

1 **Genome-scale resources for *Thermoanaerobacterium saccharolyticum***

2

3 Devin H. Currie¹, Babu Raman^{2,3}, Christopher M. Gowen^{4,5}, Timothy J. Tschaplinski²,
4 Miriam L. Land², Steven D. Brown², Sean F. Covalla¹, Dawn M. Klingeman², Zamin K.
5 Yang², Nancy L. Engle², Courtney M. Johnson², Miguel Rodriguez², A. Joe Shaw^{1,6},
6 William R. Kenealy¹, Stephen S. Fong⁴, Jonathan R. Mielenz², Brian H. Davison², David
7 A. Hogsett¹, Christopher D. Herring^{1,7,8}

8

9 1. Mascoma Corporation, 67 Etna Rd, Lebanon, NH 03766

10 2. BioEnergy Science Center, Oak Ridge National Laboratory, P.O. Box 2008, Oak
11 Ridge, TN 37831

12 3. Current affiliation: Dow AgroSciences, 9330 Zionsville Road, Indianapolis, IN 46268

13 4. Virginia Commonwealth University, P.O. Box 843028, Chemical and Life Science
14 Engineering, Richmond, Virginia 23284

15 5. Current affiliation: Centre for Applied Bioscience and Bioengineering, Department of
16 Chemical Engineering and Applied Chemistry, University of Toronto

17 6. Current affiliation: Novogy Inc., Cambridge, MA 02138

18 7. Current affiliation: Thayer School of Engineering, Dartmouth College, 14 Engineering
19 Drive, Hanover, NH 03755

20 8. To whom correspondence should be addressed

21

22 **Email Addresses**

23 Devin H. Currie - Devin.H.Currie.GR@dartmouth.edu

- 24 Babu Raman – raman.babu@gmail.com
- 25 Christopher M. Gowen - chris.gowen@utoronto.ca
- 26 Timothy J. Tschaplinski - tschaplinstj@ornl.gov
- 27 Miriam L. Land - landml@ornl.gov
- 28 Steven D. Brown - brownsd@ornl.gov
- 29 Sean F. Covalla - scovalla@mascoma.com
- 30 Dawn M. Klingeman - klingemandm@ornl.gov
- 31 Zamin K. Yang - yangz@ornl.gov
- 32 Nancy L. Engle - englenl@ornl.gov
- 33 Courtney M. Johnson - Courtney.Johnson10@fsis.usda.gov
- 34 Miguel Rodriguez - rodriguezmrj@ornl.gov
- 35 A. Joe Shaw - joeshawiv@gmail.com
- 36 William R. Kenealy - wkenealy@mascoma.com
- 37 Stephen S. Fong - ssfong@vcu.edu
- 38 Jonathan R. Mielenz - mielenzjr@ornl.gov
- 39 Brian H. Davison - davisonbh@ornl.gov
- 40 David A. Hogsett - dhogsett@opxbio.com

42 **Abstract**

43 **Background:** *Thermoanaerobacterium saccharolyticum* is a hemicellulose-degrading
44 thermophilic anaerobe that was previously engineered to produce ethanol at high yield.
45 A major project was undertaken to develop this organism into an industrial biocatalyst,
46 but the lack of genome information and resources were recognized early on as a key
47 limitation.

48
49 **Results:** Here we present a set of genome-scale resources to enable the systems level
50 investigation and development of this potentially important industrial organism.
51 Resources include a complete genome sequence for strain JW/SL-YS485, a genome-
52 scale reconstruction of metabolism, tiled microarray data showing transcription units,
53 mRNA expression data from 71 different growth conditions or timepoints and GC/MS-
54 based metabolite analysis data from 42 different conditions or timepoints. Growth
55 conditions include hemicellulose hydrolysate, the inhibitors HMF, furfural, diamide, and
56 ethanol, as well as high levels of cellulose, xylose, cellobiose or maltodextrin. The
57 genome consists of a 2.7 Mbp chromosome and a 110 Kbp megaplasmid. An active
58 prophage was also detected, and the expression levels of CRISPR genes were
59 observed to increase in association with those of the phage. Hemicellulose hydrolysate
60 elicited a response of carbohydrate transport and catabolism genes, as well as poorly
61 characterized genes suggesting a redox challenge. In some conditions, a time series of
62 combined transcription and metabolite measurements were made to allow careful study
63 of microbial physiology under process conditions. As a demonstration of the potential
64 utility of the metabolic reconstruction, the OptKnock algorithm was used to predict a set

65 of gene knockouts that maximize growth-coupled ethanol production. The predictions
66 validated intuitive strain designs and matched previous experimental results.

67

68 **Conclusion:** These data will be a useful asset for efforts to develop *T. saccharolyticum*
69 for efficient industrial production of biofuels. The resources presented herein may also
70 be useful on a comparative basis for development of other lignocellulose degrading
71 microbes, such as *Clostridium thermocellum*.

72

73 **Background**

74 Whether biomass-derived fuels play a major role in the world's energy future depends
75 on the development of technology to produce them at a cost that is competitive with
76 petroleum and other alternatives [1]. Fermentation of lignocellulose is an attractive
77 approach to this task [2, 3] and development of better fermenting organisms could
78 achieve much of the necessary cost reductions [4-6]. This represents an opportunity to
79 apply recent advances in metabolic engineering and systems biology to a problem of
80 major importance: the need for carbon-neutral fuels [7].

81

82 The thermophilic anaerobic bacteria include species with natural abilities to digest and
83 ferment the polysaccharides that make up the bulk of lignocellulosic biomass [8, 9].

84 Unfortunately, the lack of information and resources for these organisms has hindered
85 their development. *Thermoanaerobacterium saccharolyticum* is a Gram positive, low
86 G+C bacterium in the phylogenetic class "*Clostridia*" [10]. Members of the genus
87 *Thermoanaerobacterium* are thermophilic, rod shaped, chemoorganotrophic and able to

88 reduce thiosulfate to elemental sulfur. The species *T. saccharolyticum* can ferment a
89 wide array of carbohydrates, such as starch, xylan, glucose, cellobiose, xylose,
90 arabinose, mannose, and galactose, but cannot degrade crystalline cellulose [10]. Most
91 sugars are fermented to ethanol, acetic acid, lactic acid, carbon dioxide and hydrogen
92 [4]. *T. saccharolyticum* has a temperature range of 45-70°C, and pH range between 4.5-
93 7.0. The formation of endospores has not been observed in this species as they have in
94 the related genus *Clostridium*.

95

96 A variety of thermophilic enzymes of industrial utility have been isolated from *T.*
97 *saccharolyticum*, including endoxylanase, beta-xylosidase, amylopullanase and
98 glucuronidase [11-17]. A system for genetic manipulation of *T. saccharolyticum* was first
99 described by Mai et al. [18], which has been improved by the discovery of natural
100 competence [19], and the development of methods for making unmarked mutations with
101 negatively selectable markers [20]. The genes for lactate dehydrogenase, phosphate
102 transacetylase and acetate kinase were knocked out using these methods [4, 20]. The
103 result was a strain that produces ethanol at greater than 90% of theoretical yield,
104 comparable to other ethanologens such as yeast, *E. coli* or *Z. mobilis* [21, 22]. The
105 advantages that *T. saccharolyticum* has over these other biocatalysts are its elevated
106 growth temperature (matching the temperature optimum of many cellulases [21, 23]),
107 and its ability to hydrolyze hemicellulose and co-ferment the major sugars present in
108 lignocellulose [24].

109

110 Biomass is prepared for hydrolysis and fermentation by various forms of pretreatment in
111 order to expose the cellulose fibers and reduce particle size, though inhibitory
112 compounds, such as furfural and hydroxymethyl furfural (HMF) are generated in the
113 process [25]. Cost effective ethanol production requires ethanol concentrations > 40 g/L,
114 which necessitates that substrates, and by the same token their inhibitors, be present at
115 fairly high concentrations. The ability to reduce costs by increasing levels of pretreated
116 substrate is limited by the levels of inhibitors in the fermentation. While there is great
117 potential to reduce costs by developing organisms with greater tolerance to inhibitors,
118 little is known about the effects of these compounds on microbial physiology. One of the
119 goals of this project was to generate information about the effects of specific inhibitors
120 and complex inhibitor extracts from pretreated material. The project was undertaken as
121 part of a larger project to develop *T. saccharolyticum* for fermentation of pretreated
122 hardwood [26].

