}$  in other industrially relevant acetogens (Figure 2C). We  
642 believe, however, that more studies are needed for the experimental determination of TSSs and  
643 transcriptional features to facilitate a broader understanding of transcriptional regulation in  
644 acetogens.

645 Our findings have the potential to significantly advance the understanding of  
646 transcriptional regulation and metabolic engineering of the ancient metabolism of acetogens.  
647 Firstly, acetogen metabolism, which operates at the thermodynamic edge of feasibility  
648 (Schuchmann and Müller, 2014), seems to be wired for utilising less energy-consuming  
649 mechanisms (i.e. transcriptional vs. translational regulation) for operating under different  
650 conditions evidenced by the complexity of the condition-specific transcriptional architecture  
651 (Valgepea et al., 2018). More importantly, the discovery of  $P_{\text{cauto}}$  and a key positive transcription  
652 factor (TetR-family protein) in acetogens can lead to the mechanistic description of  
653 transcriptional regulation of arguably the first biochemical pathway on Earth (Fuchs, 2011;  
654 Russell and Martin, 2004; Weiss et al., 2016). In addition to expanding the fundamental  
655 understanding of a model acetogen, knowledge of the features controlling the expression of  
656 essential genes in acetogens could also contribute for the improvement of commercial gas  
657 fermentation for the sustainable production of fuels and chemicals. Increasing or modulation of  
658 the activity of the described TetR transcription factor (either through over-expression and/or  
659 protein engineering or by deleting transcriptional repressor genes) could enhance the uptake of  
660  $C_1$  substrates through the WLP and thus improve growth and/or product formation (possibly by  
661 introducing  $P_{\text{cauto}}$  in front of key genes). It could also be used as an orthologous system in other  
662 organisms, as, for instance, the TcdR system has been used in other *Clostridium* species (Minton  
663 et al., 2016; Zhang et al., 2015). Importantly, the newly discovered promoter  $P_{\text{cauto}}$  could be  
664 harnessed to couple expression of heterologous pathways to mimic those of key central  
665 metabolism enzymes, potentially alleviating the common problem of imbalanced flux throughput  
666 between heterologous and native metabolic pathways.

667

### 668 **Conflict of interest**

669 RT, MK and SDS are employed by Lanzatech. The authors declare that this study received  
670 funding from the Australian Research Council (ARC LP140100213) and LanzaTech through and  
671 ARC linkage grant. LanzaTech has interest in commercialising gas fermentation with *C.*  
672 *autoethanogenum*. RT, MK and SS were involved in experimental design, data analysis and  
673 interpretation and were involved in writing the manuscript.

674

### 675 **Author contributions**

676 (i) RL, KV, MK, RT, LN and EM designed the study and the experiments (ii) RL, RG, RP, KV,  
677 CB performed the experiments. RL, KV, RT, RP, MK, SS, LN and EM analysed and interpreted  
678 the data; (iii) RL, KV, RT and EM wrote the manuscript. All authors reviewed the manuscript.

679

### 680 **Funding**

681 This work was funded by the Australian Research Council (ARC LP140100213) in collaboration  
682 with LanzaTech. The ARC had no role in study design, data collection and interpretation, or the  
683 decision to submit the work for publication. There was no funding support from the European  
684 Union for the experimental part of the study. However, KV acknowledges support also from the  
685 European Union's Horizon 2020 research and innovation programme under grant agreement  
686 N810755.

687

## 688 **Acknowledgements**

689 The authors acknowledge support from the Queensland node of Metabolomics Australia  
690 (MA) at The University of Queensland, an NCRIS initiative under Bioplatforms Australia Pty  
691 Ltd. We thank Dr Christopher Howard for his helpful advice in the pull down/affinity  
692 chromatography assay and protein purification. We thank the following investors in LanzaTech's  
693 technology: BASF, CICC Growth Capital Fund I, CITIC Capital, Indian Oil Company, K1W1,  
694 Khosla Ventures, the Malaysian Life Sciences, Capital Fund, L. P., Mitsui, the New Zealand  
695 Superannuation Fund, Petronas Technology Ventures, Primetals, Qiming Venture Partners,  
696 Softbank China, and Suncor.

