

Design of synthetic human gut microbiome assembly and function

Ryan L. Clark¹, Bryce M. Connors^{1,3}, David M. Stevenson², Susan E. Hromada^{1,2}, Joshua J. Hamilton¹, Daniel Amador-Noguez² & Ophelia S. Venturelli^{1,2,3*}

¹Department of Biochemistry, University of Wisconsin-Madison, Madison, WI 53706

²Department of Bacteriology, University of Wisconsin-Madison, Madison, WI 53706

³Department of Chemical & Biological Engineering, University of Wisconsin-Madison, Madison, WI 53706

*To whom correspondence should be addressed: venturelli@wisc.edu

1 ABSTRACT

2 The assembly of microbial communities and functions emerge from a complex and dynamic web
3 of interactions. A major challenge in microbiome engineering is identifying organism
4 configurations with community-level behaviors that achieve a desired function. The number of
5 possible subcommunities scales exponentially with the number of species in a system, creating
6 a vast experimental design space that is challenging to even sparsely traverse. We develop a
7 model-guided experimental design framework for microbial communities and apply this method
8 to explore the functional landscape of the health-relevant metabolite butyrate using a 25-member
9 synthetic human gut microbiome community. Based on limited experimental measurements, our
10 model accurately forecasts community assembly and butyrate production at every possible level
11 of complexity. Our results elucidate key ecological and molecular mechanisms driving butyrate
12 production including inter-species interactions, pH and hydrogen sulfide. Our model-guided
13 iterative approach provides a flexible framework for understanding and predicting community
14 functions for a broad range of applications.

15 INTRODUCTION

16 Microbial communities carry out pivotal chemical transformations in nearly every environment on
17 Earth¹. Many of these processes critically impact human health and environmental sustainability,
18 including oceanic CO₂-fixation², production of growth-promoting molecules in the plant
19 rhizosphere³, and degradation of indigestible dietary substrates⁴. Microbial community dynamics
20 and functions are determined by complex and dynamic interactions between constituent
21 community members and their environment. Developing the capabilities to engineer microbiome
22 properties holds promise to address grand challenges facing human society⁵ and methods to
23 predict microbiome functions are needed to enable microbiome engineering efforts.

24 A bottom-up approach to build and characterize synthetic microcosms has key
25 advantages including reduced complexity compared to natural systems, ability to manipulate
26 environmental parameters and community membership and achieve a high temporal resolution.
27 Previous studies have leveraged synthetic microcosms of bacteria isolated from the human gut⁶
28 or soil⁷ to demonstrate that dynamic models based on pairwise interactions are predictive of multi-
29 species community assembly. In addition, modeling community assembly using pairwise
30 interactions has provided a deeper understanding of the effects of environmental factors including
31 pH⁸, dilution⁹, nutrient availability¹⁰, toxins¹¹, and temperature¹² on microbial community
32 behaviors.

33 A key challenge for predicting microbiome properties is mapping community composition
34 to community-level metabolic functions. Genome-scale metabolic models have been used to
35 predict collective metabolic outputs of microbial communities, an approach which is limited by the
36 quality of functional gene annotations and stringent assumptions¹³. Bottom-up assembly of
37 microbial consortia coupled to mathematical modeling has been used to interrogate how the
38 production or consumption of molecules changes in a community context relative to individual
39 species^{11,14,15}. However, computational frameworks to predict both community dynamics and
40 functional outputs for high-dimensional communities that mirror the complexity of natural
41 microbiomes are needed to harness the potential of microbiome engineering for diverse
42 applications.

43 A detailed and quantitative understanding of microbial interaction networks would enable
44 the design of microbial consortia with robust target functions from the bottom up. While model-
45 guided design¹⁶ has been used to identify gut microbial communities that elicit a target immune
46 response in mouse models, the complexity of the host system and the low-throughput of mouse
47 studies limits the observability of system parameters and a comprehensive understanding of