123

124 Another goal was to compare the genome of *T. saccharolyticum* to the genomes of
125 other bacteria potentially important for biofuel production, including *Clostridium*
126 *thermocellum*, an organism highly specialized for the hydrolysis of cellulose and the
127 focus of other OMICs, and systems biology efforts. This work supplements the
128 knowledge about both these important organisms and presents a comprehensive
129 resource for further investigation.

130

131 **Results and Discussion**

132 *Genome Sequencing*

133 As the sequence was being generated, there were early indications that contig 2 was in
134 fact a megaplasmid. Furthermore, early draft sequences showed that the ends of contig
135 2 overlapped. When PCR primers were designed at the ends of contig 2 facing
136 outwards, they amplified a product consistent with a circular DNA molecule. The gene
137 *Tsac_2822* on the putative megaplasmid encodes a RepB DNA replication protein with
138 high similarity to replication proteins from bacterial megaplasmids: *C. botulinum* plasmid
139 pCLI, *B. methanolicus* plasmid pBM19, and *B. weihenstephanensis* plasmid pBWB402.
140 Contig 2 was poorly represented in the initial Sanger sequence data and was observed
141 to be completely absent in strain ALK2; its loss as a complete unit further supporting its
142 identification as an extra chromosomal unit [4].

143
144 The genome contains 39 ORFs predicted to have transposase function, with 12 of these
145 concentrated in a 50 kbp region. The tool Prophage Finder [27] was used to identify two
146 regions containing genes with similarity to known phage genes. These two regions are
147 36 kbp and 42 kbp, located between ORFs *Tsac_2404* – *Tsac_2458* in contig 1 and
148 between ORFs *Tsac_2829*-*Tsac_2885* in contig 3 (the later listed under a separate
149 accession number in GenBank, CP003186). Close examination of individual reads of
150 CP003186 showed that some proceeded from the phage into contig 1 near position
151 2,009,359, suggesting a phage integration site. Contig 1 reads showed that in some, but
152 not all of them the phage was absent. In those reads where the phage was missing, the
153 sequence at 2009359-2009371 was duplicated. Primers were designed to the
154 chromosome flanking this region and in contig 3 facing outwards. All combinations of

155 primers amplified, supporting the conclusion that contig 3 is a phage that exists in both
156 integrated and circular forms at this locus (Supplemental Figure S1).

157
158 The chromosome contains a region containing 39 CRISPR repeats along with 8
159 CRISPR-associated genes. The CRISPR spacers were aligned with BLAST against the
160 genome and two of them were found to match the two putative phage regions. This
161 suggests that this strain of *T. saccharolyticum* has a history of infection and defense
162 against these two phage [28]. Analysis of *C. thermocellum* also showed possible
163 prophages and much more numerous and extensive CRISPR repeats and CRISPR-
164 associated genes, possibly related to the low transformation efficiency of *C.*
165 *thermocellum* [29]. Additional analysis across other Clostridia show further CRISPR
166 features [30].

167
168 A high percentage of genes (11.2%) have predicted functions (i.e. COG category)
169 related to carbohydrate transport and metabolism. For comparison, only 6.5% of the
170 ORFs in *C. thermocellum* ATCC27405 are assigned to this functional group. Both ABC-
171 type and phosphotransferase transporters occur. The tool dbCAN [31] was used to
172 compare all *T. saccharolyticum* protein sequences to hidden Markov models (HMMs) of
173 all protein families in the CAZY database. The program identified 73 ORFs with
174 similarity to glycosyl hydrolase HMMs, including 3 in glycosyl hydrolase family 5 with a
175 predicted function of “cellulase.” It also identified 18 proteins with similarity to Cellulose
176 Binding Module HMMs. It should be noted though that *T. saccharolyticum* does not
177 grow on crystalline cellulose such as avicel [10].

178

179 Surprisingly, a total of 67 sporulation-associated genes were identified, including *spo0A*,
180 but the strain is sporulation deficient. As with *Thermoanaerobacterium*
181 *thermosaccharolyticum* [32], *T. saccharolyticum* contains the nitrogenase genes
182 required for nitrogen fixation. Sequenced members of the related genus
183 *Thermoanaerobacter* apparently do not.

184

185 The genome contains 5 ribosomal regions, all oriented in the same direction.
186 Remarkably, the ribosomal sequences are not uniform, but rather of two types showing
187 only 95% identity in the “universal” region of the 16s subunit (Figure 1). Similar, but less
188 pronounced heterogeneity of ribosomal sequences has been noted in other firmicutes
189 [33], but has yet to be explained.

190

191 It is possible that these additional sequences confer some advantage during growth at
192 elevated temperatures. Another possibility is that these modifications decrease
193 sensitivity to an environmentally prevalent antibiotic that targets the 16s rRNA. The 16s
194 rRNA is a common target for antibiotic compounds, for example aminoglycosides [34].
195 That said, resistance-conferring mutations are frequently single base pair changes
196 rather than large insertion events [34, 35]. In addition, at least for the aminoglycosides,
197 the reported site of action is the A site near the 3’ end of the 16s rRNA [35, 36],
198 whereas these insertions are very close to the 5’ end. However, the version that
199 contains the inserts causes the 5’ and 3’ ends to no longer be located near one another,

200 as can be seen in Figure 1 panels B and C, and thus may be playing a role in
201 resistance.

202

203 *Effects of hemicellulose extract*

204 Spotted microarrays were used to examine the effect that biomass-derived hydrolysate
205 and the associated inhibitors have on *T. saccharolyticum*. In an initial experiment, cells
206 were grown in fermenters containing rich medium and a mixture of xylose and glucose
207 to mid-logarithmic phase, whereupon the cells were “shocked” by the addition of 10%
208 volume of hemicellulose extract (“washate”). Control fermentations were conducted by
209 shocking the cells with a mix of xylose and acetate at the same concentration and pH.
210 The cells continued to grow, though growth was slightly slowed. Samples were
211 analyzed before and up to 1 hour after the shock using spotted microarrays. Each
212 mRNA sample was measured relative to a genomic DNA control, and all log₂ ratios
213 given below are relative to the gDNA control [37].

214

215 When comparing the results from control reactors to those treated with washate, an
216 increasing number of genes were upregulated over time in response to washate (spots
217 above the diagonal in Figure 2). An alternate way of analyzing the same data is by
218 comparing expression levels at a given time point to those previous to the shock
219 (Supplemental Figure S2). Such comparisons versus the pre-shock culture showed
220 more scatter, most likely due to growth phase-related gene expression changes.

221

222 Most of the genes affected by the washate were upregulated, with 58 having \log_2 ratios
223 > 1 in at least one time point (Figure 3). At 5 min post-shock, a cluster of 17 genes
224 (Tsac_1270-1286) was upregulated. This cluster includes glycosyl hydrolases and
225 carbohydrate transport and catabolism genes, including three genes required for
226 arabinose utilization. At 15 min post-shock, additional genes were upregulated,
227 including those responsible for the formation of bacterial microcompartments and
228 rhamnose utilization.

229
230 cDNA samples from before and 1 hour post washate shock were also hybridized to tiled
231 Nimblegen microarrays. Compared to data from spotted arrays, the tiled array data
232 showed less noise in the lower dynamic range. Moreover, by examining the expression
233 levels visually, the boundaries of transcription units can be determined (Figure 4).

234

235 *Effects of HMF and furfural*

236 Two of the major inhibitors in washate are furfural and hydroxymethylfurfural (HMF). To
237 further investigate the specific effect these components have on *T. saccharolyticum*, we
238 performed additional “shock” experiments in which HMF and furfural were added during
239 logarithmic growth, while observing the cellular response by microarray and metabolite
240 analysis. The levels of HMF and furfural in pretreated hardwood hydrolysates is
241 approximately 0.1 g/L. We tested additions of HMF and furfural from 0.1 to 1.0 g/L and
242 found that 0.5 g/L of each showed a clear effect on growth (data not shown). Notably,
243 the effect was greatly diminished in medium containing yeast extract, so a defined
244 medium was used in this experiment. Sample processing methods for metabolite

245 analysis were validated as described in Supplemental Figure S3 and Supplemental
246 Table S1. Actively growing fermenters of *T. saccharolyticum* strain M700 at an O.D. of
247 0.6 were shocked with 0.5 g/L HMF and furfural. Samples were taken before the shock
248 and at 15 minutes, 1, 2, and 4 hours after shock. These samples were analyzed via
249 GC/MS and spotted microarrays.