697

## 698 **References**

- 699 Aklujkar, M., Leang, C., Shrestha, P. M., Shrestha, M., and Lovley, D. R. (2017).  
700 Transcriptomic profiles of *Clostridium ljungdahlii* during lithotrophic growth with syngas  
701 or H<sub>2</sub> and CO<sub>2</sub> compared to organotrophic growth with fructose. *Sci. Rep.* 7, 13135.  
702 doi:10.1038/s41598-017-12712-w.
- 703 Amman, F., Wolfinger, M. T., Lorenz, R., Hofacker, I. L., Stadler, P. F., and Findeiß, S. (2014).  
704 TSSAR: TSS annotation regime for dRNA-seq data. *BMC Bioinformatics* 15, 89.  
705 doi:10.1186/1471-2105-15-89.
- 706 Bailey, T. L., Boden, M., Buske, F. A., Frith, M., Grant, C. E., Clementi, L., et al. (2009).  
707 MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res.* 37, 202–208.  
708 doi:10.1093/nar/gkp335.
- 709 Bengelsdorf, F. R., Poehlein, A., Linder, S., Erz, C., Hummel, T., Hoffmeister, S., et al. (2016).  
710 Industrial Acetogenic Biocatalysts: A Comparative Metabolic and Genomic Analysis.  
711 *Front. Microbiol.* 7, 1–15. doi:10.3389/fmicb.2016.01036.
- 712 Bengelsdorf, F. R., Straub, M., and Dürre, P. (2013). Bacterial synthesis gas (syngas)  
713 fermentation. *Environ. Technol.* 34, 1639–1651. doi:10.1080/09593330.2013.827747.
- 714 Brown, S. D., Nagaraju, S., Utturkar, S., De Tissera, S., Segovia, S., Mitchell, W., et al. (2014).  
715 Comparison of single-molecule sequencing and hybrid approaches for finishing the genome  
716 of *Clostridium autoethanogenum* and analysis of CRISPR systems in industrial relevant  
717 *Clostridia*. *Biotechnol. Biofuels* 7, 40. doi:10.1186/1754-6834-7-40.
- 718 Burgess, R. R., and Anthony, L. (2001). How sigma docks to RNA polymerase and what sigma  
719 does. *Curr. Opin. Microbiol.* 4, 126–131. doi:10.1016/S1369-5274(00)00177-6.
- 720 Cho, B.-K., Kim, D., Knight, E. M., Zengler, K., and Palsson, B. O. (2014). Genome-scale  
721 reconstruction of the sigma factor network in *Escherichia coli*: topology and functional  
722 states. *BMC Biol.* 12, 4. doi:10.1186/1741-7007-12-4.
- 723 Cho, B.-K., Zengler, K., Qiu, Y., Park, Y. S., Knight, E. M., Barrett, C. L., et al. (2009). The  
724 transcription unit architecture of the *Escherichia coli* K-12 MG1655 genome. *Nat.*  
725 *Biotechnol.* 27, 1043–1049. doi:10.1038/nbt.1582.
- 726 Claassens, N. J., Sousa, D. Z., dos Santos, V. A. P. M., de Vos, W. M., and van der Oost, J.  
727 (2016). Harnessing the power of microbial autotrophy. *Nat. Rev. Microbiol.* 14, 692–706.  
728 doi:10.1038/nrmicro.2016.130.
- 729 Cuthbertson, L., and Nodwell, J. R. (2013). The TetR Family of Regulators. *Microbiol. Mol.*  
730 *Biol. Rev.* 77, 440–475. doi:10.1128/MMBR.00018-13.
- 731 David, L., Huber, W., Granovskaia, M., Toedling, J., Palm, C. J., Bofkin, L., et al. (2006). A  
732 high-resolution map of transcription in the yeast genome. *Proc. Natl. Acad. Sci. U. S. A.*