48 ecological factors shaping microbiome behaviors. Here, we use a data-driven approach to build
49 a model of butyrate production by complex *in vitro* communities of human-associated intestinal
50 isolates. Butyrate production is a major function of the gut microbiome associated with protection
51 from a wide range of human diseases, including arthritis¹⁷, diet-induced obesity^{18–20}, colitis^{21,22},
52 opportunistic pathogen infection²³, diabetes²⁴, and colorectal cancer²⁵. Our approach leverages
53 data-driven models to quantify interactions impacting the growth dynamics of functional
54 organisms and interpretable statistical models to quantify interactions impacting metabolic
55 activities (functional yield of butyrate per unit biomass). By modeling these two interaction types
56 separately, we demonstrate that in some contexts, accurate prediction of functional organism
57 abundance can predict function, while in other community contexts containing metabolically
58 flexible ecological driver species, interactions modifying metabolic modes must be captured to
59 predict function. We use these models to design communities of up to 25 species with a broad
60 range of butyrate production capabilities and analyze our model as well as the metabolic profiles
61 and environmental modification of designed consortia to provide key insights into metabolic
62 interactions impacting butyrate production.

63 RESULTS

64 Identifying highly functional microbial communities from the bottom-up is a major challenge
65 because the number of sub-communities exponentially increases with the system dimension²⁶.
66 To explore community design space, we develop a modeling framework to guide iterative design
67 of experiments (**Figure 1a,b**). Ecosystem functions can be modulated by selection effects,
68 defined as changes in function correlated with changes in the abundance of functional species,
69 or complementarity effects, defined as changes in the functional yield per unit biomass for each
70 functional species^{27–29}. We implement a dual module modeling framework to determine the
71 contributions of microbe-microbe interactions to each of these effect types. A community dynamic
72 model, referred to as the generalized Lotka-Volterra model (gLV), predicts community assembly
73 and the function model predicts a functional activity from community composition (**Figure 1a**).
74 The gLV model is an ordinary differential equation model that captures the temporal change in
75 species abundances due to monospecies growth parameters and inter-species interactions and
76 has been used to predict and analyze multi-species community assembly based on
77 measurements of lower-order communities⁶. Our function model consists of a regression model
78 with interaction terms mapping species abundance at a specific time point to the concentration of
79 an output metabolite. The inter-species interaction terms in the gLV model represent selection
80 effects (i.e. how one species impacts the growth of another) and the interaction terms in the
81 regression model represent complementarity effects (i.e. deviations from constant yield of a
82 function per unit biomass¹⁵) (**Figure 1a,c**). For the gLV model, we use Bayesian parameter
83 inference techniques to determine the uncertainty in our parameters based on biological and
84 technical variability in the experimental data³⁰. The composite gLV and statistical models predict
85 the probability distribution of the functional activity given an initial condition of species abundances
86 (**Figure 1a, Methods**).

87 Due to significant interest in development of defined bacterial therapeutics for human
88 health applications³¹ and the beneficial role of butyrate produced by gut microbiota on a myriad
89 of health outcomes^{17–25,32}, we sought to apply our modeling framework to understand how
90 community composition impacts butyrate production in synthetic communities of prevalent and
91 diverse human gut microbes. Butyrate production is a specialized function of a subset of species
92 in the gut (~10-25% of microbial genomes are predicted to harbor this pathway in healthy
93 individuals³³). By contrast, the production of other metabolic byproducts such as acetate and
94 lactate are distributed more broadly across members of the gut microbiome. Thus, studying
95 butyrate production allows us to investigate how ecological forces shape a function performed
96 only by a specific subset of the community. Indeed, predictably modulating a specific function

97 performed by a subset of organisms constitutes a core goal of microbiome engineering across
98 different environments^{34–36}.

99 To develop a system of microbes representing major metabolic functions in the gut, we
100 selected 25 highly prevalent bacterial species from all major phyla in the human gut microbiome³⁷.
101 This community contained 5 butyrate producing Firmicutes which have been shown to play
102 important roles in human health and protection from diseases (**Figure 2a, Table S1**). These 5
103 Firmicutes have the capability to ferment sugars and/or transform acetate to butyrate, allowing
104 the recovery of NAD⁺ for further energy generation³⁸. Additionally, *Anaerostipes caccae* (AC) can
105 ferment lactate and acetate to butyrate, generating a modest amount of energy³⁹. However, each
106 of these species can alternatively produce acetate and/or lactate as fermentation products
107 depending on the environmental context (**Figure 2b**).

