29 April 19 sbs

1    **Genomic and physiological analyses reveal that extremely thermophilic**

2    ***Caldicellulosiruptor changbaiensis* deploys unique cellulose attachment mechanisms**

3    ───────────────────────────────────

4    Asma M.A.M. Khan[‡], Carl Mendoza[‡], Valerie J. Hauk and Sara E. Blumer-Schuette*

5    Department of Biological Sciences

6    Oakland University, Rochester, MI 48309

7    [‡]A.M.A.M. and C.M. contributed equally to this manuscript

8

10

11   **Running Title:**          *Caldicellulosiruptor comparative genomics*

12

13   **Keywords:**          *Caldicellulosiruptor*, extreme thermophile, cellulase, tāpirin

14

17

18   * Address Correspondence to:          Sara E. Blumer-Schuette
19                                         Dept. of Biological Sciences
20                                         Oakland University
21                                         118 Library Drive
22                                         374 Dodge Hall
23                                         Phone: (248) 370-3168
24                                         Fax:    (248) 370-4225
25                                         Email: blumerschuette@oakland.edu

26 **ABSTRACT**

27 The genus *Caldicellulosiruptor* are extremely thermophilic, heterotrophic anaerobes that

28 degrade plant biomass using modular, multifunctional enzymes. Prior pangenome analyses

29 determined that this genus is genetically diverse, with the current pangenome remaining open,

30 meaning that new genes are expected with each additional genome sequence added. Given the

31 high biodiversity observed among the genus *Caldicellulosiruptor*, we have sequenced and

32 added a 14th species, *Caldicellulosiruptor changbaiensis*, to the pangenome. The pangenome

33 now includes 3,791 ortholog clusters, 120 of which are unique to *C. changbaiensis* and may be

34 involved in plant biomass degradation. Comparisons between *C. changbaiensis* and

35 *Caldicellulosiruptor bescii* on the basis of growth kinetics, cellulose solubilization and cell

36 attachment to polysaccharides highlighted physiological differences between the two species

37 which are supported by their respective gene inventories. Most significantly, these comparisons

38 indicated that *C. changbaiensis* possesses unique cellulose attachment mechanisms not

39 observed among the other strongly cellulolytic members of the genus *Caldicellulosiruptor*.

40

41 **INTRODUCTION**

42      The genus *Caldicellulosiruptor* is comprised of extremely thermophilic, fermentative

43 heterotrophs whose members have been isolated worldwide from terrestrial geothermal springs

44 or other thermal environments [37]. The original isolates from the genus *Caldicellulosiruptor*

45 were identified on the basis of their ability to grow on cellulose at elevated temperatures [56,54],

46 especially temperatures beyond the optimal growth temperature of *Ruminiclostridium*

47 *thermocellum* [48]. Interest in thermostable enzymes produced by this genus continues, as the

48 initial discovery of their multifunctional, modular enzymes [51,26,57,67] represented an alternate

49 paradigm to cellulosomes [2,52]. Further discoveries on the capabilities of these thermostable

50 enzymes include the unique mode of action used by the central cellulase, CelA, [8], synergistic

51 activity in ionic liquid optimized enzyme mixtures [45,46] and the creation of designer

52 cellulosomes from *Caldicellulosiruptor* catalytic domains [29]. Development of a genetics system

53 for *Caldicellulosiruptor bescii* [14,16] has also expanded the scope of work with this genus,

54 including metabolic engineering [10,12,13,50] and catalytic improvement [18,30,32,31,34,33].

55      The availability of genome sequences has precipitated deeper insights into the genus

56 *Caldicellulosiruptor*, including comparative studies which have identified biomarkers for plant

57 biomass deconstruction [6,5,23], novel insertion elements [15], genetic tractability [11], diverse

58 mechanisms involved in biomass solubilization [66,37], unique cellulose adhesins (tāpirins)

59 [5,37] and the identification of new combinations of catalytic domains [5,36,23]. Perhaps owing

60 to the unique thermal environments that this genus inhabits, their genomes appear to be

61 dynamic, as the first described *Caldicellulosiruptor* pangenome was predicted to be open [5],

62 and remained open after the addition of five additional genome sequences [36].

63      Here, we have analyzed the genome sequence of *Caldicellulosiruptor changbaiensis*,

64 isolated from a hot spring in the Changbai Mountains [3], representing the 14[th] and most recent

65 addition to the *Caldicellulosiruptor* pangenome. Past *Caldicellulosiruptor* pangenomes were

66 comprised of multiple species from most countries of origin, which allowed for prior analysis on

2

67  the basis of biogeography [5], with the exception of China and Japan [20]. Now with the addition

68  of the *C. changbaiensis* genome sequence, insights into the biogeography of isolates from

69  China and how they compare to the global *Caldicellulosiruptor* pangenome is possible.

70  Furthermore, on the basis of the open *Caldicellulosiruptor* pangenome [20,5], we hypothesize

71  that the *C. changbaiensis* genome may encode for novel substrate-binding proteins and/ or

72  plant biomass degrading enzymes. In addition to updating the *Caldicellulosiruptor* pangenome,

73  we also present differences in the growth physiology of *C. changbaiensis* versus

74  *Caldicellulosiruptor bescii,* currently the benchmark species against which most

75  *Caldicellulosiruptor* are compared for their plant biomass degrading capabilities.

76 **MATERIALS AND METHODS**

77 **Microbial strains and medium.** Freeze-dried stocks of *C. changbaiensis* strain CBS-Z

78 were obtained from the Leibniz Institute DSMZ – German Collection of Microorganisms and Cell

79 Cultures (DSMZ). Glycerol stocks of *C. bescii* DSM-6725 were obtained from the laboratory of

80 Robert M. Kelly, North Carolina State University (Raleigh, NC). Both species were cultured at

81 75°C on low osmolarity defined (LOD) medium [25] under a nitrogen headspace to maintain

82 anaerobic conditions and supplemented with carbohydrates as a carbon source. Carbohydrates

83 used as a carbon source included cellobiose (≥ 99%, Chem-Impex Int'l, Inc.), pectin (Sigma-

84 Aldrich), xylan (Sigma-Aldrich), glucomannan (NOW Foods), and microcrystalline cellulose (20

85 µm Sigmacell, Sigma-Aldrich). For genomic DNA isolation, *C. changbaiensis* was cultured

86 anaerobically at 75°C on low osmolarity complex (LOC) medium [25] with cellobiose as a carbon

87 source.

88 **Genomic DNA isolation.** Genomic DNA was isolated using the Joint Genome Institute's

89 CTAB-based protocol (https://jgi.doe.gov/user-programs/pmo-overview/protocols-sample-

90 preparation-information/jgi-bacterial-dna-isolation-ctab-protocol-2012/), with modifications. In

91 order to isolate enough DNA for sequencing, 500 ml of overnight *C. changbaiensis* culture was

92 harvested by centrifugation at 5000x*g*, 4°C for 20 minutes and resuspending the cell pellet in

93 14.8 ml of TE buffer, prior to lysis. Gel electrophoresis in 0.7% agarose was used to assess the

94 quality of genomic DNA and the concentration and purity of the sample for sequencing was

95 quantified using a NanoDrop spectrophotometer, and Qubit fluorometric assay (dsDNA HS

96 assay, Thermo Fisher). Prior to genome sequencing, a 16S rRNA gene fragment was amplified

97 from isolated genomic DNA using oligonucleotide primers (Eton Bioscience) previously

98 designed for identification of *C. changbaiensis* [3], for positive identification of *C. changbaiensis*

99 (**Table 1**). Amplicons were sent for Sanger sequencing (Eton Bioscience), using the same

100 oligonucleotide primers.

4

101     ***C. changbaiensis* genome sequencing, assembly and annotation.** The genome

102     sequence for *C. changbaiensis* [40] was assembled to 60-fold coverage from long-read Oxford

103     NanoPore (MinION) data generated in house, and short-read Illumina data generated by

104     Molecular Research, LP (MR DNA). Hybrid assembly of the complete *C. changbaiensis* genome

105     used Unicycler v0.4.7 [61], and annotation of the genome used the Prokaryotic Genome

106     Annotation Pipeline v4.7 [55] provided by the National Center for Biotechnology Information

107     (NCBI). The assembled genome and reads used for assembly of the *C. changbaiensis* genome

108     are available through NCBI BioProject accession PRJNA511150.

109     **Phylogenomic analysis.** Fourteen genome sequences from the genus

110     *Caldicellulosiruptor* were included in the phylogenomic analyses (see **Table 2** for genome

111     sequence accession numbers). Orthologous protein groups were classified using the

112     GET_HOMOLOGUES v20092018 software package [19], running OrthoMCL v1.4 [39],

113     COGtriangles v2.1 [35], or bidirectional best hits (BDBH) as determined by BLASTP [1,9].

114     Orthologous protein clusters were determined using the OrthoMCL parameters: 75% pairwise

115     coverage, maximum BLASTP E-value of 1e-5, and MCL inflation of 1.5. GET_HOMOLOGUES

116     was also used to parse the pangenome matrices comparing the *C. changbaiensis* genome

117     inventory against the recent 13 *Caldicellulosiruptor* pangenome [37] or the revised *C. bescii*

118     genome [22]. Core- (**Eq. 1**) and pangenome (**Eq. 2**) parameters were predicted after curve

119     fitting randomly sampled core- or pangenome data to functions previously described by Tettelin

120     et al., [58].

121     $$coregenes(g) = 1367 + 1668 \, exp\left(\frac{-g}{1.75}\right) \hspace{4cm} (1)$$

122     $$pangenes(g) = 2371 + 63.2(g-1) + exp\left(\frac{-2}{2.19}\right)\frac{1-exp\left(\frac{-(g-1)}{2.19}\right)}{1-exp\left(\frac{-1}{2.19}\right)} \hspace{2cm} (2)$$

123     Genome-level similarity was quantified as average nucleotide identity (ANIb) from the BLASTN+

124     alignment of 1,020 nt fragments from the 14 *Caldicellulosiruptor* genomes [49,27]. ANIb were

125    calculated by Pyani v.0.2.7, (https://github.com/widdowquinn/pyani) and percent identities were

126    plotted as a heatmap by the software package.