- 733 103, 5320–5325. doi:10.1073/pnas.0601091103.
- 734 Drake, H. L., Küsel, K., and Matthies, C. (2006). “Acetogenic Prokaryotes,” in *Prokaryotes*  
735 (*Ecophysiology and Biochemistry*), eds. M. Dworkin, E. Rosenberg, K. H. Schleifer, and E.  
736 Stackebrandt (New York: Springer), 354–420.
- 737 Dupuy, B., Mani, N., Katayama, S., and Sonenshein, A. L. (2005). Transcription activation of a  
738 UV-inducible *Clostridium perfringens* bacteriocin gene by a novel sigma factor. *Mol.*  
739 *Microbiol.* 55, 1196–1206. doi:10.1111/j.1365-2958.2004.04456.x.
- 740 Dupuy, B., and Matamouros, S. (2006). Regulation of toxin and bacteriocin synthesis in  
741 *Clostridium* species by a new subgroup of RNA polymerase  $\sigma$ -factors. *Res. Microbiol.* 157,  
742 201–205. doi:10.1016/j.resmic.2005.11.004.
- 743 Dupuy, B., Raffestin, S., Matamouros, S., Mani, N., Popoff, M. R., and Sonenshein, A. L.  
744 (2006). Regulation of toxin and bacteriocin gene expression in *Clostridium* by  
745 interchangeable RNA polymerase sigma factors. *Mol. Microbiol.* 60, 1044–1057.  
746 doi:10.1111/j.1365-2958.2006.05159.x.
- 747 Dürre, P., and Eikmanns, B. J. (2015). C1-carbon sources for chemical and fuel production by  
748 microbial gas fermentation. *Curr. Opin. Biotechnol.* 35, 63–72.  
749 doi:10.1016/j.copbio.2015.03.008.
- 750 Feklistov, A., Sharon, B. D., Darst, S. A., and Gross, C. A. (2014). Bacterial Sigma Factors: A  
751 Historical, Structural, and Genomic Perspective. *Annu. Rev. Microbiol.* 68, 357–76.  
752 doi:10.1146/annurev-micro-092412-155737.
- 753 Fouquier D’Hérouel, A., Wessner, F., Halpern, D., Ly-Vu, J., Kennedy, S. P., Serror, P., et al.  
754 (2011). A simple and efficient method to search for selected primary transcripts: Non-  
755 coding and antisense RNAs in the human pathogen *Enterococcus faecalis*. *Nucleic Acids*  
756 *Res.* 39, e46. doi:10.1093/nar/gkr012.
- 757 Fuchs, G. (2011). Alternative Pathways of Carbon Dioxide Fixation: Insights into the Early  
758 Evolution of Life? *Annu. Rev. Microbiol.* 65, 631–658. doi:10.1146/annurev-micro-090110-  
759 102801.
- 760 Gibson, D. G., Young, L., Chuang, R. Y., Venter, J. C., Hutchison, C. A., and Smith, H. O.  
761 (2009). Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nat.*  
762 *Methods* 6, 343–345. doi:10.1038/nmeth.1318.
- 763 Grant, C. E., Bailey, T. L., and Noble, W. S. (2011). FIMO: Scanning for occurrences of a given  
764 motif. *Bioinformatics* 27, 1017–1018. doi:10.1093/bioinformatics/btr064.
- 765 Green, R., and Rogers, E. J. (2013). Chemical Transformation of *E. coli*. *Methods Enzymol.*,  
766 329–336. doi:10.1016/B978-0-12-418687-3.00028-8.Chemical.
- 767 Gruber, T. M., and Gross, C. A. (2003). Multiple sigma subunits and the partitioning of bacterial  
768 transcription space. *Annu. Rev. Microbiol.* 57, 441–66.  
769 doi:10.1146/annurev.micro.57.030502.090913.
- 770 Herrmann, G., Jayamani, E., Mai, G., and Buckel, W. (2008). Energy conservation via electron-  
771 transferring flavoprotein in anaerobic bacteria. *J. Bacteriol.* 190, 784–791.  
772 doi:10.1128/JB.01422-07.
- 773 Jutras, B. L., Verma, A., and Stevenson, B. (2012). “Identification of Novel DNA-Binding  
774 Proteins Using DNA-Affinity Chromatography/Pull Down,” in *Current Protocols in*  
775 *Microbiology* (Wiley Online Library), 1–13. doi:10.1002/9780471729259.mc01f01s24.
- 776 Kappler, U., and Nouwens, A. S. (2013). The molybdoproteome of *Starkeya novella* – insights  
777 into the diversity and functions of molybdenum containing proteins in response to changing  
778 growth conditions. *Metallomics* 5, 325. doi:10.1039/c2mt20230a.

- 779 Kim, D., Perteua, G., Trapnell, C., Pimentel, H., Kelley, R., and Salzberg, S. L. (2013). TopHat2:  
780 accurate alignment of transcriptomes in the presence of insertions, deletions and gene  
781 fusions. *Genome Biol.* 14, R36. doi:10.1186/gb-2013-14-4-r36.
- 782 Li, F., Hinderberger, J., Seedorf, H., Zhang, J., Buckel, W., and Thauer, R. K. (2008). Coupled  
783 ferredoxin and crotonyl coenzyme A (CoA) reduction with NADH catalyzed by the butyryl-  
784 CoA dehydrogenase/Etf complex from *Clostridium kluyveri*. *J. Bacteriol.* 190, 843–850.  
785 doi:10.1128/JB.01417-07.
- 786 Liew, F., Martin, E., Tappel, R., Heijstra, B., Mihalcea, C., and Köpke, M. (2016). Gas  
787 Fermentation – A Flexible Platform for Commercial Scale Production of Low Carbon Fuels  
788 and Chemicals from Waste and Renewable Feedstocks. *Front. Microbiol.* 7, 694.  
789 doi:10.3389/fmicb.2016.00694.
- 790 Ljungdahl, L. G. (2009). A life with acetogens, thermophiles, and cellulolytic anaerobes. *Annu.*  
791 *Rev. Microbiol.* 63, 1–25. doi:10.1146/annurev.micro.091208.073617.
- 792 Marcellin, E., Behrendorff, J. B., Nagaraju, S., DeTissera, S., Segovia, S., Palfreyman, R., et al.  
793 (2016). Low carbon fuels and commodity chemicals from waste gases – Systematic  
794 approach to understand energy metabolism in a model acetogen. *Green Chem.* 18, 3020–  
795 3028. doi:10.1039/C5GC02708J.
- 796 Minton, N. P., Ehsaan, M., Humphreys, C. M., Little, G. T., Baker, J., Henstra, A. M., et al.  
797 (2016). A roadmap for gene system development in *Clostridium*. *Anaerobe* 41, 104–112.  
798 doi:10.1016/j.anaerobe.2016.05.011.
- 799 Molitor, B., Marcellin, E., and Angenent, L. T. (2017). Overcoming the energetic limitations of  
800 syngas fermentation. *Curr. Opin. Chem. Biol.* 41, 84–92. doi:10.1016/j.cbpa.2017.10.003.
- 801 Molitor, B., Richter, H., Martin, M. E., Jensen, R. O., Juminaga, A., Mihalcea, C., et al. (2016).  
802 Carbon recovery by fermentation of CO-rich off gases - Turning steel mills into  
803 biorefineries. *Bioresour. Technol.* 215, 386–396. doi:10.1016/j.biortech.2016.03.094.
- 804 Nagarajan, H., Sahin, M., Nogales, J., Latif, H., Lovley, D. R., Ebrahim, A., et al. (2013).  
805 Characterizing acetogenic metabolism using a genome-scale metabolic reconstruction of  
806 *Clostridium ljungdahlii*. *Microb. Cell Fact.* 12, 118. doi:10.1186/1475-2859-12-118.
- 807 Nitschke, W., and Russell, M. J. (2011). Redox bifurcations: Mechanisms and importance to life  
808 now, and at its origin. *BioEssays* 34, 106–109. doi:10.1002/bies.201100134.
- 809 Patrik, D. (2006). What are DNA sequence motifs? *Nat. Biotechnol.* 24, 423–425.
- 810 Raffestin, S., Dupuy, B., Marvaud, J. C., and Popoff, M. R. (2005). BotR/A and TetR are  
811 alternative RNA polymerase sigma factors controlling the expression of the neurotoxin and  
812 associated protein genes in *Clostridium botulinum* type A and *Clostridium tetani*. *Mol.*  
813 *Microbiol.* 55, 235–249. doi:10.1111/j.1365-2958.2004.04377.x.
- 814 Ragsdale, S. W. (1991). Enzymology of the acetyl-CoA pathway of CO<sub>2</sub> fixation. *Crit. Rev.*  
815 *Biochem. Mol. Biol.* 26, 261–300. doi:10.3109/10409239109114070.
- 816 Ragsdale, S. W. (1997). The eastern and western branches of the Wood/Ljungdahl pathway: how  
817 the east and west were won. *Biofactors* 6, 3–11. doi:10.1002/biof.5520060102.
- 818 Ragsdale, S. W. (2008). Enzymology of the Wood-Ljungdahl pathway of acetogenesis. *Ann. N.*  
819 *Y. Acad. Sci.* 1125, 129–136. doi:10.1196/annals.1419.015.
- 820 Ragsdale, S. W., and Pierce, E. (2008). Acetogenesis and the Wood-Ljungdahl pathway of CO<sub>2</sub>  
821 fixation. *Biochim. Biophys. Acta* 1784, 1873–1898. doi:10.1016/j.bbapap.2008.08.012.
- 822 Redl, S., Sukumara, S., Ploeger, T., Wu, L., Jensen, T. Ø., Nielsen, A. T., et al. (2017).  
823 Thermodynamics and economic feasibility of acetone production from syngas using the  
824 thermophilic production host *Moorella thermoacetica*. *Biotechnol. Biofuels* 10, 150.