250

251 A total of 40 different metabolites were tracked over the time course of the experiment
252 (Figure 5, Supplemental Figure S4, Supplemental Table S2). Almost all metabolites
253 decreased after 15 minutes of exposure to HMF and furfural, with the exception of
254 hydroxymethylfurfural and citric acid. Hydroxymethylfurfural, presumably resulting from
255 the reduction of HMF, increased steadily throughout the 4 hours that metabolites were
256 tracked. HMF and furfural were almost entirely metabolized after 16 hours. It is notable
257 that glucose-6-phosphate is among the many metabolites that decrease as the result of
258 HMF and furfural addition. This suggests that the inhibition occurs very early in the
259 glycolysis pathway, either at glucose transport or its phosphorylation.

260

261 Microarray analysis of the same fermentations showed large expression differences in
262 the phage loci between replicates during growth in the presence of HMF and furfural
263 (Supplemental Figure S5). Other non-phage genes were also observed to change
264 sympathetically with the phage genes, including the aforementioned CRISPRs.

265

266 In order to determine if some of the same genes were affected by the addition of HMF /
267 furfural as by washate, a comparison of the two datasets was performed. The \log_2 ratio

268 difference was calculated and analyzed via *t*-tests using the control from the same time
269 point as reference for washate shock and using the pre-shock as reference for HMF /
270 furfural shock. The 15 minute and 60 minute time points were considered for each, and
271 the greater log₂ ratio or significance value was used. In the washate shock experiment
272 502 genes were significantly affected (P value < 0.05) in either the 15 or 60 minute time
273 points, and in the HMF / furfural shock experiment, 414 genes were affected in either
274 the 15 or 60 minute time points. Between the two sets of significant genes, 88 were in
275 common (Supplemental Table S3). Of these, 40 had a log₂ difference in either
276 experiment greater than 0.7, and 9 had a log₂ difference greater than 0.7 in both.
277 Among these notable genes upregulated after both types of shock are members of a
278 gene cluster related to sulfur assimilation (Tsac_1655-1665), alanine dehydrogenase
279 (Tsac_2175) and NADPH-dependent methylglyoxal reductase (Tsac_1406). It should
280 be noted, however, that this comparison is less than ideal in that different media and
281 strains were used and that phage activity was noted in half of the HMF+furfural shock
282 samples.

283

284 A wealth of other microarray and metabolite data were generated (Table 1). Note that at
285 each timepoint listed in Table 1, multiple biological replicates were usually generated.
286 Data are available as Supplemental Files 1-5.

287

288 **Genome Scale *in silico* metabolic model**

289 Genome-scale constraint-based metabolic models are useful tools for exploring the
290 metabolic capabilities of an organism and for integrating bioinformatics data sets with

291 the metabolic network. A genome-scale model for *Thermoanaerobacterium*
292 *saccharolyticum* was built for this study based on its genomic content, current literature,
293 and experimental data (Supplemental File 6). An initial reaction list was built by
294 comparing its genetic content with that of the related bacterium *Clostridium*
295 *thermocellum*, for which a curated metabolic model already exists [38]. To do this, a
296 BLAST search was performed using the genes included in the *C. thermocellum* model
297 versus the *T. saccharolyticum* ORF predictions. This resulted in an initial set of 425
298 reactions with gene-reaction mappings to serve as the foundation of the reconstruction.
299 Additionally, metabolic pathways for xylose and sorbitol metabolism were added, and
300 cellulose breakdown reactions were removed. The bifurcating ferredoxin:NAD
301 oxidoreductase described by Wang *et al.* was added as well [39]. A number of other
302 changes were made based on biochemical evidence, and additional gap filling was
303 performed as described in the methods section to generate a working model. These
304 changes are detailed in Supplementary File 7. The final model contains 528 metabolites
305 and 516 reactions associated with 315 genes. A comparison of these statistics to the *C.*
306 *thermocellum* model is shown in Table 2.

307

308 **Model Validation**

309 Although the metabolic network composition at this stage was consistent with available
310 information based on genome annotation and experimental observations, the resulting
311 flux space remained highly underdetermined. This is a consistent challenge facing all
312 constraint-based models, because many thermodynamic and regulatory effects cannot
313 be captured in the stoichiometric network. In particular for the *T.saccharolyticum* model,

314 the diversity of hydrogenase systems hosted by this organism, left unconstrained,
315 provide the network with many ways to efficiently regenerate cofactors, allowing
316 unrealistic levels of flux towards acetate and hydrogen. From a thermodynamic
317 standpoint, actual allowable fluxes through these reactions are limited by many factors,
318 including the intracellular concentrations of the cofactors, the concentration of hydrogen,
319 intra- and extracellular pH, and the reduction potential of ferredoxin. This problem is
320 complicated further by the kinetics and expression levels of the responsible enzymes. In
321 the absence of the necessary parameters to formulate these constraints, we decided on
322 a top-down approach to replicate experimental observations by making some of the
323 hydrogenase reactions irreversible and by limiting the overall hydrogen production to
324 observed yields. In a previous study [40], the four gene operon *hfs* coding for the
325 reaction ferredoxin hydrogenase was found to be the primary hydrogen producer *in*
326 *vivo*, whereas the other hydrogenase genes tested were found to contribute only slightly
327 or not at all to hydrogen production. Reflecting this in the model, the energy-conserving
328 hydrogenase (ECH) was blocked, and the bifurcating hydrogenase (BIFH2) and the
329 NADH hydrogenase (NADH2) were forced to be irreversible in the direction of hydrogen
330 uptake. Additionally, total hydrogen export was limited to a yield of 0.9 M H₂:M glucose
331 to reflect the *in vivo* measurements [40]. These modifications had a dramatic impact on
332 the predicted performance of the model by limiting the amount of reducing equivalents
333 that could be sent to hydrogen production, thereby shifting some carbon flux from
334 acetate to ethanol and other organic acids. Further experimentation with hydrogenase
335 constraints may prove useful to help understand how electron and carbon flow are
336 related in this and other mixed-acid fermentors.

337

338 Previous metabolic engineering efforts on *T. saccharolyticum* [4, 40] have explored two
339 distinct strategies for improving ethanol yield: a carbon-centric approach that focuses on
340 eliminating competing carbon fluxes at the pyruvate branch point, and an electron-
341 centric approach that disrupts the cell's ability to produce hydrogen as a highly-reduced
342 electron acceptor. Each of these strategies was shown to improve ethanol production
343 to varying degrees. A phenotypic phase plane analysis was performed to explore the
344 metabolic implications of these knockout strategies. Figure 6 shows the optimal growth
345 surfaces for these knockouts over the complete ranges of possible carbon uptake and
346 ethanol production rates. In the wild-type seen in Figure 6-A, optimal growth can occur
347 across a wide range of ethanol flux values, limited by the maximum glucose uptake rate.
348 Knocking out the lactate dehydrogenase (LDH) and phosphotransacetylase (PTA)
349 reactions eliminates stoichiometrically equivalent solutions, leading to a maximum
350 growth rate that is coupled to high ethanol production (Figure 6-B). The coupling of
351 ethanol flux to growth rate was found to be much stronger, however, in the electron-
352 centered strategy (Figure 6-C), which removed the reactions for LDH and ferredoxin
353 hydrogenase (HFS). This knockout strategy greatly limits the available solution space
354 and tightly dictates the ethanol yield at some penalty to the maximum growth rate. This
355 finding is consistent with experimental results, which found a lower overall growth rate
356 and cellobiose uptake rate in the *ldh-hfs* deletion strain when compared to the wild-type
357 or *ldh-pta-ack* deletion strain [40]. However, the strong coupling of ethanol production to
358 growth rate in the *ldh-hfs* knockout strain implies that it may be a good candidate for
359 adaptive evolution to improve ethanol productivity.