127    **Growth kinetics on polysaccharides.** *C. bescii* or *C. changbaiensis* were revived from

128    -80˚C glycerol stocks for growth curve analysis on microcrystalline cellulose, xylan, pectin or

129    glucomannan. Glycerol stocks (1 ml) were subcultured into 50 ml LOD medium for 3

130    consecutive subcultures using 2% (v/v) inoculum at each passage. Revived cultures were then

131    transferred (2% [v/v] inoculum) to LOD medium containing a 1:1 ratio of maltose (*C. bescii*) or

132    cellobiose (*C. changbaiensis*) to polysaccharide. The 1:1 mixture was then passaged (2% [v/v]

133    inoculum) three times successively in LOD medium with polysaccharide, only. Cultures for

134    growth curves were inoculated at a starting cell density of 1 x $10^6$ cells ml$^{-1}$ in 200 ml LOD plus

135    the respective polysaccharide. Biological replicates were used for each growth phase

136    experiment. Cell counting used epifluorescence microscopy at 1000x total magnification and a

137    counting reticle as described previously [28]. Cells were fixed in a final volume of 1.1 ml

138    gluteraldehyde (2.5% [v/v] in water) prior to incubation with acridine orange (1 g l$^{-1}$) and

139    approximately 5 ml sterilized water and thoroughly mixed. Stained cells were then vacuum

140    filtered through a polycarbonate 0.22 µm filter (GE). Samples were counted using a 10x10

141    reticle a total of ten times. Cell counts were averaged for calculation of cell density (cells ml$^{-1}$).

142    Doubling times are described as the number of hours per generation during exponential growth,

143    calculated as Δtime divided by the number of generations.

144    **Microcrystalline cellulose solubilization.** Solubilization of microcrystalline cellulose

145    followed protocols established by Zurawski *et al.*, [66] with modifications. *C. bescii* or *C.*

146    *changbaiensis* were cultured in serum bottles with 50 ml of LOD medium supplemented with

147    0.6g of microcrystalline cellulose (20µm Sigmacell) at a starting cell density of $10^6$ cells ml$^{-1}$.

148    Cultures were then incubated without shaking at 75˚C for seven days, after which the remaining

149    microcrystalline cellulose was harvested by centrifugation at 6000 xg, 4˚C for 15 min in a swing

150    bucket rotor. The cellulose pellet was washed four times in sterile, deionized water and air dried

6

151    at 75°C until the weight of the microcrystalline cellulose did not change. Uninoculated LOD

152    served as an abiotic control. Percent solubilization is reported as the difference in substrate

153    weight divided by the starting weight multiplied by 100. All experimental conditions were

154    measured in triplicate and significance was determined by a t-test (p-value < 0.05).

155        **Cell attachment assays.** *C. bescii* and *C. changbaiensis* cell cultures were grown to

156    early stationary phase on either xylan or cellulose (1 g $l^{-1}$) as the carbon source, and cell

157    densities were calculated before harvesting at 5000 xg for 10 minutes at room temperature.

158    Cells were resuspended and concentrated ten-fold in the binding buffer (50 mM sodium

159    phosphate, pH 7.2) to a 10-fold density of approximately 1-2 x $10^9$ cells $ml^{-1}$ for cells cultured on

160    xylan or 1 x $10^8$ cells $mil^{-1}$ for cells cultured on cellulose. For each treatment condition, 1.2 ml of

161    *C. bescii* or *C. changbaiensis* planktonic cells in binding buffer were added to a 1.5 ml

162    microcentrifuge tube, and supplemented with 10 mg of washed substrate (experimental

163    condition: xylan or cellulose), or no substrate for the negative control. All assay tubes were

164    incubated at room temperature for one hour with gentle rotary shaking at 100 rpm. After

165    incubation, planktonic cells were enumerated as described above for the growth curves. Each

166    binding assay was repeated six times. Two-sample t-tests were used to analyze the data using

167    the R studio statistics package v.3.3.3 [47].

168

7

169     **RESULTS AND DISCUSSION**

170          **Phylogenomic analysis of the *C. changbaiensis* genome.** With the addition of the

171     fourteenth *Caldicellulosiruptor* genome [40], we sought to define an updated core- and

172     pangenome. Three different algorithms: OrthoMCL [39], bidirectional best hit and COGtriangles

173     [35] were used to classify orthologous clusters for pangenome analysis (**Table S1**). Of the three,

174     the clusters formed by OrthoMCL resulted in an estimated core- and pangenome with the lowest

175     residual standard errors, and are reported here (**Fig. 1**). Overall, there are 120 unique protein

176     clusters identified in the *C. changbaiensis* genome when compared to the prior

177     *Caldicellulosiruptor* pangenome [37], 75 of which were annotated as hypothetical proteins.

178     Further transcriptomic and proteomic studies may aid in the identification of the function of these

179     unique hypothetical proteins. By adding a 14th genome, the *Caldicellulosiruptor* core genome

180     was reduced to 1,367 orthologous clusters (see **Eq. 1**), however, the pangenome (3,791 genes)

181     continues to expand at an estimated rate of 63.2 genes per additional genome (**Eq. 2**, **Fig. 1**)

182     highlighting the plasticity of the *Caldicellulosiruptor* pangenome.

183          In contrast to previously released genome sequences from New Zealand [36], *C.*

184     *changbaiensis* exhibits a similar pattern of biogeography based on average nucleotide identity

185     (ANIb). As expected, *Caldicellulosiruptor* sp. F32, isolated from compost in China [63], and *C.*

186     *naganoensis*, isolated from a hot spring in Japan [56] shared higher percent identity levels with

187     *C. changbaiensis*, along with *C. saccharolyticus*, isolated from a hot spring in New Zealand (**Fig.**

188     **2, Table S2**). All species that *C. changbaiensis* shared the highest ANI with have been

189     described and confirmed as being strongly cellulolytic, implying that the *C. changbaiensis*

190     genome would also encode for a glucan degradation locus (GDL). Despite the high level of

191     ANIb, based on the open *Caldicellulosiruptor* pangenome, we expected to find new genes

192     involved in carbohydrate metabolism and possibly GDL arrangements.

193          ***C. changbaiensis* exhibits different abilities to grow on polysaccharides versus *C.***

194     ***bescii*.** In order to benchmark the ability of *C. changbaiensis* to grow on plant-related

8

195    polysaccharides, we compared its doubling time during exponential growth on representative

196    plant polysaccharides to *C. bescii* (**Table 3**). Doubling times (generation time) were calculated

197    from cell densities measured during exponential growth. Overall, *C. changbaiensis* grows slower

198    on microcrystalline cellulose than *C. bescii,* with a 38% larger doubling time during growth on

199    crystalline cellulose, however, both cultures grew at similar rates on xylan. On both

200    glucomannan, and pectin, *C. changbaiensis* grew faster with 35% lower doubling times (**Table**

201    **3**). The differential ability of *C. changbaiensis* and *C. bescii* to grow on pectin and glucomannan

202    is not unexpected, as the differential ability from one species to another to hydrolyze and

203    metabolize plant biomass, comprised of polysaccharides such as xylan, pectin and

204    glucomannan, was previously observed, in one case *C. saccharolyticus* grew slower on plant

205    biomass versus *C. bescii* [62] and *C. kronotskyensis* [66] and another observation where *C.*

206    *danielii* grew approximately 50% faster than *C. bescii*, *C. morganii* and *C. naganoensis* on plant

207    biomass [36].

208        When comparing the genomes of *C. changbaiensis* and *C. bescii*, *C. changbaiensis*

209    encodes for 411 genes not shared with *C. bescii*, 120 of which are unique to the genus. We

210    expect that the differences in growth rates on carbohydrates to be related to differences in gene

211    inventory. In fact, the *C. changbaiensis* gene inventory encoding for carbohydrate active

212    enzymes includes 13 genes not found in the *C. bescii* genome, including an annotated β-

213    mannanase (glycoside hydrolase [GH] family 26) and two mannooligosaccharide

214    phosphorylases (GH130). This additional β-mannanase and phosphorylases likely contribute to

215    the enhanced growth of *C. changbaiensis* on glucomannan (**Table 3**).

216        The lower doubling time on pectin is surprising, however, given that *C. changbaiensis*

217    does not encode for the pectinase cluster that is located in the *C. bescii* genome immediately

218    downstream of the GDL. *C. bescii* gene deletion strains lacking the pectinase cluster were

219    impaired in their growth on both pectin-rich plant biomass and pectin [17], indicating that *C.*

220    *changbaiensis* has evolved alternate mechanisms to deconstruct or metabolize pectin.

221    Screening the *C. changbaiensis* genome for pectin-related enzymes did not identify any genes

222    encoding for polysaccharide lyases (PL) that were unique in comparison to *C. bescii*, however

223    genes encoding for representatives from GH family 43, 51 (α-L-arabinofuranosidases) and 95 (α

224    -fucosidase) were present. One scenario is that these enzymes participate in the hydrolysis of

225    carbohydrate sidechains from pectin [44]. Another plausible explanation is that *C.*

226    *changbaiensis* has evolved to import and efficiently ferment a broader range of carbohydrates

227    released during growth on plant biomass, including uronic acids, and/ or the deoxy sugars

228    fucose and rhamnose. While *C. bescii* may rely on its enzymatic repertoire to deconstruct plant

229    biomass, it may not metabolize all types of carbohydrates that are released, similar to *R.*

230    *thermocellum* which produces xylanases, but does not metabolize xylose [42,43].

231            **Organization of the *C. changbaiensis* genome degradation locus.** *C. changbaiensis*

232    was originally described as strongly cellulolytic [3] and accordingly, its genome encodes for a

233    GDL that shares a similar organization with other strongly cellulolytic members of the genus.

234    since *C. bescii* was able to grow at a faster rate on microcrystalline cellulose than *C.*

235    *changbaiensis* (**Table** 3), we opted to focus on the comparison of GDL between these two

236    species. The GDL from both species is remarkably similar, with only CelD possessing a different

237    arrangement of catalytic and non-catalytic domains (GH10-CBM3-GH5) from *C. changbaiensis*,

238    and truncated versions of CelE (GH9-CBM3-GH5) and CelF (GH74-CBM3) present (**Fig. 3**).

239    Prior *in vitro* biochemical analyses on the synergy of cellulase mixtures from *C. bescii* had

240    observed that a mixture of three cellulases, CelA, CelC and CelE (ACE cellulases) worked

241    synergistically to hydrolyze cellulose as well as a mixture of all six *C. bescii* cellulases [21]. One

242    could hypothesize, then, that members of the genus *Caldicellulosiruptor* that possess all three of

243    these enzymes would be among the most cellulolytic. Three additional species, *C.*

244    *kronotskyensis*, *C. danielii*, and *C. naganoensis* also share a similar organization of their GDL

245    [36], including the presence of CelA, CelC and CelE. The contributions of CelD and CelF to

246    cellulose hydrolysis or solubilization are low [22,21] and likely not to impact the ability of *C.*

247    *changbaiensis* to efficiently hydrolyze cellulose.