- 825 doi:10.1186/s13068-017-0827-8.
- 826 Russell, M. J., and Martin, W. (2004). The rocky roots of the acetyl-CoA pathway. *Trends*  
827 *Biochem. Sci.* 29, 358–363. doi:10.1016/j.tibs.2004.05.007.
- 828 Sauer, U., Santangelo, J. D., Treuner, A., Buchholz, M., and Durre, P. (1995). Sigma factor and  
829 sporulation genes in *Clostridium*. *FEMS Microbiol. Rev.* 17, 331–340. doi:10.1016/0168-  
830 6445(95)00005-W.
- 831 Sauer, U., Treuner, A., Buchholz, M., Santangelo, J. D., and Durre, P. (1994). Sporulation and  
832 primary sigma factor homologous genes in *Clostridium acetobutylicum*. *J. Bacteriol.* 176,  
833 6572–6582.
- 834 Schuchmann, K., and Müller, V. (2014). Autotrophy at the thermodynamic limit of life: a model  
835 for energy conservation in acetogenic bacteria. *Nat. Rev. Microbiol.* 12, 809–821.  
836 doi:10.1038/nrmicro3365.
- 837 Sharma, C. M., Hoffmann, S., Darfeuille, F., Reignier, J., Findeiß, S., Sittka, A., et al. (2010).  
838 The primary transcriptome of the major human pathogen *Helicobacter pylori*. *Nature* 464,  
839 250–255. doi:10.1038/nature08756.
- 840 Shin, J., Song, Y., Jeong, Y., and Cho, B. K. (2016). Analysis of the core genome and pan-  
841 genome of autotrophic acetogenic bacteria. *Front. Microbiol.* 7, 1531.  
842 doi:10.3389/fmicb.2016.01531.
- 843 Shine, J., and Dalgarno, L. (1974). The 3'-terminal sequence of *Escherichia coli* 16S ribosomal  
844 RNA: complementarity to nonsense triplets and ribosome binding sites. *Proc. Natl. Acad.*  
845 *Sci. U. S. A.* 71, 1342–6. doi:10.1073/pnas.71.4.1342.
- 846 Sineva, E., Savkina, M., and Ades, S. E. (2017). Themes and variations in gene regulation by  
847 extracytoplasmic function (ECF) sigma factors. *Curr. Opin. Microbiol.* 36, 128–137.  
848 doi:10.1016/j.mib.2017.05.004.
- 849 Song, Y., Shin, J., Jeong, Y., Jin, S., Lee, J.-K., Kim, D. R., et al. (2017). Determination of the  
850 Genome and Primary Transcriptome of Syngas Fermenting Eubacterium limosum ATCC  
851 8486. *Sci. Rep.* 7, 13694. doi:10.1038/s41598-017-14123-3.
- 852 Staroń, A., Sofia, H. J., Dietrich, S., Ulrich, L. E., Liesegang, H., and Mascher, T. (2009). The  
853 third pillar of bacterial signal transduction: Classification of the extracytoplasmic function  
854 (ECF)  $\sigma$  factor protein family. *Mol. Microbiol.* 74, 557–581. doi:10.1111/j.1365-  
855 2958.2009.06870.x.
- 856 Tan, Y., Liu, J., Chen, X., Zheng, H., and Li, F. (2013). RNA-seq-based comparative  
857 transcriptome analysis of the syngas-utilizing bacterium *Clostridium ljungdahlii* DSM  
858 13528 grown autotrophically and heterotrophically. *Mol. Biosyst.* 9, 2775–84.  
859 doi:10.1039/c3mb70232d.
- 860 Valgepea, K., de Souza Pinto Lemgruber, R., Abdalla, T., Binos, S., Takemori, N., Takemori, A.,  
861 et al. (2018). H<sub>2</sub> drives metabolic rearrangements in gas-fermenting *Clostridium*  
862 autoethanogenum. *Biotechnol. Biofuels* 11, 55. doi:10.1186/s13068-018-1052-9.
- 863 Valgepea, K., de Souza Pinto Lemgruber, R., Meaghan, K., Palfreyman, R. W., Abdalla, T.,  
864 Heijstra, B. D., et al. (2017a). Maintenance of ATP Homeostasis Triggers Metabolic Shifts  
865 in Gas-Fermenting Acetogens. *Cell Syst.* 4, 505–515. doi:10.1016/j.cels.2017.04.008.
- 866 Valgepea, K., Loi, K. Q., Behrendorff, J. B., de Souza Pinto Lemgruber, R., Plan, M., Hodson,  
867 M. P., et al. (2017b). Arginine deiminase pathway provides ATP and boosts growth of the  
868 gas-fermenting acetogen *Clostridium autoethanogenum*. *Metab. Eng.* 41, 202–211.  
869 doi:10.1016/j.ymben.2017.04.007.
- 870 Walker, K. A., and Osuna, R. (2002). Factors affecting start site selection at the *Escherichia coli*