360

361 We attempted to determine if any other knockout strain designs would maximize ethanol
362 production at an optimal growth rate. The bilevel optimization algorithm OptKnock [41]
363 was used to search for knockout strain designs that would improve production of
364 ethanol by coupling it to improved growth rate. When OptKnock searches a maximum of
365 2 reaction knockouts, optimal ethanol production is predicted when knocking out LDH
366 and HFS. When allowing three reactions knockouts, OptKnock finds a marginal
367 improvement by deleting LDH, HFS, and glutamate dehydrogenase (GLUD). Removal
368 of GLUD forces the cell to use the reactions glutamate synthase (GLUS) and glutamine
369 synthetase (GLNS) in order to incorporate ammonium, spending an additional mole of
370 ATP per mole of ammonium (Table 3). This inefficiency predicts only a marginal
371 improvement in ethanol production of 0.3% over the Δ LDH- Δ HFS strain (Figure 7).

372

373 *Conclusions*

374 Here we report the first genome scale study of the industrially important bacterium *T.*
375 *saccharolyticum*. This work informs and supports not only the study of fundamental
376 microbial physiology, but also its potential applications in this organism. The resources
377 presented herein will facilitate further efforts to engineer *T. saccharolyticum* for the
378 production of biofuels. In addition, ongoing engineering efforts in other organisms to
379 increase inhibitor tolerance and ethanol yield and titer may benefit from these data.

380

381 **Materials and Methods**

382 *Strains*

383 The tiled microarrays were designed based on wild-type *T. saccharolyticum* YS485
384 DSM 8691 (Deutsche Sammlung von Mikroorganismen und Zellkulturen GmbH,
385 Germany). The microarrays and metabolite profiling were performed using engineered
386 and evolved ethanologenic strains of *T. saccharolyticum*, described previously [26].

387

388 *Growth Media*

389 MTC media [42] was supplemented with various concentrations of glucose, xylose, and
390 mixtures of hemicellulose extract or acetic acid mixed with xylose. These concentrations
391 are noted for each experiment. TS5 media was developed specifically for *T.*
392 *saccharolyticum*. It is similar to the previously published media TSC1 [20] but with only
393 0.5 g/l KH_2PO_4 and with 0.5 g/l tryptone. The full media formulation per liter is: Solution I
394 (yeast extract 8.5 g, 1.85 $(\text{NH}_4)_2\text{SO}_4$, 0.05 g FeSO_4 , 0.5 g KH_2PO_4 , 1 g $\text{MgCl}_2 \cdot 6 \text{H}_2\text{O}$,
395 0.05 g CaCl_2 , 0.5 g Tryptone, 2 g Trisodium citrate $\cdot 2 \text{H}_2\text{O}$, 800 ml H_2O) and Solution II
396 (10 g Xylose, 200 ml H_2O). These are autoclaved separately to avoid burning the
397 xylose, then mixed.

398

399 *Hemicellulose extract*

400 Hemicellulose extract, or 'washate', for the microarray and metabolite profiling
401 experiments was prepared by suspending mixed hardwood pretreated with steam in an
402 Andritz horizontal plug-flow reactor to severity 3.8 in water at 30% solids. It was then
403 autoclaved for one hour, and the liquid was removed from the solids by vacuum filtration
404 using Whatman Grade No. 1 Filter Paper (Whatman Ltd, Kent, UK). It was then brought
405 to pH 6.0 using NH_4OH . The extract contained 11.52 g/L monomeric xylose, 0.89 g/L

406 glucose, 0.84 g/L lactate, 3.54 g/L acetate, 0.56 g/L HMF, 0.26 g/L furfural, and various
407 other inhibitors. Additional carbohydrate was present but not analyzed due to its
408 oligomeric state or lack of standards for analysis.

409

410 *Genome Sequencing of T. saccharolyticum YS485*

411 The genome of *T. saccharolyticum* JW/SL-YS485 was generated over a 4 year span by
412 a variety of techniques. Initially, a clone library was constructed and Sanger sequenced
413 to 8x coverage. Clones were selected for additional sequencing to close gaps, and
414 additional sequence data was generated with the 454 platform. The assembled draft
415 was then aligned to the complete genomes of *T. tengcongensis* and *T.*
416 *pseudoethanolicus*, allowing the contigs to be ordered and oriented to each other. PCR
417 primers were designed at the ends of contigs and used to amplify across gaps,
418 consisting mostly of ribosomal regions. These PCR products were Sanger sequenced
419 and used to manually close the genome. Finally, differences between the Sanger and
420 454 data were resolved by examining sequence data from various strains sequenced
421 with Illumina technology. Open reading frames were predicted and functions annotated
422 using an annotation pipeline at Oak Ridge National Laboratory. The genome sequence
423 has been assigned GenBank accession numbers CP003184.1, CP003185.1 and
424 CP003186.1 (the genome, the mega-pasmid, and the phage, respectively).

425

426 *Phage Confirmation*

427 Primers were designed to confirm the presence of a phage in contig 3 which is present
428 in both integrated and circular forms (C17_near_endF: CTGCCCGTGGAAACATCTAAT,

429 C17_near_startR: GTTGGTTCTGCCCTGTTTGT, C15_int_siteF:
430 TTTGCACCGCCATTTAAGAG, C15_int_siteR: ACGGTGATGAAGAAGCGAAA,
431 C18_near_startR: AATTCGGCATGTGTTGGAT). PCR was performed on genomic DNA
432 from *T. saccharolyticum* strain YS485 and *T. saccharolyticum* M700 and products were
433 separated on a 1% agarose gel.

434

435 *Spotted Microarray Construction, Hybridization and Analysis*

436 Spotted oligonucleotide microarrays were essentially constructed and hybridized as
437 described previously [43, 44]. Briefly, DNA sequences that represented predicted-
438 protein encoding sequences were obtained for the *T. saccharolyticum* YS485 genome
439 (NCBI GenBank accession numbers CP003184.1, CP003185.1 and CP003186.1) and
440 70-mer oligonucleotide probes were designed using the CommOligo software [45]. The
441 original genome sequence was in draft format and 2,627 oligonucleotide probes were
442 designed for 2,667 putative CDS, representing 98.5% of the predicted-protein encoding
443 sequences for the draft genome. Subsequently, refinements were made as the genome
444 sequence was closed. Oligonucleotides commercially synthesized without modification
445 (Integrated DNA Technologies, Coralville, Iowa) in 96-well stock plates and transferred
446 to 384-well printing plates in a final concentration of 50% DMSO using a BioMek FX
447 liquid handling robot (Beckman-Coulter, Fullerton, CA). Probes were then spotted onto
448 UltraGAPS glass slides (Corning Life Sciences, Corning, NY) using a BioRobotics
449 Microgrid II microarrayer (Genomic Solutions, Ann Arbor, MI) in a dust-free clean room
450 maintained at 21°C and 50% relative humidity. Spotted DNA was stabilized on slides by

451 ultraviolet cross-linking using a UV 1800 Stratlinker (Stratagene, La Jolla, CA)
452 according to slide manufacturer's instructions (Corning Life Sciences).

453
454 Total RNA was purified using an RNeasy Plus Mini Kit (Qiagen), which was used as
455 template to generate cDNA copies labeled with Cy5-dUTP (Amersham Biosciences,
456 Piscataway, NJ). Labeled genomic DNA (Cy3) was used as a control and as the
457 common reference to co-hybridize with labeled RNA (Cy5) samples for each slide and
458 microarray hybridization and washing conditions have been described elsewhere [43,
459 44, 46]. Microarray images were scanned using a ScanArray Express (PerkinElmer)
460 scanner, and spot signal, quality, and background fluorescent intensities were quantified
461 using ImaGene version 6.0 (Biodiscovery, Marina Del Rey, CA). Outlier detection,
462 background correction, normalizations and log ratios were generated as described
463 previously [46], except that the workflow was conducted using JMP Genomics (SAS
464 Institute Inc.) with custom scripts.