108 Due to lack of defined media that universally support growth of gut microbes, most *in vitro*
109 studies use rich media, making it difficult to interrogate the effects of unknown components on
110 community behaviors⁴⁰. To maximize our knowledge of the substrates available to the
111 communities in our experiments and to simplify the metabolite quantification, we developed a
112 single chemically defined medium that supports the growth of all species in monoculture with the
113 exception of *Faecalibacterium prausnitzii* (FP) (**Methods**). We measured time-resolved growth of
114 each species and constructed a gLV null model that assumed no inter-species interactions. Our
115 results demonstrated a wide variety of growth dynamics within each phylum, including disparate
116 growth rates and carrying capacities (**Figure 2c**). Using this system, we implemented an iterative
117 design, test, learn (DTL) cycle (**Figure 1b**) to explore a vast community design space and explore
118 an ecosystem functional landscape.

119 *Butyrate production impacted by selection and complementarity*

120 For the first cycle of our iterative DTL approach, we sought to decipher interactions impacting
121 butyrate production in pairwise communities, with the goal of understanding how these
122 interactions combine to shape community assembly and butyrate production in higher complexity
123 communities. We grew each pairwise community containing at least one butyrate producer (the
124 focal species of our system⁴¹) and measured species abundance and the concentrations of
125 organic acid fermentation products (including butyrate, lactate, succinate and acetate) after 48
126 hours. Based on previous studies using pairwise communities to predict higher complexity
127 community behaviors^{6,14,42}, we hypothesized that these measurements would provide a highly
128 informative dataset to develop an initial model that captured inter-species interactions shaping
129 selection and complementarity effects in the system.

130 Based on our data, we first considered to what extent butyrate production was impacted
131 by selection effects and complementarity effects using a model-free approach (**Figure 2d**). For
132 each pairwise community, the selection effect was computed as the difference between the
133 expected butyrate concentration assuming constant butyrate yield and the monoculture butyrate
134 concentration (**Figure S1**). The complementarity effect was defined as the difference between the
135 measured butyrate concentration of the community and the expected butyrate concentration
136 assuming constant yield (**Figure S1, Methods**). Negative selection effects influenced all butyrate
137 producers except FP, which did not grow in monoculture (**Figure 2d, inset**). Compared to the
138 other butyrate producers, *Roseburia intestinalis* (RI) exhibited the largest negative selection
139 effects, while AC tended to display positive complementarity effects. In sum, both selection and
140 complementarity can modulate butyrate production, highlighting the utility of a building a
141 composite model that captures both types of effects (**Figure 1a**). Further, the model-free
142 approach to determining the contributions of selection and complementarity effects cannot be
143 applied to communities containing multiple butyrate-producers. Therefore, our modeling approach
144 can elucidate the selection and complementarity effects in communities with functional
145 redundancy, representing real systems³³.

146 To enable prediction of butyrate concentration in higher complexity communities, we used
147 the data from the monoculture and pairwise community experiments to train our model (M1). The
148 inferred gLV inter-species interaction network showed many negative interactions, including
149 strong negative interactions impacting the growth of RI, consistent with the dominance of negative
150 selection effects in our model-free analysis of RI pairwise communities (**Figure 2e**). A network
151 representation of the parameters in the regression model indicated that AC had significantly more
152 pairwise interaction terms than the other butyrate producers, consistent with the major role of
153 positive complementarity effects in our model-free analysis of AC pairwise communities (**Figure**
154 **2f**).

155 *Model trained on pairwise consortia predicts 3-5 species community behaviors*