871        fis promoter. *J. Bacteriol.* 184, 4783–4791. doi:10.1128/JB.184.17.4783-4791.2002.  
872 Weiss, M. C., Sousa, F. L., Mrnjavac, N., Neukirchen, S., Roettger, M., Nelson-Sathi, S., et al.  
873        (2016). The physiology and habitat of the last universal common ancestor. *Nat. Microbiol.*  
874        1, 16116. doi:10.1038/nmicrobiol.2016.116.  
875 Wood, H. G. (1991). Life with CO or CO<sub>2</sub> and H<sub>2</sub> as a source of carbon and energy. *FASEB J.*  
876        5, 156–163.  
877 Yang et al, (2016) The Snake with the Scorpion’s Sting: Novel Three-Finger Toxin Sodium  
878        Channel Activators from the Venom of the Long-Glanded Blue Coral Snake (*Calliophis*  
879        *bivirgatus*). *Toxins* 8(10), 303  
880 Zhang, Y., Grosse-Honebrink, A., and Minton, N. P. (2015). A universal mariner transposon  
881        system for forward genetic studies in the genus *Clostridium*. *PLoS One* 10, 1–21.  
882        doi:10.1371/journal.pone.0122411.  
883 Zheng, X., Hu, G., She, Z., and Zhu, H. (2011). Leaderless genes in bacteria: clue to the  
884        evolution of translation initiation mechanisms in prokaryotes. *BMC Genomics* 12, 361.  
885        doi:10.1186/1471-2164-12-361.

886  
887

## 888 **Figure legends**

889

890 **Figure 1.** Characteristics of transcriptional and translational architecture in *C. autoethanogenum*.  
891 (A) Our dRNA-Seq approach generated genome-wide TSS maps through the comparison of  
892 libraries enriched for processed (TAP-) and primary (TAP+) transcripts. (B) Classification of  
893 TSSs for syngas and fructose as: primary, within 250 nt upstream of an annotated gene; internal,  
894 within an annotated gene; antisense, on the opposite strand of an annotated gene; orphan, not  
895 assigned to any of the previous classes. (C) Distribution of primary TSSs per gene for syngas and  
896 fructose. (D) Nucleotide base preference for transcription initiation from primary TSSs on  
897 syngas. +1 denotes the position of the TSS. (E) Distribution of 5'UTR lengths for primary TSSs  
898 for syngas and fructose. (F) The Shine-Dalgarno sequence AGGAGG is highly conserved within  
899 9-14 nt upstream of the first start codon. Sequencing reads were processed with the TSSAR  
900 software (Amman et al., 2014) for automated *de novo* determination of TSSs from dRNA-Seq  
901 data using the following parameters: p-Value 1e-3, Noise threshold 10, Merge range 5. The  
902 Shine-Dalgarno sequence was searched 30 nt upstream of annotated genes (CP006763.1 and  
903 NC\_022592.1) using the MEME software (Bailey et al., 2009) and the same parameters as for  
904 promoter motif search, except for -nmotifs 10, -maxw 30. See Methods for details.