465
466 *Washate Shock Microarrays*
467 *T. saccharolyticum* ALK2 was grown overnight shaking in bottles with 50 ml MTC + 5
468 g/L glucose and xylose to an optical density of 0.5 at 600 nm. 25 ml from these bottles
469 were inoculated into 1 L of MTC media + 2.7 g/L xylose + 4.6 g/L glucose.
470 Fermentations were performed in duplicate at 1 L volume in Sartorius BiostatA+
471 reactors maintained at pH 5.8, 55°C, stirring 150 rpm, and purged with N₂/CO₂ prior to
472 inoculation. Upon reaching an OD of 0.6, 100 ml of either a control solution (11.5 g/L
473 xylose, 3.5 g/L acetic acid, with pH adjusted to 6.0 with NH₄OH) or washate was added.

474 Samples were taken at time 0, 5, 15, and 60 minutes after shock. The samples were
475 mixed with 30 ml RNAprotect bacteria reagent (QIAGEN Corp, Valencia CA) and left at
476 room temperature for 5 minutes. The samples were then centrifuged at 4000 rpm for 10
477 min at 4°C. The pellets were then resuspended in 1 ml SET buffer and stored at -80°C.

478

479 *HMF and Furfural Shock Microarrays*

480 *T. saccharolyticum* M700 was grown overnight in bottles with 50 ml Defined TS5 media
481 shaking at 55°C. 25 ml from these bottles were inoculated into 4 reactors containing 1 L
482 of Defined TS5 and maintained at pH 5.8, 55°C, stirring 150 rpm, and under constant
483 purging with N₂/CO₂. The reactors were grown to an O.D. of 0.06 at which point 0.5 g/l
484 each of HMF and furfural were added to two of the reactors, leaving the second two as
485 controls. Samples were taken at times 0, 15 minutes, 1, 2, and 4 hours after addition.
486 Two sets of samples were taken at each time point, one for microarrays and one for
487 proteomics. The samples for microarray analysis were mixed with 30 ml RNAprotect
488 bacteria reagent (QIAGEN Corp, Valencia CA) and left at room temperature for 5
489 minutes. The samples were then centrifuged at 4000 rpm for 10 min at 4°C. The pellets
490 were then resuspended in 1 ml SET buffer and stored at -80°C. The samples for the
491 metabolite profiling assays were centrifuged at 4°C at 4000 rpm for 10 minutes,
492 supernatants were poured off and the pellets were frozen at -80°C

493

494 *Tiled Microarrays*

495 Tiled microarrays were performed by Nimblegen Corporation (Madison, WI).

496

497 *Metabolite profiling*

498 Metabolites from *T. saccharolyticum* culture pellets and hydrolysates were analyzed as
499 trimethylsilyl (TMS) derivatives by gas chromatography-mass spectrometry (GC/MS)
500 using electron impact (EI) ionization, as described previously (Yang et al. 2009). Briefly,
501 aliquots of culture supernatants (50 μ L to 2 mL) and sorbitol (aqueous internal standard
502 added to yield 10 – 60 ng per μ L injected) were transferred by pipette to a vial and
503 stored at -20°C until analyzed. Microbial pellets were fast-frozen in liquid nitrogen and
504 stored at -80 °C until analyzed. Frozen pellets were weighed and added to 10 mL 80%
505 ethanol containing sorbitol as internal standard. Cell pellets were ruptured by sonication
506 with temperature maintained below 0 °C, and cell debris separated from the extract by
507 centrifugation at 4°C, and 2 mL were dried in a stream of N₂ prior to silylation. The
508 hydrolysate samples were thawed and also concentrated to dryness under a stream of
509 N₂. The internal standard was added to correct for subsequent differences in
510 derivatization efficiency and changes in sample volume during heating. Dried extracts
511 were dissolved in 500 μ L of silylation–grade acetonitrile followed by the addition of 500
512 μ L N-methyl-N-trimethylsilyltrifluoroacetamide (MSTFA) with 1% trimethylchlorosilane
513 (TMCS) (Thermo Scientific, Bellefonte, PA), and samples then heated for 1 h at 70 °C to
514 generate TMS derivatives. After 2-3 days, 1- μ L aliquots were injected into an Agilent
515 Technologies Inc. (Santa Clara, CA) 5975C inert XL gas chromatograph-mass
516 spectrometer, fitted with an Rtx-5MS with Integra-guard (5% diphenyl/95% dimethyl
517 polysiloxane) 30 m x 250 μ m x 0.25 μ m film thickness capillary column. The standard
518 quadrupole GC/MS was operated in the EI (70 eV) ionization mode, with 6 full-spectrum
519 (50-650 Da) scans per second. Gas (helium) flow was 1.3 mL per minute with the

520 injection port configured in the splitless mode. The injection port, MS Source, and MS
521 Quad temperatures were 250 °C, 230 °C, and 150 °C, respectively. The initial oven
522 temperature was held at 50 °C for 2 min and was programmed to increase at 20 °C per
523 min to 325 °C and held for another 11 min, before cycling back to the initial conditions. A
524 large user-created database (>1800 spectra) of mass spectral electron ionization (EI)
525 fragmentation patterns of TMS-derivatized compounds, as well as the Wiley Registry
526 8th Edition combined with NIST 05 mass spectral database, were used to identify the
527 metabolites of interest to be quantified. Peaks were reintegrated and reanalyzed using a
528 key selected ion, characteristic m/z fragment, rather than the total ion chromatogram, to
529 minimize integrating co-eluting metabolites. The extracted peaks of known metabolites
530 were scaled back up to the total ion current using predetermined scaling factors. The
531 scaling factor for the internal standard was used for unidentified metabolites. Peaks
532 were quantified by area integration and the concentrations were normalized to the
533 quantity of the internal standard recovered, volume of sample processed, derivatized,
534 and injected. Three to six replicate samples were analyzed per time point, and the
535 metabolite data were averaged at a given time point. Unidentified metabolites were
536 denoted by their retention time as well as key mass-to-charge (m/z) ratios.

537

538 *Constraint-based modeling of Thermoanaerobacterium saccharolyticum*

539 Initial construction of the *Thermoanaerobacterium saccharolyticum* reaction list was
540 based on the previously published model of the closely related species *Clostridium*
541 *thermocellum* [38, 47]. This was accomplished by using BLAST to search for genes in
542 *T. saccharolyticum* that were homologous to the genes represented in the *C.*

543 *thermocellum* model. Further refinement to the model was done by manual curation,
544 incorporating available biochemical and genetic information. The resulting reaction list
545 was not yet able to produce flux through the biomass reaction using appropriate
546 exchange boundary conditions, so additional gap filling was required. This was
547 accomplished through the use of a novel gap filling algorithm called FBA-gap [48] which
548 proposes a minimal set of reaction additions necessary to support biomass production.
549 These reactions are sourced from a reaction database populated using the reaction lists
550 of available stoichiometric models.

551
552 Flux balance analysis (FBA) [49], was used throughout the reconstruction and analysis
553 of the *T. saccharolyticum* model to simulate optimal growth. Modeling work was done
554 using the COBRA toolbox for Matlab [50, 51] along with custom methods and the
555 Gurobi Optimizer. OptKnock [41] was used to search for knockout strains that would
556 putatively couple ethanol production with an improved growth rate. An implementation
557 of OptKnock is available in the COBRA toolbox for MATLAB.

558
559 *Statistical analysis*

560 All statistical analyses were performed using R statistical software [52] and the package
561 gplots [53].

562

563 **Availability of supporting data**

564 The data sets supporting the results of this article are included within the article (and its
565 additional files).

566

567 **List of abbreviations**

568 HMF- hydroxymethyl furfural

569 ORF - Open Reading Frame

570 CRISPR - Clustered Regularly Interspaced Short Palindromic Repeats

571 HMM - Hidden Markov Model

572 ECH - Energy-Conserving Hydrogenase

573 BIFH2 - Bifurcating Hydrogenase

574 NADH2 - NADH Hydrogenase

575 LDH - Lactate Dehydrogenase

576 PTA - Phosphotransacetylase

577 HFS - Ferredoxin Hydrogenase

578 GLUD - Glutamate Dehydrogenase

579 GLUS - Reactions Glutamate Synthase

580 GC/MS - Gas Chromatography-Mass Spectrometry

581 MSTFA - N-methyl-N-trimethylsilyltrifluoroacetamide

582 BLAST - Basic Local Alignment Search Tool

583 FBA - Flux Balance Analysis

584

585 **Competing Interests**

586 DHC was partially supported by funding from Mascoma Corporation during his doctoral

587 work at Dartmouth College. SFC, AJS and CDH were salaried employees of Mascoma

588 Corporation and are listed as inventors on several patent filings related to

589 *Thermoanaerobacterium saccharolyticum* by Mascoma Corporation. WRK was a
590 salaried employee of Mascoma Corporation. DAH is a stockholder of Mascoma
591 Corporation and was a salaried employee of Mascoma Corp. and Dartmouth College.
592 He is listed as an inventor on several patent filings related to *T. saccharolyticum* by
593 Mascoma Corporation and Dartmouth College, each of which owns or has applied for
594 patents related to *T. sacchraolyticum*. JRM DAH is a stockholder of Mascoma
595 Corporation. All other authors have no competing interests.