248        Indeed, *C. changbaiensis* can solubilize microcrystalline cellulose (**Fig. 4**), however the

249    amount of cellulose solubilized was 22.4% lower than the amount solubilized by *C. bescii*, which

250    is similar to the performance of *C. saccharolyticus* when compared to *C. bescii* [24.8% lower,

251    66]. This result begs the question if the mere presence of the ACE cellulases is sufficient to

252    meet the *C. bescii* benchmark for hydrolysis of cellulose. One explanation could be that the *C.*

253    *changbaiensis* CelE ortholog may not be as efficient in cellulose hydrolysis since it is lacking

254    two CBM3 modules. However, the nearly equal reduction of cellulose solubilization by both *C.*

255    *bescii* gene deletion strains incapable of producing CelA-CelC versus CelA-CelE does not

256    support this possibility [22]. Furthermore, CelE truncations that possessed the GH9 catalytic

257    domain and three or two CBM3 domains were equally capable of microcrystalline cellulose

258    hydrolysis [53], making it unlikely that the loss of a CBM3 domain from the *C. changbaiensis*

259    CelA ortholog hampered its activity.

260        Alternately, sequence divergence of ACE cellulase orthologs may play a larger role in

261    the catalytic capacity of cellulolytic members from the genus *Caldicellulosiruptor.* Of the ACE

262    cellulases, CelA is a key player, supported by its unique hydrolysis mechanism [8], the severe

263    reduction in cellulose hydrolysis by *C. bescii* celA gene deletion mutant [65,22], and biochemical

264    analysis of GDL enzyme synergy [21]. Prior comparison of CelA orthologs from *C. bescii* and *C.*

265    *danielii* found *Cb*CelA to be a superior enzyme [36], indicating that GDL sequences have

266    diverged during speciation, making it likely that the ACE cellulases from *C. changbaiensis* may

267    not demonstrate the same catalytic efficiency as *C. bescii*.

268        **Attachment of *C. bescii* and *C. changbaiensis* to plant polysaccharides.** Aside from

269    comparisons of catalytic ability, we also compared the ability of *C. changbaiensis* versus *C.*

270    *bescii* planktonic cells to bind to insoluble substrates (xylan and cellulose). A decrease in the

271    planktonic cell density (PCD) after exposure to the substrate compared to the PCD of the

11

272   negative controls without substrate is indicative of cells binding to the substrate. Surprisingly, we

273   saw no such decrease in PCD for *C. changbaiensis* cultured on xylan after incubation with

274   cellulose or xylan (**Figs. 5A and B**). This inability of *C. changbaiensis* to attach to xylan or

275   cellulose after growth on xylan is surprising, given that xylan is a major polysaccharide

276   constituent of lignocellulose, and would likely serve as a chemical signal. Since no xylan or

277   cellulose attachment proteins are produced in response to growth on xylan, *C. changbaiensis*

278   appears to act as a specialist, responding only to cellulose. Regardless, when *C. changbaiensis*

279   is grown on cellulose, it maintains an ability to attach to cellulose (29% cells attached), which is

280   slightly lower than the relative amount of *C. bescii* cells attached to cellulose (33% attached,

281   **Fig. 5C**). Surprisingly, when *C. bescii* cells cultured on xylan were tested for attachment to

282   either xylan or cellulose there was a significant decrease in (PCD) of indicating that *C. bescii*

283   cells grown on xylan are producing proteins capable of attaching to xylan (33% attachment, **Fig.**

284   **5A**) or cellulose (68% attachment, **Fig. 5B**). While we expected to see cells from cultures grown

285   on xylan attaching to xylan, interestingly, *C. bescii* cell attachment was most pronounced when

286   cells were grown on xylan and incubated with cellulose (**Fig. 5B**). The ability of *C. bescii* to

287   attach to cellulose (**Figs. 5B and C**), is in large part due to the presence of tāpirins, since a *C.*

288   *bescii* tāpirin deletion mutant was severely impaired in cellular attachment to cellulose [37].

289          **The *C. changbaiensis* genome encodes for atypical tāpirin genes.** Another notable

290   difference observed between *C. changbaiensis* and *C. bescii* during growth on cellulose is the

291   lack of floc formation by *C. changbaiensi*s (**Fig. 6**). Based on this discrepancy between *C.*

292   *changbaiensis* and *C. bescii***,** we examined the genomic context of the type IV pilus locus

293   encoded by the *C. changbaiensis* genome (**Fig. 7**). The T4P locus is found in the genome in all

294   members of the *Caldicellulosiruptor*, and is also located upstream of the GDL in the genomes of

295   strongly cellulolytic species [5,4]. Most notably, while a full T4P locus is present in the *C.*

296   *changbaiensis* genome, classical tāpirin genes are absent which encode for proteins that bind

297   with high affinity to cellulose [4,37]. Instead, two genes with little, to no homology to the classical

298    tāpirins are located directly downstream of the T4P locus which we will refer to as atypical

299    tāpirins. The proteins encoded for by these genes are not unique to *C. changbaiensis*, as both

300    *C. acetigenus* and *C. ownesensis* also encode for these atypical tāpirins. All three species

301    encode for two atypical tāpirins: a hypothetical protein (Genbank accession: WP_127352232.1)

302    and a von Willebrand Factor A protein (Genbank accession: WP_127352233.1) (yellow arrows,

303    **Fig. 7**). While *C. changbaiensis* shares a similar genomic context at the 3' end of the T4P locus,

304    the atypical *C. changbaiensis* tāpirins are not close orthologs, as they share 74.33% and

305    68.01% amino sequence similarity with the first and second atypical tāpirins encoded *C.*

306    *owensensis*. Prior proteomics data collected from cellulose-bound, supernatant and whole cell

307    lysate protein fractions determined that both atypical tāpirins are produced by *C. owensensis* in

308    response to cellulose [5], supporting their potential role in cell attachment to cellulose.

309         This observed sequence divergence between the atypical tāpirins from strongly and

310    weakly cellulolytic species is similar to the tāpirin encoded for by *C. hydrothermalis* which

311    shares little amino acid sequence homology with classical tāpirins, but shares a similar tertiary

312    structure, and is capable of occupying more sites on crystalline cellulose in comparison to

313    classical tāpirins [37]. Production of tāpirins with an affinity to cellulose likely plays a role in the

314    ability of weakly cellulolytic members of the genus to adhere to cellulose and benefit from the

315    cellooligosaccharides released by the action of cellulases [60]. The atypical tāpirins, originally

316    only observed in the genomes of weakly cellulolytic species, may also serve as cellulose

317    adhesins, however, further in-depth biochemical characterization of both atypical tāpirin proteins

318    is required to confirm their function.

**CONCLUSIONS**

Overall, the *Caldicellulosiruptor* pangenome remains open, and is expected to gain approximately 63 new genes with each additional species sequenced (**Fig. 1A**). The addition of a second species isolated from China indicates that the diversity of *Caldicellulosiruptor* species from this region is higher than those isolated from Iceland, however, the level of observed diversity is not as high as those species isolated from Kamchatka, Russia or New Zealand on the basis of ANIb (**Fig. 2**). *C. changbaiensis* encodes for a GDL (**Fig. 3**) similar in organization as *C. bescii*, however is not as cellulolytic as *C. bescii* on the basis of doubling time (**Table 3**) and cellulose solubilization (**Fig. 4**). However, *C. changbaiensis* does appear to have a broader metabolic appetite for uronic acids or deoxy sugars. *C. changbaiensis* also fails to form a floc during growth on microcrystalline cellulose (**Fig. 6**), a phenotype previously described for C. bescii [64], however both species are capable of attaching to cellulose (**Fig. 5**). Interestingly, *C. bescii* retains an ability to attach to cellulose when previously grown on xylan, while *C. changbaiensis* does not (**Fig. 5B**) indicating that the two species respond differently to soluble carbohydrates present in their environment. Tāpirins were previously demonstrated to be key cellulose adhesins for strongly [4] to weakly cellulolytic [37] members of the genus *Caldicellulosiruptor*. Surprisingly, *C. changbaiensis* does not encode for the classical tāpirins, and instead encodes for atypical tāpirins, one of which possesses a von Willebrand type A protein domain (**Fig. 7**). These atypical tāpirins are homologous to those encoded for by weakly cellulolytic *C. owensensis* and *C. acetigenus*, however this may not indicate that the atypical tāpirins are not involved in attachment to cellulose, as the divergent classical tāpirin encoded for by *C. hyrothermalis* binds at a high density to cellulose [37]. The combined lack of classical tāpirins, along with the ability to attach to cellulose indicates that *C. changbaiensis* evolved a unique strategy to attach to cellulose. Further study on the biophysical properties of these atypical tāpirins is warranted to assess their ability to interact with plant polysaccharides, including cellulose.

14

345 **FIGURE LEGENDS**

346 **Figure 1. Core- and pangenome size estimates calculated from random sampling of 14**

347 ***Caldicellulosiruptor* genomes. (a)** Fitted curve of the estimated *Caldicellulosiruptor* core

348 genome from 10 random samples of genomes up to n=14. The current size of the core genome

349 is 1367 orthologous clusters. **(b)** Fitted curve of the estimated *Caldicellulosiruptor* pangenome

350 from 10 random samples of genomes up to n=14. The *Caldicellulosiruptor* pangenome remains

351 open and has increased to 3791 genes. The rate of growth for the pangenome is 63.2 new

352 genes per genome sequenced. Core- and pangenome estimates were calculated from the

353 equations reported by Tettelin *et al.*, [58] using GET_HOMOLOGUES software [19].