905  
906

907 **Figure 2.** *In silico* determination of genome-wide promoter motifs in *C. autoethanogenum*. (A)  
908 The top-3 promoter motifs for primary TSSs are shared among syngas and fructose. The height  
909 of the letter indicates its relative frequency at the given position within the motif. Refer to Tables  
910 S5-8 for all the determined motifs and their assigned TSSs. The mutated nucleotides used in the  
911 *in vivo* assay for P<sub>cauto</sub> motif are also shown. We show the nucleotide position relative to the TSS  
912 in all top3 motifs (B) The new promoter motif (P<sub>cauto</sub>) is assigned with TSSs of essential genes in  
913 acetogens. Motifs with the lowest p-value for syngas are shown. Refer to Tables S2 and 5-8 for  
914 all the TSSs and genes associated with P<sub>cauto</sub>. (C) The P<sub>cauto</sub> motif is represented in other  
915 industrially relevant acetogens. Occurrence in each acetogen relative to *C. autoethanogenum* is  
916 normalised with the number of annotated genes. To determine promoter motifs in *C.*

917 *autoethanogenum*, we searched for consensus sequence motifs 50 nt upstream of primary TSSs  
 918 using the MEME software (Bailey et al., 2009) with the following parameters: -dna, -max size  
 919 10000000, -mod zoops, -nmotifs 50, -minw 4, -maxw 50, -revcomp, -oc.

920  
 921 **Figure 3.** DNA-protein binding assay shows specific binding of *C. autoethanogenum* RNAP  
 922 subunits and a selenium transferase to the new promoter motif. (A) Overview of the DNA-  
 923 protein binding assay (i.e. the promoter pull down/DNA affinity chromatography method (Jutras  
 924 et al., 2012)). (B) Separation of proteins specifically bound to the TATAAT motif (for gene  
 925 CAETHG\_3424) or the new promoter motif (for gene CAETHG\_1617) with gel electrophoresis  
 926 and identification using mass spectrometry. The alpha and beta subunits of the RNAP  
 927 (CAETHG\_1920 and 1954-55) were successfully identified for both the new promoter motif  
 928 (CAETHG\_1615 and CAETHG\_1617) and the TATAAT motif control. Technical replicate  
 929 denotes replicate of the DNA-protein binding assay (panel A) (data not shown for  
 930 CAETHG\_1615).

931  
 932 **Figure 4.** TetR-family transcriptional regulator (CAETHG\_0459) and  $\sigma^A$  (CAETHG\_2917)  
 933 activate expression from the new promoter motif. (A) *In vivo* experiment using *E. coli* cells  
 934 carrying the pACYC plasmid with the new promoter-GFPuV fusion report in *trans* with a pET  
 935 plasmid carrying each of the candidates. The experiment was conducted with either 0.0 mM or  
 936 1.0 mM IPTG. Only in the presence of TetR-family protein (CAETHG\_0459) and  $\sigma^A$   
 937 (CAETHG\_2917) the fluorescence intensity normalized per OD (FI/OD) is statistically  
 938 significantly different (p-value <0.01) compared to the control system (with no candidate  
 939 protein). 1 Cells harbouring the PET\_ (Negative control with no candidate gene); 2 Selenium  
 940 transferase (CAETHG\_2839); 3 TetR-family protein (CAETHG\_0459); 4 TetR-family protein  
 941 (CAETHG\_0936); 5 GntR (CAETHG\_3915); 6  $\sigma^A$  (CAETHG\_2917); 7 Short version (130 bp)  
 942 of pAC\_P<sub>cauto</sub>30C\_gfp and TetR-family protein (CAETHG\_0459); 8 Mutated version of the  
 943 promoter region (pAC\_P<sub>cauto</sub>30C\_gfp) by introducing nucleotide changes as follow:  
 944 ctggagcaggtttgtagttgcagtaactggtcaata, changed to **ccatcaaaggtctta**aagttgcagtaactggtcaata and  
 945 TetR-family protein (CAETHG\_0459); 9 Short version (130 bp) of pAC\_P<sub>cauto</sub>30C\_gfp and  $\sigma^A$   
 946 (B). Cells carrying the TetR-family protein (CAETHG\_0459) grown in LB-agar plate with 1 mM  
 947 IPTG were visualized under microscopy for fluorescence (GFP) visualization. (C) Protein-  
 948 protein interaction assay. TetR-family protein (CAETHG\_0459) was incubated with *E. coli* RNA  
 949 polymerase Core enzyme. Lane 1: Marker (Thermo #26614); Lane 2: *E. coli* RNA polymerase  
 950 Core Enzyme; Lane 3: *E. coli* RNA polymerase Core incubated with Ni<sup>+</sup> agarose beads and  
 951 washed; Lane 4: Purified TetR-family protein (CAETHG\_0459); Lane 5: Ni<sup>+</sup> agarose beads  
 952 coupled with TetR-family protein (CAETHG\_0459); Lane 6: Ni<sup>+</sup> agarose beads coupled with  
 953 TetR-family protein (CAETHG\_0459) incubated with RNA polymerase Core and washed; Lane  
 954 7: Marker