596

597 **Author Contributions**

598 DH Currie analyzed data and drafted the manuscript, B Raman performed microarray
599 work and was the main point of contact at ONRL, CM Gowen performed the metabolic
600 reconstruction and drafted the manuscript, TJ Tschaplinski performed metabolite
601 profiling analyses, ML Land performed genome analysis and annotation, SD Brown
602 designed the microarrays and assisted with analysis, SF Covalla performed
603 fermentations and processed samples, DM Klingeman processed microarray samples,
604 ZK Yang processed microarray samples, NL Engle processed metabolite sample
605 extracts, CM Johnson processed microarray samples, M Rodriguez assisted in
606 performing fermentations, AJ Shaw assisted in metabolic analysis, WR Kenealy
607 supervised fermentation work, SS Fong supervised metabolic reconstruction and
608 analysis, JR Mielenz coordinated work at ORNL and helped design experiments, BH
609 Davison coordinated work at ORNL and helped design experiments, DA Hogsett
610 supervised the project, CD Herring directed the project and drafted the manuscript.

611

612 **Acknowledgements**

613 ORNL annotated the genome sequence with funding from Mascoma Corp. The U.S.
614 Department of Energy's Office of Energy Efficiency and Renewable Energy Office in the
615 BioEnergy Technologies Office provided support for the experimental work and analysis
616 under award # GO17057. Manuscript preparation was supported by the BioEnergy
617 Science Center. The BioEnergy Science Center is a U.S. Department of Energy
618 Bioenergy Research Center supported by the Office of Biological and Environmental
619 Research in the DOE Office of Science. Oak Ridge National Laboratory is managed by
620 UT-Battelle, LLC, for the DOE under Contract DE-AC05-00OR22725.

621

622

623 **References**

624

- 625 1. Sims REH, Mabee W, Saddler JN, Taylor M: **An overview of second**
626 **generation biofuel technologies.** *Bioresource Technol* 2010, **101**(6):1570-
627 1580.
- 628 2. Olson DG, McBride JE, Joe Shaw A, Lynd LR: **Recent progress in**
629 **consolidated bioprocessing.** *Curr Opin Biotechnol* 2012, **23**(3):396-405.
- 630 3. Lynd LR, van Zyl WH, McBride JE, Laser M: **Consolidated bioprocessing of**
631 **cellulosic biomass: an update.** *Curr Opin Biotechnol* 2005, **16**(5):577-583.
- 632 4. Shaw AJ, Podkaminer KK, Desai SG, Bardsley JS, Rogers SR, Thorne PG,
633 Hogsett DA, Lynd LR: **Metabolic engineering of a thermophilic bacterium to**

- 634 **produce ethanol at high yield.** *Proc Natl Acad Sci U S A* 2008, **105**(37):13769-
635 13774.
- 636 5. Shaw JA, Covalla SF, Miller BB, Firliet BT, Hogsett DA, Herring CD: **Urease**
637 **expression in a *Thermoanaerobacterium saccharolyticum* ethanologen**
638 **allows high titer ethanol production.** *Metab Eng* 2012, **14**:528–532.
- 639 6. Parawira W, Tekere M: **Biotechnological strategies to overcome inhibitors in**
640 **lignocellulose hydrolysates for ethanol production: review.** *Crit Rev*
641 *Biotechnol* 2011, **31**(1):20-31.
- 642 7. Ragauskas AJ, Williams CK, Davison BH, Britovsek G, Cairney J, Eckert CA,
643 Frederick WJ, Hallett JP, Leak DJ, Liotta CL *et al*: **The path forward for**
644 **biofuels and biomaterials.** *Science* 2006, **311**(5760):484-489.
- 645 8. Lee J: **Biological conversion of lignocellulosic biomass to ethanol.** *J*
646 *Biotechnol* 1997, **56**(1):1-24.
- 647 9. Chang T, Yao S: **Thermophilic, lignocellulolytic bacteria for ethanol**
648 **production: current state and perspectives.** *Appl Microbiol Biotechnol* 2011,
649 **92**(1):13-27.
- 650 10. Lee YE, Jain MK, Lee CY, Lowe SE, Zeikus JG: **Taxonomic Distinction of**
651 **Saccharolytic Thermophilic Anaerobes - Description of**
652 **Thermoanaerobacterium-Xylanolyticum Gen-Nov, Sp-Nov, and**
653 **Thermoanaerobacterium-Saccharolyticum Gen-Nov, Sp-Nov -**
654 **Reclassification of Thermoanaerobium-Brockii, Clostridium-**
655 **Thermosulfurogenes, and Clostridium-Thermohydrosulfuricum E100-69 as**
656 **Thermoanaerobacter-Brockii Comb-Nov, Thermoanaerobacterium-**

- 657 **Thermosulfurigenes Comb-Nov, and Thermoanaerobacter-**
658 **Thermohydrosulfuricus Comb-Nov, Respectively - and Transfer of**
659 **Clostridium-Thermohydrosulfuricum 39e to Thermoanaerobacter-**
660 **Ethanolicus. *Int J Syst Bacteriol* 1993, **43**(1):41-51.**
- 661 11. Podkaminer KK, Guss AM, Trajano HL, Hogsett DA, Lynd LR: **Characterization**
662 **of xylan utilization and discovery of a new endoxylanase in**
663 **Thermoanaerobacterium saccharolyticum through targeted gene deletions.**
664 *Appl Environ Microbiol* 2012.
- 665 12. Vocadlo DJ, Wicki J, Rupitz K, Withers SG: **Mechanism of**
666 **Thermoanaerobacterium saccharolyticum beta-xylosidase: kinetic studies.**
667 *Biochemistry* 2002, **41**(31):9727-9735.
- 668 13. Bronnenmeier K, Meissner H, Stocker S, Staudenbauer WL: **alpha-D-**
669 **glucuronidases from the xylanolytic thermophiles Clostridium stercorarium**
670 **and Thermoanaerobacterium saccharolyticum. *Microbiology* 1995, **141** (Pt**
671 **9):2033-2040.**
- 672 14. Ramesh MV, Podkovyrov SM, Lowe SE, Zeikus JG: **Cloning and sequencing**
673 **of the Thermoanaerobacterium saccharolyticum B6A-RI apu gene and**
674 **purification and characterization of the amylopullulanase from Escherichia**
675 **coli. *Appl Environ Microbiol* 1994, **60**(1):94-101.**
- 676 15. Lee YE, Lowe SE, Zeikus JG: **Gene cloning, sequencing, and biochemical**
677 **characterization of endoxylanase from Thermoanaerobacterium**
678 **saccharolyticum B6A-RI. *Appl Environ Microbiol* 1993, **59**(9):3134-3137.**

- 679 16. Lee YE, Zeikus JG: **Genetic organization, sequence and biochemical**
680 **characterization of recombinant beta-xylosidase from**
681 **Thermoanaerobacterium saccharolyticum strain B6A-Ri.** *J Gen Microbiol*
682 1993, **139 Pt 6**:1235-1243.
- 683 17. Lee YE, Lowe SE, Zeikus JG: **Regulation and Characterization of Xylanolytic**
684 **Enzymes of Thermoanaerobacterium-Saccharolyticum B6a-Ri.** *Appl Environ*
685 *Microbiol* 1993, **59(3)**:763-771.
- 686 18. Mai V, Wiegel J: **Advances in development of a genetic system for**
687 **Thermoanaerobacterium spp.: expression of genes encoding hydrolytic**
688 **enzymes, development of a second shuttle vector, and integration of genes**
689 **into the chromosome.** *Appl Environ Microbiol* 2000, **66(11)**:4817-4821.
- 690 19. Shaw AJ, Hogsett DA, Lynd LR: **Natural competence in *Thermoanaerobacter***
691 **and *Thermoanaerobacterium* species.** *Appl Environ Microbiol* 2010,
692 **76(14)**:4713-4719.
- 693 20. Shaw AJ, Covalla SF, Hogsett DA, Herring CD: **Marker removal system for**
694 ***Thermoanaerobacterium saccharolyticum* and development of a**
695 **markerless ethanologen.** *Appl Environ Microbiol* 2011, **77(7)**:2534-2536.
- 696 21. Dien BS, Cotta MA, Jeffries TW: **Bacteria engineered for fuel ethanol**
697 **production: current status.** *Appl Microbiol Biotechnol* 2003, **63(3)**:258-266.
- 698 22. Jarboe LR, Grabar TB, Yomano LP, Shanmugan KT, Ingram LO: **Development**
699 **of ethanologenic bacteria.** *Adv Biochem Eng Biotechnol* 2007, **108**:237-261.