354

355 **Figure 2. Heatmap representation of the average nucleotide identity for 14 genome**

356 **sequenced species from the genus *Caldicellulosiruptor.*** Average nucleotide identity (ANIb)

357 was calculated on the basis of legacy BLASTn sequence identity over 1020nt sequence

358 fragments. ANIb values of all 14 genomes are represented by a heat plot ranging from blue

359 (75%< ANIb <90%), white (90%< ANIb <95%) to red (ANIb >95%). Pyani

360 (https://github.com/widdowquinn/pyani) was used to calculate ANIb values and generate the

361 clustered heatmap. Hierarchal cluster dendrograms were generated on the basis of similar ANIb

362 values across each species. ANIb values are reported in Table S1. Calace, *C. acetigenus*;

363 Cbes, *C. bescii*; Calcha, *C. changbaiensis*; Caldan, *C. danielii*; Calhy, *C. hydrothermalis*; Calkr,

364 *C. kristjanssonii*; Calkro, *C. kronotskyensis*; Calla, *C. lactoaceticus*; Calmo, *C. morganii*; Calna,

365 *C. naganoensis*; COB47, *C. obsidiansis*; Calow, *C. owensensis*; Csac, *C. saccharolyticus*; F32,

366 *C.* sp. F32.

367

368 **Figure 3. Modular multifunctional enzymes encoded for by the glucan degradation locus.**

369 Glucan degradation loci were selected on the basis of the presence of "ACE" cellulases. ACE

370 cellulases: CelA, CelC and CelE. Circles represent the glycoside hydrolase (GH) domains,

15

371    rectangles represent the carbohydrate binding module (CBM) domains. GH5, green circles;

372    GH9, red circles; GH10, violet circles; GH 44, blue circles; GH48, grey circles; GH74, orange

373    circles. CBM3, grey rectangles; CBM22, pink rectangles.

374

375    **Figure 4. Solubilization of microcrystalline cellulose by *C. bescii* and *C. changbaiensis*.**

376    Uninoculated control, indicates abiotic cellulose solubilization in LOD medium. Error bars

377    represent standard error (n=3). Similar letters over columns denote $p < 0.05$ as determined by a

378    t-test.

379

380    **Figure 5. Comparison of the ability of *C. bescii* or *C. changbaiensis* planktonic cells to**

381    **attach to polysaccharides.** Titles above bar charts indicate the carbon source for growth/

382    binding substrate. **(a, b)** When cells are grown on xylan, only planktonic *C. bescii* cells were

383    able to attach to xylan or cellulose. **(c)** Cells grown on cellulose as the carbon source and

384    exposed to cellulose as the binding substrate. Planktonic cell densities (PCD), enumerated by

385    epifluorescence microscopy are plotted on the y-axis. Green columns indicate PCD without

386    binding substrate and purple columns indicate PCD with the binding substrate. * indicates $p <$

387    0.01 as determined by a t-test. All assays had n=6 biological replicates.

388

389    **Figure 6. Flocculation of *C. bescii* cells cultured on chemically defined medium and**

390    **microcrystalline cellulose. (a)** Formation of a floc of *C. bescii* cells around microcrystalline

391    cellulose (diameter, 20µm) while planktonic *C. changbaiensis* cells (cloudiness) are visible. **(b)**

392    Same serum bottles as in "A", however the bottles were vigorously mixed. The *C. bescii* floc

393    remains fairly stable, while both microcrystalline cellulose and cells are mixed in the *C.*

394    *changbaiensis* culture.

395

16

396     **Figure 7. Genomic context for the location of the tāpirins from strongly to weakly**

397     **cellulolytic *Caldicellulosiruptor* species.** Different colors represent the classical versus

398     atypical tāpirins. Blue arrows: Cbes tāpirin 1 (Gen bank accession: YP_002573732) and Cbes

399     tāpirin 2 (Gen bank accession: YP_002573731). Green arrow: Calhy tāpirin 1 (Gen bank

400     accession number: YP_003992006). Yellow arrows: Calcha tāpirin 1 (Gen bank accession:

401     WP_127352232.1) and 2 (Gen bank accession: WP_127352233.1), and Calow tāpirin 1 (Gen

402     bank accession: YP_004002936) and 2 (Gen bank accession r: YP_004002935). Grey

403     rectangles indicate the presence of the GDL downstream of the tāpirins. Atypical tāpirin 1 is

404     annotated as a hypothetical protein and atypical tāpirin 2 is annotated as a von Willebrand

405     factor A protein. Cbes, *C. bescii*; Calhy, *C. hydrothermalis*; Calcha, *C. changbaiensis* and

406     Calow, *C. owensensis*. Peach rectangles represent the type IV pilus locus directly upstream of

407     the tāpirins. Arrows indicate tāpirin 1 and 2. Numbers in the tāpirin arrows indicate the amino

408     acid length.

17

409   **Conflict of Interest:** A.M.A.M., C.M., V.J.H. and S.E.B.-S. declare that they have no conflict of

410   interest.

411

| Table 1. Oligonucleotide primers used for *Caldicellulosiruptor* 16S rRNA gene fragment amplification | | 412 |
|---|---|---|
| **Primer Name** | **Primer Sequence (5' to 3')** | **Source** |
| 8F-207 | AGAGTTTGATCCTGGCTCAG | [3] |
| Caldi-R-208 | GTACGGCTACCTTGTTACG | |

| Table 2. *Caldicellulosiruptor* genome sequences included in the updated pangenome analysis | | |
|---|---|---|
| **Species Name** | **NCBI RefSeq Accession** | **Reference** |
| *C. acetigenus* | GCF_000421725.1 | [41] |
| *C. bescii* | GCF_000022325.1 | [22] |
| *C. changbaiensis* | GCF_003999255.1 | [40] |
| *C. danielii* | GCF_000955725.1 | [38,36] |
| *C. hydrothermalis* | GCF_000166355.1 | [7,5] |
| *C. kristjanssonii* | GCF_000166695.1 | [7,5] |
| *C. kronotskyensis* | GCF_000166775.1 | [7,5] |
| *C. lactoaceticus* | GCF_000193435.2 | [7,5] |
| *C. morganii* | GCF_000955745.1 | [38,36] |
| *C. naganoensis* | GCF_000955735.1 | [38,36] |
| *C. obsidiansis* | GCF_000145215.1 | [24] |
| *C. owensensis* | GCF_000166335.1 | [7,5] |
| *C. saccharolyticus* | GCF_000016545.1 | [59] |
| *C.* str. F32 | GCF_000404025.1 | [63] |

413

414

20

| Table 3. Doubling time of *C. changbaiensis* or *C. bescii* grown on plant polysaccharides | | |
|---|---|---|
| **Polysaccharide** | **g$_{Cbes}$ (hr)** | **g$_{Calcha}$ (hr)** |
| Microcrystalline cellulose | 3.93 ± 0.157 | 5.43 ± 0.304 |
| Beechwood xylan | 2.55 ± 0.211 | 2.54 ± 0.428 |
| Glucomannan | 3.22 ± 0.62 | 2.08 ± 0.025 |
| Pectin | 3.48 ± 0.224 | 2.26 ± 0.167 |

415                                    References
416

417   1. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ (1997)

418       Gapped BLAST and PSI-BLAST: a new generation of protein database search

419       programs. Nucleic Acids Res 25:3389-3402

420   2. Artzi L, Bayer EA, Morais S (2017) Cellulosomes: bacterial nanomachines for dismantling

421       plant polysaccharides. Nat Rev Microbiol 15:83-95. doi:10.1038/nrmicro.2016.164

422   3. Bing W, Wang H, Zheng B, Zhang F, Zhu G, Feng Y, Zhang Z (2015) *Caldicellulosiruptor*

423       *changbaiensis* sp. nov., a cellulolytic and hydrogen-producing bacterium from a hot

424       spring. Int J Syst Evol Microbiol 65:293-297. doi:10.1099/ijs.0.065441-0

425   4. Blumer-Schuette SE, Alahuhta M, Conway JM, Lee LL, Zurawski JV, Giannone RJ, Hettich

426       RL, Lunin VV, Himmel ME, Kelly RM (2015) Discrete and structurally unique proteins

427       (tāpirins) mediate attachment of extremely thermophilic *Caldicellulosiruptor* species to

428       cellulose. J Biol Chem 290:10645-10656. doi:10.1074/jbc.M115.641480

429   5. Blumer-Schuette SE, Giannone RJ, Zurawski JV, Ozdemir I, Ma Q, Yin Y, Xu Y, Kataeva I,

430       Poole FL, 2nd, Adams MW, Hamilton-Brehm SD, Elkins JG, Larimer FW, Land ML,

431       Hauser LJ, Cottingham RW, Hettich RL, Kelly RM (2012) *Caldicellulosiruptor* core and

432       pangenomes reveal determinants for noncellulosomal thermophilic deconstruction of

433       plant biomass. J Bacteriol 194:4015-4028. doi:10.1128/JB.00266-12

434   6. Blumer-Schuette SE, Lewis DL, Kelly RM (2010) Phylogenetic, microbiological, and glycoside

435       hydrolase diversities within the extremely thermophilic, plant biomass-degrading genus

436       *Caldicellulosiruptor*. Appl Environ Microbiol 76:8084-8092. doi:10.1128/AEM.01400-10

437   7. Blumer-Schuette SE, Ozdemir I, Mistry D, Lucas S, Lapidus A, Cheng JF, Goodwin LA,

438       Pitluck S, Land ML, Hauser LJ, Woyke T, Mikhailova N, Pati A, Kyrpides NC, Ivanova N,

439       Detter JC, Walston-Davenport K, Han S, Adams MW, Kelly RM (2011) Complete

440       genome sequences for the anaerobic, extremely thermophilic plant biomass-degrading

441       bacteria *Caldicellulosiruptor hydrothermalis*, *Caldicellulosiruptor kristjanssonii*,

22

442    *Caldicellulosiruptor kronotskyensis*, *Caldicellulosiruptor owensensis*, and

443    *Caldicellulosiruptor lactoaceticu*s. J Bacteriol 193:1483-1484. doi:10.1128/JB.01515-10

444    8. Brunecky R, Alahuhta M, Xu Q, Donohoe BS, Crowley MF, Kataeva IA, Yang SJ, Resch MG,

445    Adams MW, Lunin VV, Himmel ME, Bomble YJ (2013) Revealing nature's cellulase

446    diversity: the digestion mechanism of *Caldicellulosiruptor bescii* CelA. Science

447    342:1513-1516. doi:10.1126/science.1244273

448    9. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL (2009)

449    BLAST+: architecture and applications. BMC Bioinformatics 10:421. doi:10.1186/1471-

450    2105-10-421

451    10. Cha M, Chung D, Elkins JG, Guss AM, Westpheling J (2013) Metabolic engineering of

452    *Caldicellulosiruptor bescii* yields increased hydrogen production from lignocellulosic

453    biomass. Biotechnol Biofuel 6:1-8. doi:10.1186/1754-6834-6-85

454    11. Chung D, Cha M, Farkas J, Westpheling J (2013) Construction of a stable replicating shuttle

455    vector for *Caldicellulosiruptor* species: Use for extending genetic methodologies to other

456    members of this genus. PLoS One 8:e62881. doi:10.1371/journal.pone.0062881

457    12. Chung D, Cha M, Guss AM, Westpheling J (2014) Direct conversion of plant biomass to

458    ethanol by engineered *Caldicellulosiruptor bescii*. Proceedings of the National Academy

459    of Sciences:201402210. doi:10.1073/pnas.1402210111

460    13. Chung D, Cha M, Snyder EN, Elkins JG, Guss AM, Westpheling J (2015) Cellulosic ethanol

461    production via consolidated bioprocessing at 75 °C by engineered *Caldicellulosiruptor*

462    *bescii*. Biotechnol Biofuel 8:163. doi:10.1186/s13068-015-0346-4

463    14. Chung D, Farkas J, Huddleston JR, Olivar E, Westpheling J (2012) Methylation by a unique

464    α-class N4-cytosine methyltransferase Is required for DNA transformation of

465    *Caldicellulosiruptor bescii* DSM6725. PLoS One 7:e43844.