955  
 956 **Table 1.** *C. autoethanogenum* proteins annotated as transcriptional regulators uniquely binding  
 957 to the new promoter motif P<sub>cauto</sub>

Gene ID <sup>a</sup>	Gene ID <sup>b</sup>	Gene product annotation <sup>a</sup>
CAETHG_RS02185	CAETHG_0459	TetR/AcrR family transcriptional regulator
CAETHG_RS04465	CAETHG_0936	TetR/AcrR family transcriptional regulator
CAETHG_RS19205	CAETHG_3915	GntR family transcriptional regulator

958 <sup>a</sup> Annotation NC\_022592.1.

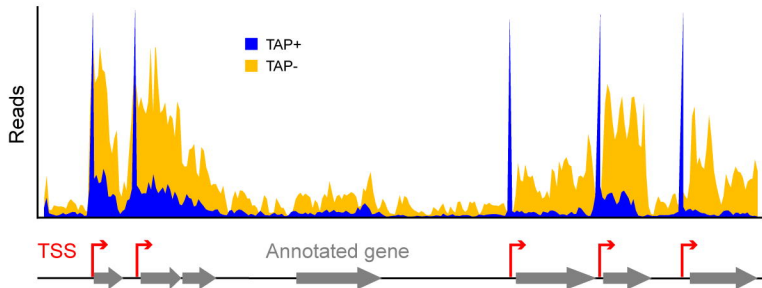
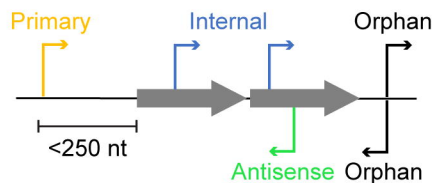
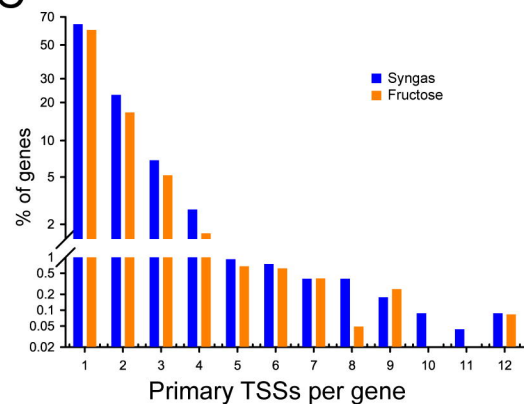
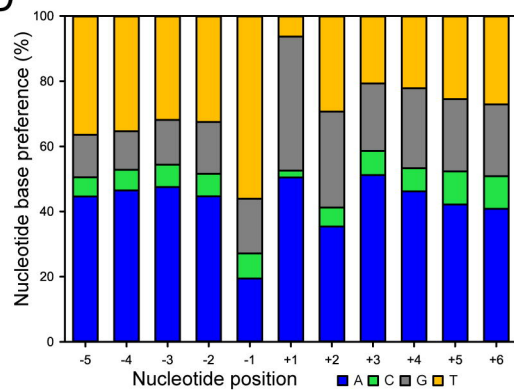
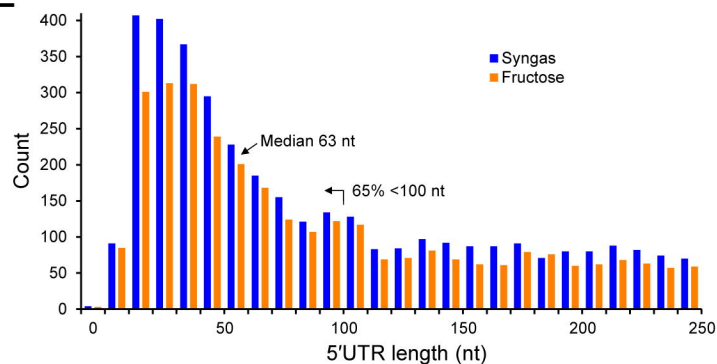
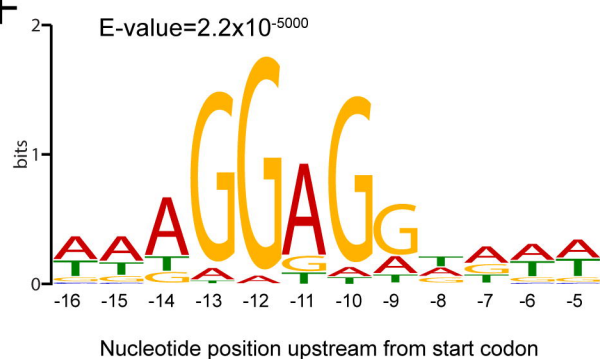
959 <sup>b</sup> Annotation CP006763.1 (Brown et al., 2014).