- 700 23. Olofsson K, Bertilsson M, Liden G: **A short review on SSF - an interesting**
701 **process option for ethanol production from lignocellulosic feedstocks.**
702 *Biotechnol Biofuels* 2008, **1**(1):7.
- 703 24. Tsakraklides V, Shaw AJ, Miller BB, Hogsett DA, Herring CD: **Carbon catabolite**
704 **repression in *Thermoanaerobacterium saccharolyticum*.** *Biotechnol Biofuels*
705 2012, **5**(1):85.
- 706 25. Liu S, Lu H, Hu R, Shupe A, Lin L, Liang B: **A sustainable woody biomass**
707 **biorefinery.** *Biotechnol Adv* 2012, **30**(4):785-810.
- 708 26. **Final Report on Development of *Thermoanaerobacterium saccharolyticum***
709 **for the conversion of lignocellulose to ethanol.**
710 [<http://www.osti.gov/servlets/purl/1033560/>]
- 711 27. Bose M, Barber RD: **Prophage Finder: a prophage loci prediction tool for**
712 **prokaryotic genome sequences.** *In Silico Biol* 2006, **6**(3):223-227.
- 713 28. Bhaya D, Davison M, Barrangou R: **CRISPR-Cas systems in bacteria and**
714 **archaea: versatile small RNAs for adaptive defense and regulation.** *Annu*
715 *Rev Genet* 2011, **45**:273-297.
- 716 29. Olson DG, Lynd LR: **Transformation of *Clostridium thermocellum* by**
717 **electroporation.** *Methods Enzymol* 2012, **510**:317-330.
- 718 30. Brown SD, Nagaraju S, Utturkar S, De Tissera S, Segovia S, Mitchell W, Land
719 ML, Dassanayake A, Koepke M: **Comparison of single-molecule sequencing**
720 **and hybrid approaches for finishing the genome of *Clostridium***
721 **autoethanogenum and analysis of CRISPR systems in industrial relevant**
722 ***Clostridia*.** *Biotechnology for Biofuels* 2014, **7**.

- 723 31. Yin Y, Mao X, Yang J, Chen X, Mao F, Xu Y: **dbCAN: a web resource for**
724 **automated carbohydrate-active enzyme annotation.** *Nucleic Acids Res* 2012,
725 **40**(Web Server issue):W445-451.
- 726 32. Bogdahn M, Kleiner D: **N₂ fixation and NH₄⁺ assimilation in the thermophilic**
727 **anaerobes *Clostridium thermosaccharolyticum* and *Clostridium***
728 **thermoautotrophicum.** *Arch Microbiol* 1986, **144**(1):102-104.
- 729 33. Shimizu T, Ohshima S, Ohtani K, Hoshino K, Honjo K, Hayashi H: **Sequence**
730 **heterogeneity of the ten rRNA operons in *Clostridium perfringens*.** *Syst Appl*
731 *Microbiol* 2001, **24**(2):149-156.
- 732 34. Beauclerk AA, Cundliffe E: **Sites of action of two ribosomal RNA methylases**
733 **responsible for resistance to aminoglycosides.** *J Mol Biol* 1987, **193**(4):661-
734 671.
- 735 35. Recht MI, Fourmy D, Blanchard SC, Dahlquist KD, Puglisi JD: **RNA sequence**
736 **determinants for aminoglycoside binding to an A-site rRNA model**
737 **oligonucleotide.** *J Mol Biol* 1996, **262**(4):421-436.
- 738 36. Kotra LP, Haddad J, Mobashery S: **Aminoglycosides: perspectives on**
739 **mechanisms of action and resistance and strategies to counter resistance.**
740 *Antimicrob Agents Chemother* 2000, **44**(12):3249-3256.
- 741 37. Yang Y, Zhu M, Wu L, Zhou J: **Assessment of data processing to improve**
742 **reliability of microarray experiments using genomic DNA reference.** *Bmc*
743 *Genomics* 2008, **9 Suppl 2**:S5.

- 744 38. Roberts SB, Gowen CM, Brooks JP, Fong SS: **Genome-scale metabolic**
745 **analysis of *Clostridium thermocellum* for bioethanol production.** *BMC Syst*
746 *Biol* 2010, **4**:31.
- 747 39. Wang S, Huang H, Moll J, Thauer RK: **NADP⁺ reduction with reduced**
748 **ferredoxin and NADP⁺ reduction with NADH are coupled via an electron-**
749 **bifurcating enzyme complex in *Clostridium kluyveri*.** *J Bacteriol* 2010,
750 **192**(19):5115-5123.
- 751 40. Shaw AJ, Hogsett DA, Lynd LR: **Identification of the [FeFe]-hydrogenase**
752 **responsible for hydrogen generation in *Thermoanaerobacterium***
753 ***saccharolyticum* and demonstration of increased ethanol yield via**
754 **hydrogenase knockout.** *J Bacteriol* 2009, **191**(20):6457-6464.
- 755 41. Burgard AP, Pharkya P, Maranas CD: **Optknock: a bilevel programming**
756 **framework for identifying gene knockout strategies for microbial strain**
757 **optimization.** *Biotechnol Bioeng* 2003, **84**(6):647-657.
- 758 42. Hogsett DA: **Cellulose hydrolysis and fermentation by *Clostridium***
759 ***thermocellum* for the production of ethanol.** Hanover: Dartmouth College;
760 1995.
- 761 43. Brown SD, Raman B, McKeown CK, Kale SP, He ZL, Mielenz JR: **Construction**
762 **and evaluation of a *Clostridium thermocellum* ATCC 27405 whole-genome**
763 **oligonucleotide microarray.** *Appl Biochem Biotechnol* 2007, **137**:663-674.
- 764 44. Chhabra SR, He Q, Huang KH, Gaucher SP, Alm EJ, He Z, Hadi MZ, Hazen TC,
765 Wall JD, Zhou J *et al*: **Global analysis of heat shock response in**
766 ***Desulfovibrio vulgaris* Hildenborough.** *J Bacteriol* 2006, **188**:1817–1828.

- 767 45. Li X, He Z, Zhou J: **Selection of optimal oligonucleotide probes for**
768 **microarrays using multiple criteria, global alignment and parameter**
769 **estimation.** *Nucleic Acids Res* 2005, **33**:6114-6123.
- 770 46. Mukhopadhyay A, He Z, Alm EJ, Arkin AP, Baidoo EE, Borglin SC, Chen W,
771 Hazen TC, He Q, Holman H-Y *et al*: **Salt Stress in *Desulfovibrio vulgaris***
772 **Hildenborough: an Integrated Genomics Approach.** *Journal of Bacteriology*
773 2006, **188**(11):4068-4078.
- 774 47. Gowen CM, Fong SS: **Genome-scale metabolic model integrated with**
775 **RNAseq data to identify metabolic states of *Clostridium thermocellum*.**
776 *Biotechnol J* 2010, **5**(7):759-767.
- 777 48. Brooks JP, Burns WP, Fong SS, Gowen CM, Roberts SB: **Gap detection for**
778 **genome-scale constraint-based models.** *Adv Bioinformatics* 2012,
779 **2012**:323472.
- 780 49. Edwards JS, Ramakrishna R, Palsson BO: **Characterizing the metabolic**
781 **phenotype: a phenotype phase plane analysis.** *Biotechnol Bioeng* 2002,
782 **77**(1):27-36.
- 783 50. Schellenberger J, Que R, Fleming RM, Thiele I, Orth JD, Feist AM, Zielinski DC,
784 Bordbar A, Lewis NE, Rahmanian S *et al*: **Quantitative prediction of cellular**
785 **metabolism with constraint-based models: the COBRA Toolbox v2.0.** *Nat*
786 *Protoc* 2011, **6**(9):1290-1307.
- 787 51. Becker SA, Feist AM, Mo ML, Hannum G, Palsson BO, Herrgard MJ:
788 **Quantitative prediction of cellular metabolism with constraint-based**
789 **models: the COBRA Toolbox.** *Nat Protoc* 2007, **2**(3):727-738.