466    doi:10.1371/journal.pone.0043844

467    15. Chung D, Farkas J, Westpheling J (2013) Detection of a novel active transposable element

468        in *Caldicellulosiruptor hydrothermalis* and a new search for elements in this genus. J Ind

469        Microbiol Biotechnol 40:517-521. doi:10.1007/s10295-013-1244-z

470    16. Chung D, Farkas J, Westpheling J (2013) Overcoming restriction as a barrier to DNA

471        transformation in *Caldicellulosiruptor* species results in efficient marker replacement.

472        Biotechnol Biofuel 6:1-9. doi:10.1186/1754-6834-6-82

473    17. Chung D, Pattathil S, Biswal AK, Hahn MG, Mohnen D, Westpheling J (2014) Deletion of a

474        gene cluster encoding pectin degrading enzymes in *Caldicellulosiruptor bescii* reveals an

475        important role for pectin in plant biomass recalcitrance. Biotechnol Biofuel 7:147.

476        doi:10.1186/s13068-014-0147-1

477    18. Chung D, Young J, Cha M, Brunecky R, Bomble YJ, Himmel ME, Westpheling J (2015)

478        Expression of the *Acidothermus cellulolyticus* E1 endoglucanase in *Caldicellulosiruptor*

479        *bescii* enhances its ability to deconstruct crystalline cellulose. Biotechnol Biofuel 8:113.

480        doi:10.1186/s13068-015-0296-x

481    19. Contreras-Moreira B, Vinuesa P (2013) GET_HOMOLOGUES, a versatile software package

482        for scalable and robust microbial pangenome analysis. Appl Environ Microbiol 79:7696-

483        7701. doi:10.1128/AEM.02411-13

484    20. Conway JM, Crosby JR, Hren AP, Southerland RT, Lee LL, Lunin VV, Alahuhta P, Himmel

485        ME, Bomble YJ, Adams MWW, Kelly RM (2018) Novel multidomain, multifunctional

486        glycoside hydrolases from highly lignocellulolytic *Caldicellulosiruptor* species. AIChE J

487        64:4218-4228. doi:10.1002/aic.16354

488    21. Conway JM, Crosby JR, McKinley BS, Seals NL, Adams MWW, Kelly RM (2018) Parsing in

489        vivo and in vitro contributions to microcrystalline cellulose hydrolysis by multidomain

490        glycoside hydrolases in the *Caldicellulosiruptor bescii* secretome. Biotechnol Bioeng 0.

491        doi:10.1002/bit.26773

492    22. Conway JM, McKinley BS, Seals NL, Hernandez D, Khatibi PA, Poudel S, Giannone RJ,

493         Hettich RL, Williams-Rhaesa AM, Lipscomb GL, Adams MWW, Kelly RM (2017)

494         Functional analysis of the glucan degradation locus in *Caldicellulosiruptor bescii* reveals

495         essential roles of component glycoside hydrolases in plant biomass deconstruction. Appl

496         Environ Microbiol 83:e01828-01817. doi:10.1128/AEM.01828-17

497    23. Dam P, Kataeva I, Yang S-J, Zhou F, Yin Y, Chou W, Poole FL, Westpheling J, Hettich R,

498         Giannone R, Lewis DL, Kelly R, Gilbert HJ, Henrissat B, Xu Y, Adams MWW (2011)

499         Insights into plant biomass conversion from the genome of the anaerobic thermophilic

500         bacterium *Caldicellulosiruptor bescii* DSM 6725. Nucleic Acids Res 39:3240 -3254.

501         doi:10.1093/nar/gkq1281

502    24. Elkins JG, Lochner A, Hamilton-Brehm SD, Davenport KW, Podar M, Brown SD, Land ML,

503         Hauser LJ, Klingeman DM, Raman B, Goodwin LA, Tapia R, Meincke LJ, Detter JC,

504         Bruce DC, Han CS, Palumbo AV, Cottingham RW, Keller M, Graham DE (2010)

505         Complete genome sequence of the cellulolytic thermophile *Caldicellulosiruptor*

506         *obsidiansis* OB47T. J Bacteriol 192:6099-6100. doi:10.1128/JB.00950-10

507    25. Farkas J, Chung D, Cha M, Copeland J, Grayeski P, Westpheling J (2013) Improved growth

508         media and culture techniques for genetic analysis and assessment of biomass utilization

509         by *Caldicellulosiruptor bescii*. J Ind Microbiol Biotechnol 40:1-9. doi:10.1007/s10295-

510         012-1202-1

511    26. Gibbs MD, Saul DJ, Lüthi E, Bergquist PL (1992) The beta-mannanase from" *Caldocellum*

512         *saccharolyticum*" is part of a multidomain enzyme. Appl Environ Microbiol 58:3864–3867

513    27. Goris J, Konstantinidis KT, Klappenbach JA, Coenye T, Vandamme P, Tiedje JM (2007)

514         DNA-DNA hybridization values and their relationship to whole-genome sequence

515         similarities. Int J Syst Evol Microbiol 57:81-91. doi:10.1099/ijs.0.64483-0

516    28. Hobbie JE, Daley RJ, Jasper S (1977) Use of nuclepore filters for counting bacteria by

517         fluorescence microscopy. Appl Environ Microbiol 33:1225-1228

29. Kahn A, Moraïs S, Galanopoulou AP, Chung D, Sarai NS, Hengge N, Hatzinikolaou DG, Himmel ME, Bomble YJ, Bayer EA (2019) Creation of a functional hyperthermostable designer cellulosome. Biotechnol Biofuel 12:44. doi:10.1186/s13068-019-1386-y

30. Kim S-K, Chung D, Himmel ME, Bomble YJ, Westpheling J (2016) Heterologous expression of family 10 xylanases from *Acidothermus cellulolyticus* enhances the exoproteome of *Caldicellulosiruptor bescii* and growth on xylan substrates. Biotechnol Biofuel 9. doi:10.1186/s13068-016-0588-9

31. Kim S-K, Chung D, Himmel ME, Bomble YJ, Westpheling J (2017) Heterologous expression of a β-d-glucosidase in *Caldicellulosiruptor bescii* has a surprisingly modest effect on the activity of the exoproteome and growth on crystalline cellulose. J Ind Microbiol Biotechnol 44:1643-1651. doi:10.1007/s10295-017-1982-4

32. Kim S-K, Chung D, Himmel ME, Bomble YJ, Westpheling J (2017) In vivo synergistic activity of a CAZyme cassette from A*cidothermus cellulolyticus* significantly improves the cellulolytic activity of the *C. bescii* exoproteome. Biotechnol Bioeng 114:2474-2480. doi:10.1002/bit.26366

33. Kim S-K, Chung D, Himmel ME, Bomble YJ, Westpheling J (2019) Heterologous co-expression of two β-glucanases and a cellobiose phosphorylase resulted in a significant increase in the cellulolytic activity of the *Caldicellulosiruptor bescii* exoproteome. J Ind Microbiol Biotechnol. doi:10.1007/s10295-019-02150-0

34. Kim S-K, Himmel ME, Bomble YJ, Westpheling J (2018) Expression of a cellobiose phosphorylase from *Thermotoga maritima* in *Caldicellulosiruptor bescii* improves the phosphorolytic pathway and results in a dramatic increase in cellulolytic activity. Appl Environ Microbiol 84:e02348-02317. doi:10.1128/AEM.02348-17

35. Kristensen DM, Kannan L, Coleman MK, Wolf YI, Sorokin A, Koonin EV, Mushegian A (2010) A low-polynomial algorithm for assembling clusters of orthologous groups from

543    intergenomic symmetric best matches. Bioinformatics 26:1481-1487.

544    doi:10.1093/bioinformatics/btq229

545    36. Lee LL, Blumer-Schuette SE, Izquierdo JA, Zurawski JV, Loder AJ, Conway JM, Elkins JG,

546    Podar M, Clum A, Jones PC, Piatek MJ, Weighill DA, Jacobson DA, Adams MWW, Kelly

547    RM (2018) Genus-wide assessment of lignocellulose utilization in the extremely

548    thermophilic genus *Caldicellulosiruptor* by genomic, pangenomic, and metagenomic

549    analyses. Appl Environ Microbiol 84. doi:10.1128/AEM.02694-17

550    37. Lee LL, Hart WS, Lunin VV, Alahuhta M, Bomble YJ, Himmel ME, Blumer-Schuette SE,

551    Adams MWW, Kelly RM (2019) Comparative biochemical and structural analysis of

552    novel cellulose binding proteins (tāpirins) from extremely thermophilic

553    *Caldicellulosiruptor* species. Appl Environ Microbiol 85:e01983-01918.