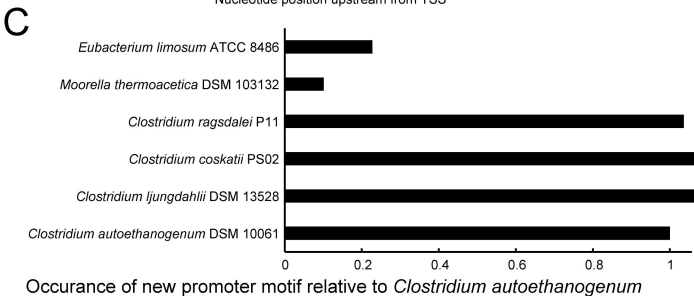
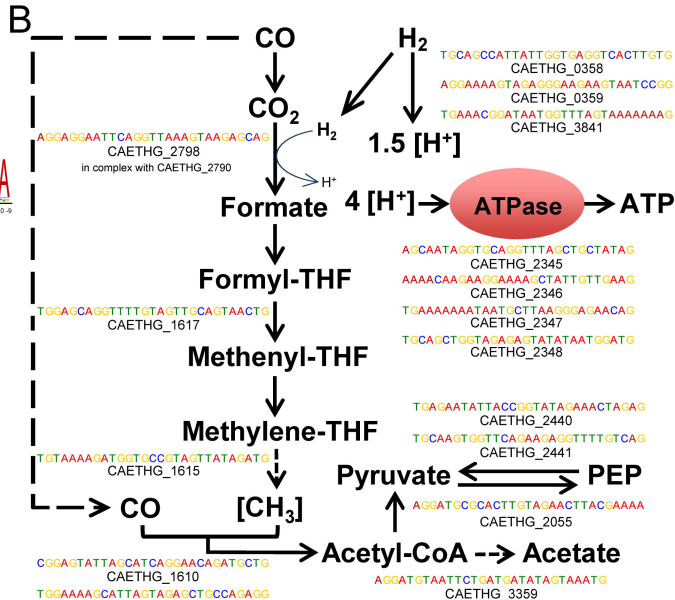
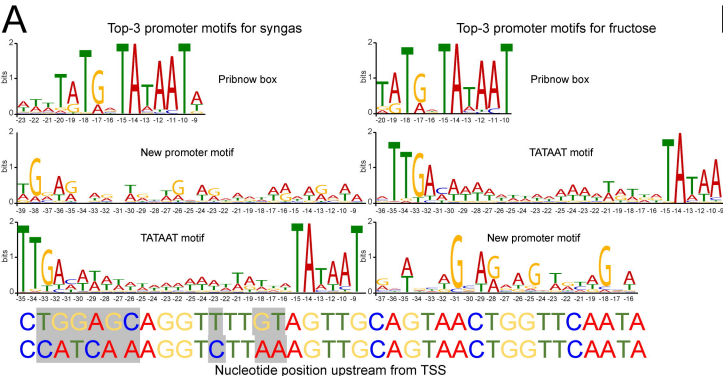
960

961 **Datasets are in a publicly accessible repository.**

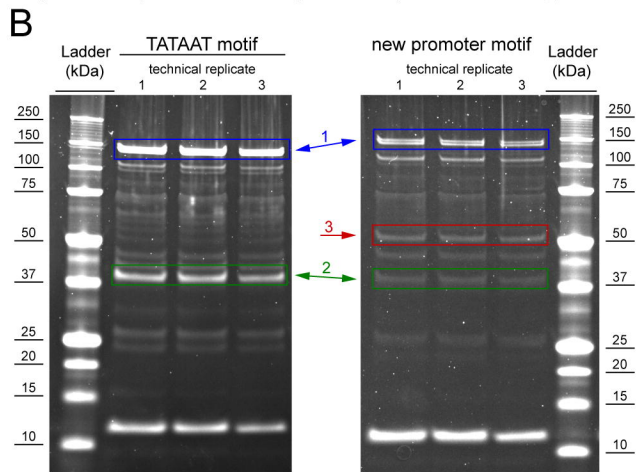
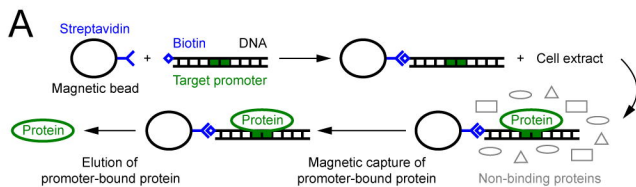
962 dRNA-Seq data have been deposited in the NCBI Gene Expression Omnibus depository under  
963 accession number GSE108700.

964 Re-annotation of *C. autoethanogenum* genome was deposited in the NCBI GenBank  
965 Third Party Annotation database under accession number BK010482.

**A****B****C****D****E****F**







Proteins identified by mass spectrometry

1-RNAP beta subunit (CAETHG\_1955) +  
RNAP beta' subunit (CAETHG\_1954)  
2-RNAP alpha subunit (CAETHG\_1920) +  
α' subunit (CAETHG\_2917)

1-RNAP beta subunit (CAETHG\_1955) +  
RNAP beta' subunit (CAETHG\_1954)  
2-RNAP alpha subunit (CAETHG\_1920)  
3-L-seryl-tRNA(Sec) selenio transferase (CAETHG\_2839)