790 52. R Core Team: **R: A language and environment for statistical computing**. In.
791 Vienna, Austria; 2012.

792 53. Gregory R. Warnes BB, Lodewijk Bonebakker, Robert Gentleman,, Wolfgang
793 Huber Andy Liaw TL, Martin Maechler, Arni, Magnusson SM, Marc Schwartz and
794 Bill Venables: **gplots: Various R programming tools for plotting data**. In.,
795 2.11.0 edn; 2012.

796 54. Zuker M: **Mfold web server for nucleic acid folding and hybridization**
797 **prediction**. *Nucleic Acids Res* 2003, **31**(13):3406-3415.

798

799

800

801 **Figure 1.** A comparison between the two versions of the 16s mRNA found in *T.*
802 *saccharolyticum*. A) an alignment and consensus sequence for a heterogeneous
803 segment of the five 16S ribosomal components found in *T. saccharolyticum*. B) Mfold
804 prediction of the structure of the shorter 16S mRNA. [54]. C) Mfold prediction of the
805 structure of the longer 16S mRNA.

806

807 **Figure 2.** Time points between 5 and 60 minutes post-shock with hemicellulose extract.
808 The horizontal axis represents \log_2 of the control xylose+acetate expression level
809 (mRNA:gDNA ratio), while the vertical axis represents the hemicellulose extract-treated
810 expression level. All data are the average of duplicate experiments with the exception of
811 the 5 minutes post hemicellulose extract shock which is in triplicate.

812

813 **Figure 3.** Heat map of hierarchical clustering of genes that change in expression level
814 upon the addition of washate with a P value of <0.01 and with a \log_2 ratio >1.0 in at
815 least one time point. The range of \log_2 mRNA:gDNA ratios is given in the color key.

816

817 **Figure 4.** Example of data from Nimbegen tiled microarrays (bottom) showing
818 transcription units correlated to open reading frames (top).

819

820 **Figure 5.** Inhibitor shock. A) Plot showing the addition of HMF and furfural in culture
821 supernatants and the temporary disruption of growth. B) Plot showing the levels of
822 intracellular citric acid and hydroxymethylfurfural, as well as the average of all other

823 metabolites. C) A heat map of a hierarchical clustering of the concentration of all
824 monitored intracellular metabolites over the course of the 4 hour experiment.

825

826 **Figure 6.** Phenotypic phase planes for *T. saccharolyticum* high-ethanol knock out
827 strains. The maximum growth rate is shown as a surface over a range of fluxes for
828 glucose uptake and ethanol production. The wild-type surface (A) shows the maximum
829 growth rate occurring equally across a wide range of ethanol production rates, while the
830 phase planes for the $\Delta ldh-ptd$ strain (B) and the $\Delta ldh-hfs$ strain (C) demonstrate that the
831 potential solution space is trimmed in a way that couples maximum growth to high
832 ethanol yield.

833

834 **Figure 7.** Growth envelope for various ethanol strain designs during growth on glucose.
835 $\Delta LDH-\Delta HFS$ and $\Delta HFS-\Delta LDH-\Delta GLUD$ were both identified by OptKnock as being
836 optimal designs for ethanol production.

837

838 **Table 1.** Summary of microarray and metabolomics data sets.

839

840 **Table 2.** A comparison between the number of components in the models generated for
841 *T. saccharolyticum* and *C. thermocellum*.

842

843 **Table 3.** Relevant reactions in ethanol producing knockout strain designs.

| Strain | Medium | Condition | Microarray timepoints analyzed | Metabolite timepoints analyzed |
|---------------|---------------|---------------------------------|---------------------------------------|---------------------------------------|
| ALK2 | MTC | arabinose | 1 | |
| ALK2 | MTC | cellobiose | 2 | 5 |
| ALK2 | MTC | cellobiose: nitrogen-limited | 4 | 4 |
| ALK2 | MTC | cellobiose fed-batch: N-limited | 4 | 5 |
| ALK2 | MTC | enzymatic hydrolysate | 1 | |
| ALK2 | MTC | glucose-arabinose-mannose | 4 | |
| ALK2 | MTC | glucose-xylose | 7 | |
| ALK2 | MTC | glucose-xylose-cellobiose | 1 | |
| ALK2 | MTC | glucose-xylose-ethanol | 1 | |
| ALK2 | MTC | glucose-xylose-acetate shock | 3 | |
| ALK2 | MTC | glucose-xylose-washate shock | 3 | |
| ALK2 | MTC | pretreated hardwood SSF | 3 | |
| ALK2 | MTC | xylose fed-batch | 4 | 5 |
| M0355 | MTC | glucose-xylose | 1 | |

| | | | | |
|-----------|------|-----------------------------|---|---|
| M0355 | MTC | glucose-xylose-ethanol | 1 | |
| M0355 | MTC | pretreated hardwood SSF | 1 | |
| M0521 | MTC | pretreated hardwood SSF | 2 | |
| M0700 | MTC | glucose-xylose | 2 | |
| M0700 | MTC | glucose-xylose-ethanol | 2 | |
| | | glucose-xylose HMF+furfural | | |
| M0700 | TSD | shock | 5 | 5 |
| M1151 | TSC3 | cellobiose-maltodextrin | 4 | 4 |
| | | xylose-detoxified washate | | |
| M1151 | TSC4 | shock | 7 | 3 |
| M1291/144 | | | | |
| 2 | TSC4 | sigmacell SSF | 2 | 6 |
| M1732 | TSC7 | xylose-diamide shock | 6 | 5 |

846

847

848

849

850 **Table 2.**

| | <i>T. saccharolyticum</i> | <i>C. thermocellum</i> ^a |
|--------------------------------------|---------------------------|-------------------------------------|
| Genes | 315 | 432 |
| Metabolites | 503 | 525 |
| Reactions | 515 | 577 |
| - Gene-associated | 461 (90%) | 463 (80%) |
| - Biomass | 1 | 1 |
| - Non-gene associated [cytosolic] | 43 | 60 |
| - Non-gene associated [transport] | 11 | 54 |

851 ^a[38]

852

853

854

855 **Table 3.**

| ID | Reaction Name | Formula | Gene | Association |
|--------------|-------------------------------------|---|--|--------------------|
| GLNS | Glutamine synthetase | glu-L[c] + ATP[c] + NH4[c] - > ADP[c] + Pi[c] + H[c] + gln- L[c] | Tsac_2029 | |
| GLUD | Glutamate y dehydrogenase (NADP) | NADP[c] + H2O[c] + glu-L[c] <=> H[c] + NADPH[c] + NH4[c] + akg[c] | Tsac_2172 | |
| GLUSy | Glutamate synthase(NADPH) | H[c] + NADPH[c] + gln-L[c] + akg[c] -> NADP[c] + 2 glu- L[c] | Tsac_1234 | |
| LDH_L | L-lactate dehydrogenase | lac-L[c] + NAD[c] <=> NADH [c] + H[c] + pyr[c] | Tsac_0179 | |
| PTAr | Phosphotransacetylase | Pi[c] + AcCoA[c] <=> CoA[c] + actp[c] | Tsac_1744 | |
| HFS | Ferredoxin hydrogenase | Fdred[c] + 2 h[c] <==> Fdox + H2[c] | Tsac_1550 & Tsac_1551 & Tsac_1552 & Tsac_1553 | |

856