554    doi:10.1128/AEM.01983-18

555    38. Lee LL, Izquierdo JA, Blumer-Schuette SE, Zurawski JV, Conway JM, Cottingham RW,

556    Huntemann M, Copeland A, Chen IM, Kyrpides N, Markowitz V, Palaniappan K, Ivanova

557    N, Mikhailova N, Ovchinnikova G, Andersen E, Pati A, Stamatis D, Reddy TB, Shapiro

558    N, Nordberg HP, Cantor MN, Hua SX, Woyke T, Kelly RM (2015) Complete genome

559    sequences of *Caldicellulosiruptor* sp. strain Rt8.B8, *Caldicellulosiruptor* sp. strain

560    Wai35.B1, and "*Thermoanaerobacter cellulolyticus*". Genome Announc 3.

561    doi:10.1128/genomeA.00440-15

562    39. Li L, Stoeckert CJ, Jr., Roos DS (2003) OrthoMCL: identification of ortholog groups for

563    eukaryotic genomes. Genome Res 13:2178-2189. doi:10.1101/gr.1224503

564    40. Mendoza C, Blumer-Schuette SE (2019) Complete genome sequence of *Caldicellulosiruptor*

565    *changbaiensis* CBS-Z, an extremely thermophilic, cellulolytic bacterium isolated from a

566    hot spring in China. Microbiol Resour Announc 8. doi:10.1128/MRA.00021-19

567    41. Mukherjee S, Seshadri R, Varghese NJ, Eloe-Fadrosh EA, Meier-Kolthoff JP, Goker M,

568    Coates RC, Hadjithomas M, Pavlopoulos GA, Paez-Espino D, Yoshikuni Y, Visel A,

569       Whitman WB, Garrity GM, Eisen JA, Hugenholtz P, Pati A, Ivanova NN, Woyke T, Klenk

570       HP, Kyrpides NC (2017) 1,003 reference genomes of bacterial and archaeal isolates

571       expand coverage of the tree of life. Nat Biotechnol 35:676-+. doi:10.1038/nbt.3886

572  42. Ng TK, Ben-Bassat A, Zeikus JG (1981) Ethanol production by thermophilic bacteria:

573       Fermentation of cellulosic substrates by cocultures of *Clostridium thermocellum* and

574       *Clostridium thermohydrosulfuricum.* Appl Environ Microbiol 41:1337-1343

575  43. Ng TK, Zeikus JG (1981) Comparison of extracellular cellulase activities of *Clostridium*

576       *thermocellum* LQRI and *Trichoderma reese*i QM9414. Appl Environ Microbiol 42:231-

577       240

578  44. Numan MT, Bhosle NB (2006) Alpha-L-arabinofuranosidases: the potential applications in

579       biotechnology. J Ind Microbiol Biotechnol 33:247-260. doi:10.1007/s10295-005-0072-1

580  45. Park JI, Kent MS, Datta S, Holmes BM, Huang Z, Simmons BA, Sale KL, Sapra R (2011)

581       Enzymatic hydrolysis of cellulose by the cellobiohydrolase domain of CelB from the

582       hyperthermophilic bacterium *Caldicellulosiruptor saccharolyticus.* Bioresour Technol

583       102:5988-5994. doi:10.1016/j.biortech.2011.02.036

584  46. Park JI, Steen EJ, Burd H, Evans SS, Redding-Johnson AM, Batth T, Benke PI,

585       D'Haeseleer P, Sun N, Sale KL, Keasling JD, Lee TS, Petzold CJ, Mukhopadhyay A,

586       Singer SW, Simmons BA, Gladden JM (2012) A thermophilic ionic liquid-tolerant

587       cellulase cocktail for the production of cellulosic biofuels. PLoS One 7:e37010.

588       doi:10.1371/journal.pone.0037010

589  47. R Core T (2015) R: A language and environment for statistical computing. R Foundation for

590       Statistical Computing. Vienna, Austria

591  48. Reynolds PH, Sissons CH, Daniel RM, Morgan HW (1986) Comparison of cellulolytic

592       activities in clostridium thermocellum and three thermophilic, cellulolytic anaerobes. Appl
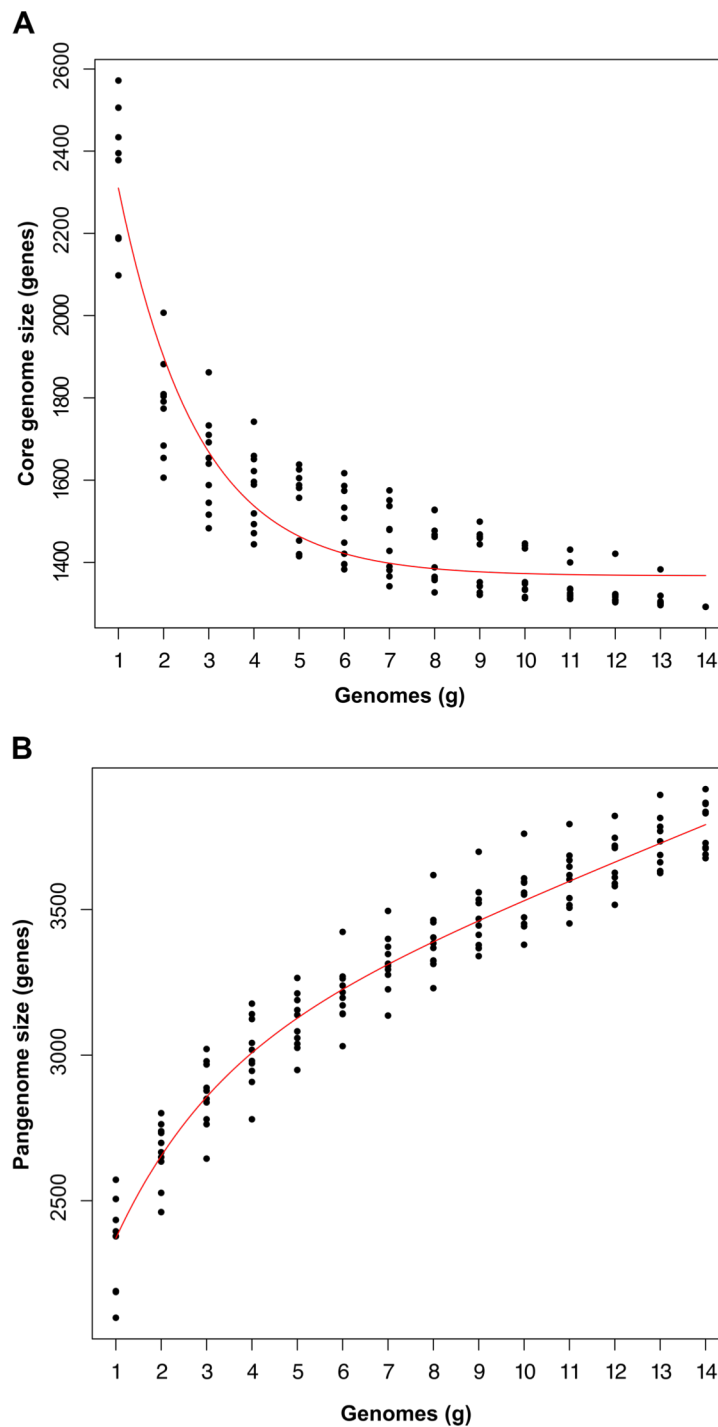
593       Environ Microbiol 51:12-17

594    49. Richter M, Rossello-Mora R (2009) Shifting the genomic gold standard for the prokaryotic

595         species definition. Proc Natl Acad Sci U S A 106:19126-19131.

596         doi:10.1073/pnas.0906412106

597    50. Sander K, Chung D, Hyatt D, Westpheling J, Klingeman DM, Rodriguez M, Engle NL,

598         Tschaplinski TJ, Davison BH, Brown SD Rex in *Caldicellulosiruptor bescii*: Novel regulon

599         members and its effect on the production of ethanol and overflow metabolites.

600         MicrobiologyOpen 0:e00639. doi:10.1002/mbo3.639

601    51. Saul DJ, Williams LC, Grayling RA, Chamley LW, Love DR, Bergquist PL (1990) celB, a

602         gene coding for a bifunctional cellulase from the extreme thermophile "*Caldocellum*

603         *saccharolyticum*". Appl Environ Microbiol 56:3117-3124

604    52. Smith SP, Bayer EA, Czjzek M (2017) Continually emerging mechanistic complexity of the

605         multi-enzyme cellulosome complex. Curr Opin Struct Biol 44:151-160.

606         doi:10.1016/j.sbi.2017.03.009

607    53. Su X, Mackie RI, Cann IKO (2012) Biochemical and mutational analyses of a multidomain

608         cellulase/mannanase from *Caldicellulosiruptor bescii*. Appl Environ Microbiol 78:2230.

609         doi:10.1128/AEM.06814-11

610    54. Svetlichnyi VA, T. P. Svetlichnaya, N. A. Chernykh, and G. A. Zavarzin (1990) *Anaerocellum*

611         *thermophilum* gen. nov sp. nov. an extremely thermophilic cellulolytic eubacterium

612         isolated from hot-springs in the Valley of Geysers. Microbiology 59:598-604

613    55. Tatusova T, DiCuccio M, Badretdin A, Chetvernin V, Nawrocki EP, Zaslavsky L, Lomsadze

614         A, Pruitt KD, Borodovsky M, Ostell J (2016) NCBI prokaryotic genome annotation

615         pipeline. Nucleic Acids Res 44:6614-6624. doi:10.1093/nar/gkw569

616    56. Taya M, Hinoki H, Kobayashi T (1985) Tungsten requirement of an extremely thermophilic,

617         cellulolytic anaerobe (strain NA10). Agric Biol Chem 49:2513-2515.

618         doi:http://dx.doi.org/10.1080/00021369.1985.10867120

619    57. Te'o VS, Saul DJ, Bergquist PL (1995) celA, another gene coding for a multidomain

620         cellulase from the extreme thermophile *Caldocellum saccharolyticum*. Appl Microbiol

621         Biotechnol 43:291-296

622    58. Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, Angiuoli SV, Crabtree

623         J, Jones AL, Durkin AS, Deboy RT, Davidsen TM, Mora M, Scarselli M, Margarit y Ros I,

624         Peterson JD, Hauser CR, Sundaram JP, Nelson WC, Madupu R, Brinkac LM, Dodson

625         RJ, Rosovitz MJ, Sullivan SA, Daugherty SC, Haft DH, Selengut J, Gwinn ML, Zhou L,

626         Zafar N, Khouri H, Radune D, Dimitrov G, Watkins K, O'Connor KJ, Smith S, Utterback

627         TR, White O, Rubens CE, Grandi G, Madoff LC, Kasper DL, Telford JL, Wessels MR,

628         Rappuoli R, Fraser CM (2005) Genome analysis of multiple pathogenic isolates of

629         *Streptococcus agalactiae*: implications for the microbial "pan-genome". Proc Natl Acad

630         Sci U S A 102:13950-13955. doi:10.1073/pnas.0506758102

631    59. van de Werken HJG, Verhaart MRA, VanFossen AL, Willquist K, Lewis DL, Nichols JD,

632         Goorissen HP, Mongodin EF, Nelson KE, van Niel EWJ, Stams AJM, Ward DE, de Vos

633         WM, van der Oost J, Kelly RM, Kengen SWM (2008) Hydrogenomics of the extremely

634         thermophilic bacterium *Caldicellulosiruptor saccharolyticus*. Appl Environ Microbiol

635         74:6720-6729. doi:10.1128/Aem.00968-08

636    60. Wang Z-W, Hamilton-Brehm SD, Lochner A, Elkins JG, Morrell-Falvey JL (2011)

637         Mathematical modeling of hydrolysate diffusion and utilization in cellulolytic biofilms of

638         the extreme thermophile *Caldicellulosiruptor obsidiansis*. Bioresour Technol 102:3155-

639         3162. doi:10.1016/j.biortech.2010.10.104

640    61. Wick RR, Judd LM, Gorrie CL, Holt KE (2017) Unicycler: Resolving bacterial genome

641         assemblies from short and long sequencing reads. PLoS Comput Biol 13:e1005595.

642         doi:10.1371/journal.pcbi.1005595

643    62. Yang SJ, Kataeva I, Hamilton-Brehm SD, Engle NL, Tschaplinski TJ, Doeppke C, Davis M,

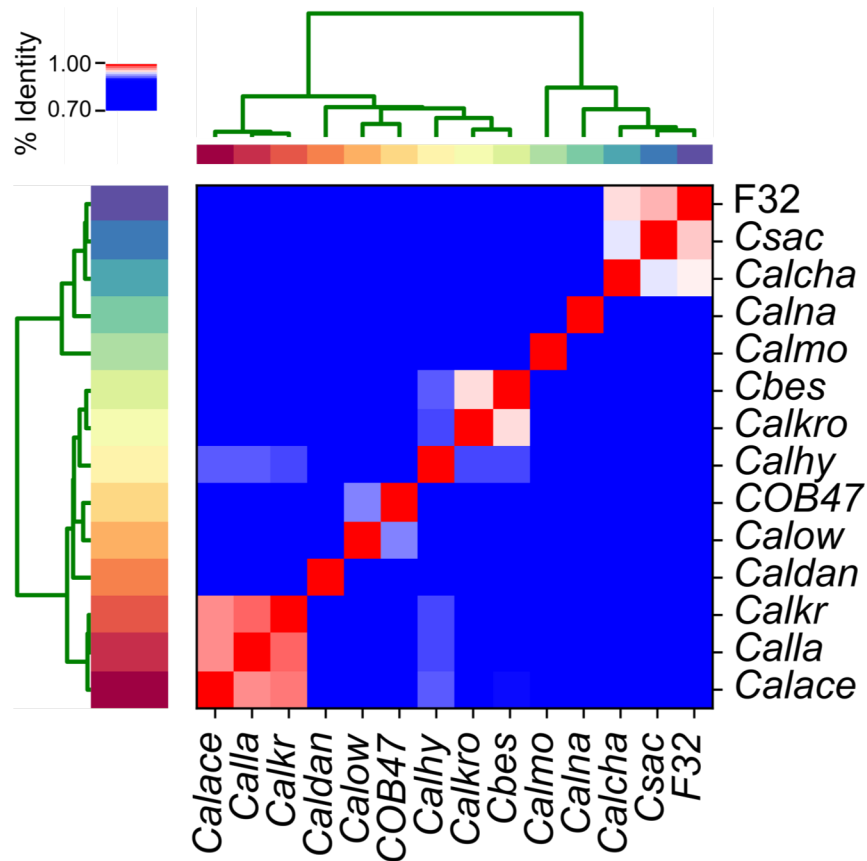644         Westpheling J, Adams MWW (2009) Efficient degradation of lignocellulosic plant

645     biomass, without pretreatment, by the thermophilic anaerobe "*Anaerocellum*

646     *thermophilum*" DSM 6725. Applied and Environmental Microbiology 75:4762-4769.

647     doi:10.1128/Aem.00236-09

648     63. Ying Y, Meng D, Chen X, Li F (2013) An extremely thermophilic anaerobic bacterium

649     *Caldicellulosiruptor* sp. F32 exhibits distinctive properties in growth and xylanases during

650     xylan hydrolysis. Enzyme Microb Technol 53:194-199.

651     doi:10.1016/j.enzmictec.2013.04.004

652     64. Yokoyama H, Yamashita T, Morioka R, Ohmori H (2014) Extracellular secretion of

653     noncatalytic plant cell wall-binding proteins by the cellulolytic thermophile

654     Caldicellulosiruptor bescii. J Bacteriol 196:3784-3792. doi:10.1128/JB.01897-14

655     65. Young J, Chung D, Bomble YJ, Himmel ME, Westpheling J (2014) Deletion of

656     *Caldicellulosiruptor bescii* CelA reveals its crucial role in the deconstruction of

657     lignocellulosic biomass. Biotechnol Biofuel 7:142. doi:10.1186/s13068-014-0142-6

658     66. Zurawski JV, Conway JM, Lee LL, Simpson HJ, Izquierdo JA, Blumer-Schuette S, Nookaew

659     I, Adams MW, Kelly RM (2015) Comparative analysis of extremely thermophilic

660     *Caldicellulosiruptor* species reveals common and unique cellular strategies for plant

661     biomass utilization. Appl Environ Microbiol 81:7159-7170. doi:10.1128/AEM.01622-15

662     67. Zverlov V, Mahr S, Riedel K, Bronnenmeier K (1998) Properties and gene structure of a

663     bifunctional cellulolytic enzyme (CelA) from the extreme thermophile '*Anaerocellum*

664     *thermophilum*' with separate glycosyl hydrolase family 9 and 48 catalytic domains.

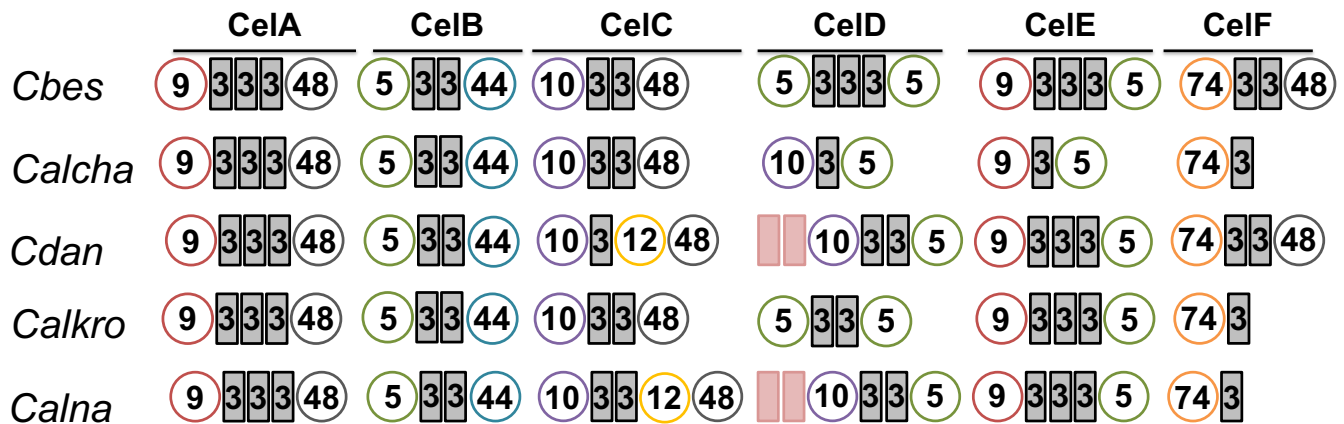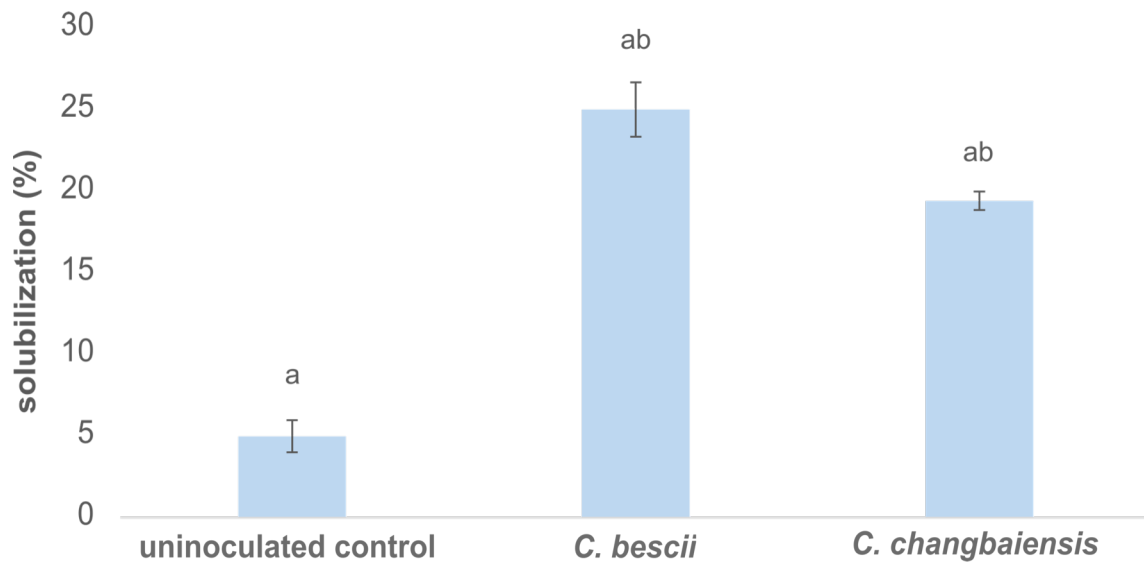665     Microbiology 144:457-465. doi:10.1099/00221287-144-2-457

666

**Figure 1. Core- and pangenome size estimates calculated from random sampling of 14 *Caldicellulosiruptor* genomes. (A)** Fitted curve of the estimated *Caldicellulosiruptor* core genome from 10 random samples of genomes up to n=14. The current size of the core genome is 1367 orthologous clusters. **(B)** Fitted curve of the estimated *Caldicellulosiruptor* pangenome from 10 random samples of genomes up to n=14. The *Caldicellulosiruptor* pangenome remains open and has increased to 3791 genes. The rate of growth for the pangenome is 63.2 new genes per genome sequenced. Core- and pangenome estimates were calculated from the equations reported by Tettelin *et al.*, [58] using GET_HOMOLOGUES software [19].
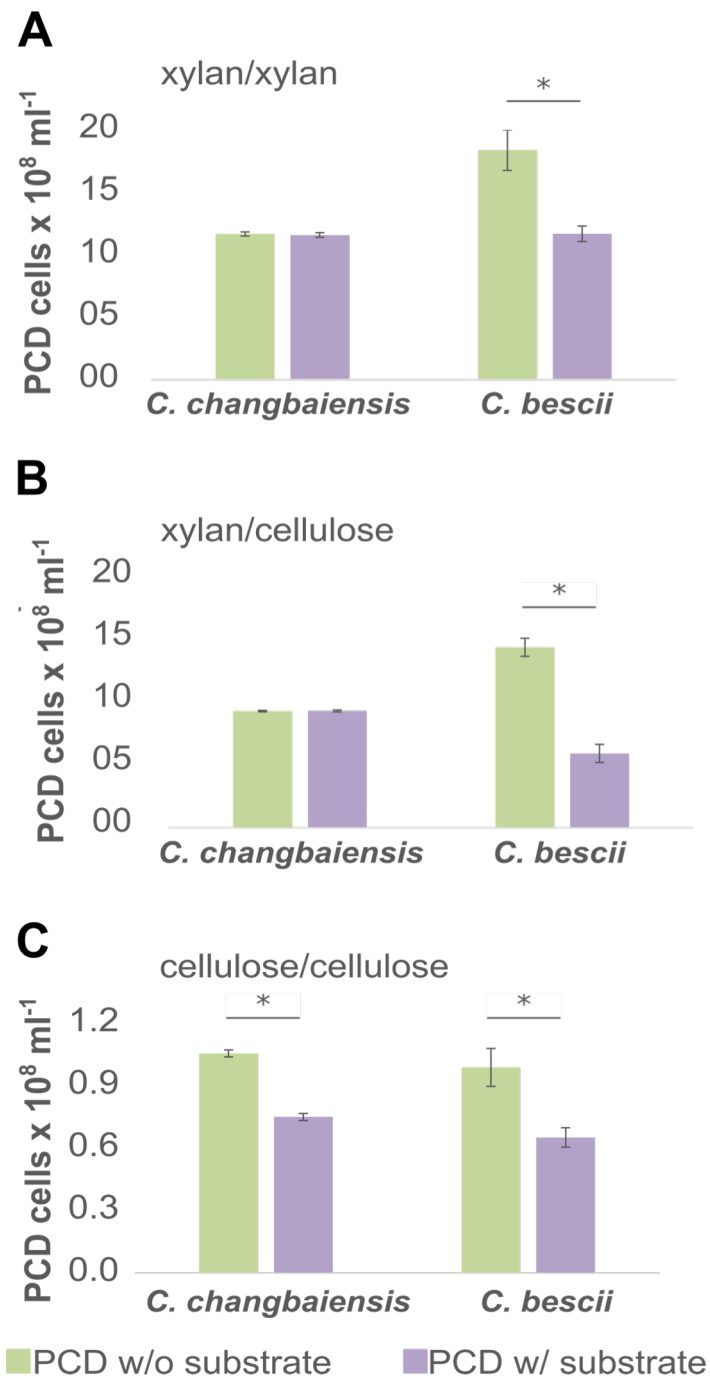
**Figure 2. Heatmap representation of the average nucleotide identity for 14 genome sequenced species from the genus *Caldicellulosiruptor*.** Average nucleotide identity (ANIb) was calculated on the basis of legacy BLASTn sequence identity over 1020nt sequence fragments. ANIb values of all 14 genomes are represented by a heat plot ranging from blue (75%< ANIb <90%), white (90%< ANIb <95%) to red (ANIb >95%). Pyani (https://github.com/widdowquinn/pyani) was used to calculate ANIb values and generate the clustered heatmap. Hierarchal cluster dendrograms were generated on the basis of similar ANIb values across each species. ANIb values are reported in **Table S1**. Calace, *C. acetigenus*; Cbes, *C. bescii*; Calcha, *C. changbaiensis*; Caldan, *C. danielii*; Calhy, *C. hydrothermalis*; Calkr, *C. kristjanssonii*; Calkro, *C. kronotskyensis*; Calla, *C. lactoaceticus*; Calmo, *C. morganii*; Calna, *C. naganoensis*; COB47, *C. obsidiansis*; Calow, *C. owensensis*; Csac, *C. saccharolyticus*; F32, *C.* sp. F32.
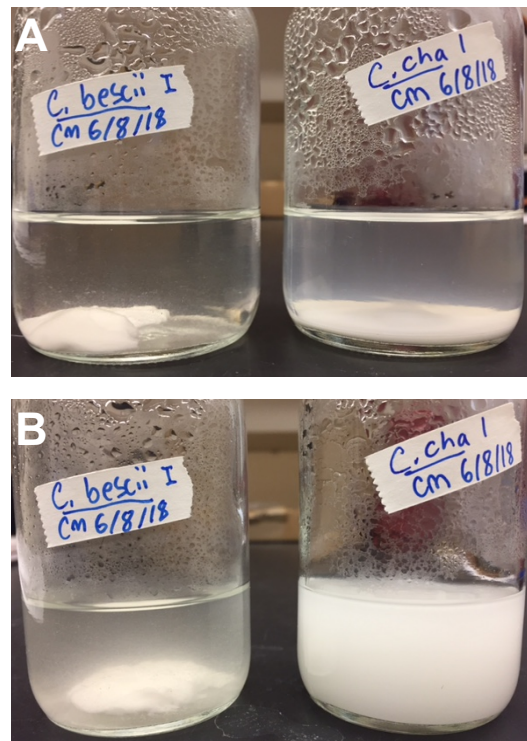
**Figure 3. Modular multifunctional enzymes encoded for by the glucan degradation locus.** Glucan degradation loci were selected on the basis of the presence of "ACE" cellulases. ACE cellulases: CelA, CelC and CelE. Circles represent the glycoside hydrolase (GH) domains, rectangles represent the carbohydrate binding module (CBM) domains. GH5, green circles; GH9, red circles; GH10, violet circles; GH 44, blue circles; GH48, grey circles; GH74, orange circles. CBM3, grey rectangles; CBM22, pink rectangles.
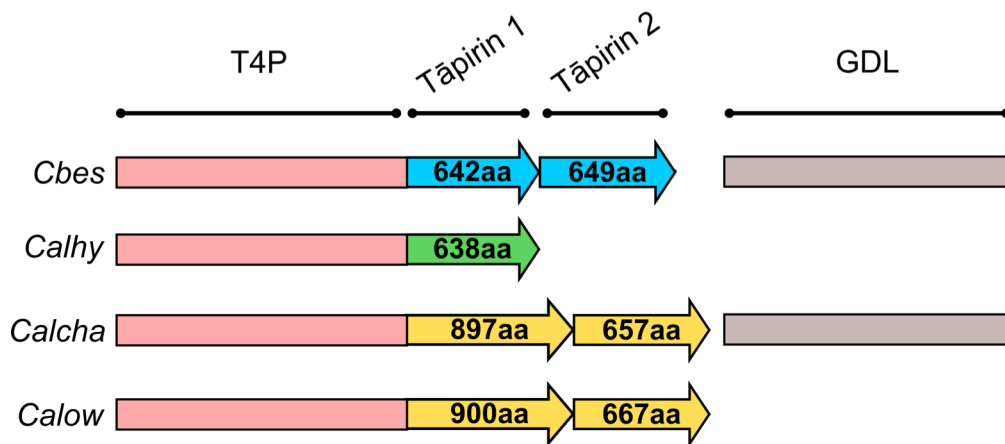
**Figure 4. Solubilization of microcrystalline cellulose by *C. bescii* and *C. changbaiensis*.** Uninoculated control, indicates abiotic cellulose solubilization in LOD medium. Error bars represent standard error (n=3). Similar letters over columns denote $p < 0.05$ as determined by a t-test.

**Figure 5. Comparison of the ability of *C. bescii* or *C. changbaiensis* planktonic cells to attach to polysaccharides.** Titles above bar charts indicate the carbon source for growth/ binding substrate. **(A, B)** When cells are grown on xylan, only planktonic *C. bescii* cells were able to attach to xylan or cellulose. **(C)** Cells grown on cellulose as the carbon source and exposed to cellulose as the binding substrate. Planktonic cell densities (PCD), enumerated by epifluorescence microscopy are plotted on the y-axis. Green columns indicate PCD without binding substrate and purple columns indicate PCD with the binding substrate. * indicates p < 0.01. All assays had n=6 biological replicates.

**Figure 6. Flocculation of *C. bescii* cells cultured on chemically defined medium and microcrystalline cellulose. (A)** Formation of a floc of *C. bescii* cells around microcrystalline cellulose (diameter, 20μm) while planktonic *C. changbaiensis* cells (cloudiness) are visible. **(B)** Same serum bottles as in "A", however the bottles were vigorously mixed. The *C. bescii* floc remains fairly stable, while both microcrystalline cellulose and cells are mixed in the *C. changbaiensis* culture.

**Figure 7. Genomic context for the location of the tāpirins from strongly to weakly cellulolytic** *Caldicellulosiruptor* **species.** Different colors represent the classical versus atypical tāpirins. Blue arrows: Cbes tāpirin 1 (Gen bank accession: YP_002573732) and Cbes tāpirin 2 (Gen bank accession: YP_002573731). Green arrow: Calhy tāpirin 1 (Gen bank accession number: YP_003992006). Yellow arrows: Calcha tāpirin 1 (Gen bank accession: WP_127352232.1) and 2 (Gen bank accession: WP_127352233.1), and Calow tāpirin 1 (Gen bank accession: YP_004002936) and 2 (Gen bank accession r: YP_004002935). Grey rectangles indicate the presence of the GDL downstream of the tāpirins. Atypical tāpirin 1 is annotated as a hypothetical protein and atypical tāpirin 2 is annotated as a von Willebrand factor A protein. Cbes, *C. bescii*; Calhy, *C. hydrothermalis*; Calcha, *C. changbaiensis* and Calow, *C. owensensis*. Peach rectangles represent the type IV pilus locus directly upstream of the tāpirins. Arrows indicate tāpirin 1 and 2. Numbers in the tāpirin arrows indicate the amino acid length.