156 To test our model's ability to predict function in communities with an incremental increase in
157 complexity, we implemented a second DTL cycle with the goal of mapping the functional
158 landscape of 3-5 member communities. Model M1 was informed only by pairwise communities
159 that contained at least one butyrate producer and was thus naive to all interactions between non-
160 butyrate producers. Therefore, we needed to make some assumptions to enable prediction of
161 multi-species consortia containing combinations of non-butyrate producers. Based on patterns
162 observed in previous gLV model parameter sets⁶, we hypothesized that unmeasured interactions
163 could be estimated based on the trends in measured interactions across phylogenetic
164 relatedness. Therefore, we used a matrix imputation method to estimate interaction parameters
165 for unmeasured interactions in the gLV model (**Methods**). The resulting model was used to predict
166 the probability distributions of butyrate production for all 3-5 species communities containing at
167 least one butyrate producer (46,588 communities). The predicted butyrate production varied
168 substantially between the combinations of butyrate-producer groups (**Figure 3a**). To evaluate the
169 ability of our model to predict the behaviors of butyrate producers in a variety of community
170 contexts, we experimentally characterized 156 communities that spanned a broad range of
171 predicted butyrate concentrations across the butyrate-producer groups (**Figure 3a**). The model
172 prediction exhibited good agreement with the rank order of butyrate production (Spearman
173 $\rho=0.84$, $p=9*10^{-43}$), though moderately overpredicted the magnitude on average (**Figure 3b**).

174 The quality of these predictions demonstrated that our initial dataset was sufficient to build
175 a model predicting broad trends in butyrate production but suggested that additional data was
176 required to predict specific outliers. To understand the factors contributing to deviations between
177 predicted and measured butyrate concentrations, we updated our models to decipher key inter-
178 species interactions that model M1 failed to capture yielding model M2. The gLV model from M2
179 contained many new negative interaction parameters (46 new negative interactions, 15 conserved
180 negative interactions) and sparse positive interactions (3 new and 2 conserved positive
181 interactions) out of 386 possible observed interspecies interaction parameters in the designed
182 set, primarily between non-butyrate producers (**Figure 3c**, **Figure S2**). The regression model in
183 M2 highlighted significant complementarity effects in the 3-5 member communities, with strong
184 negative interactions between *Desulfovibrio piger* (DP)-AC and AC-*Eubacterium rectale* (ER)
185 consistent with model M1 (**Figure 2f**). We used the regression model with the experimental
186 abundance measurements to quantify the magnitude and variability of each complementarity
187 interaction across the experimentally measured communities due to differences in species
188 growth. These data showed that some interactions consistently modified butyrate production in
189 the presence of the species pair (e.g. DP-AC), whereas the contributions of other interactions to
190 butyrate production varied across communities (e.g. ER-RI, *Clostridium asparagiforme* (CG)-AC)
191 (**Figure 3d**). Equipped with this updated model, we set out to explore our model's experimental
192 design capabilities for communities of even greater complexity (i.e. >10 species).

193 *Model-guided exploration of complex community design space*

194 One approach to determining the contributions of constituent community members to ecosystem
195 behaviors involves characterization of the full and all single-species dropout consortia^{6,26}. In our
196 system, all 24 and the 25 member communities exhibited similar low butyrate production (~10-15
197 mM Butyrate), except for a moderate increase of butyrate in the DP-lacking community (~22 mM
198 Butyrate) and a large decrease in the AC-lacking community (~2 mM Butyrate) (**Figure S3**). Many
199 of the 3-5 member communities displayed higher butyrate production than the highest complexity
200 communities (**Figure 3b**), suggesting that high complexity communities may trend towards an
201 undesired low butyrate producing state. Additionally, the concentrations of all measured organic
202 acids spanned a much smaller range in the 24 and 25 member communities than the low
203 complexity (1-5 member) assemblages (**Figure S3b**), further supporting the notion that
204 communities of increasing complexity may trend toward a similar functional state. Indeed, a key
205 challenge in engineering microbial communities is the tendency to assemble to compositional
206 attractors and resist change, due to a multitude of abiotic and biotic interactions⁴³⁻⁴⁶. In our
207 system, implementing a standard approach of analyzing the highest complexity set of consortia
208 failed to elucidate diverse community metabolic states, highlighting the utility of the model to
209 design sub-communities that span the functional space.

210 To address this challenge, we used our model M2 to design complex communities (>10
211 species) that deviated from the observed trend towards low butyrate production. Since the human
212 gut microbiome exhibits functional redundancy in butyrate pathways³³, we used model M2 to
213 simulate the assembly of all communities containing all five butyrate producers to map species
214 abundance to butyrate concentration (1,048,575 communities). Based on the hypothesis that
215 high-complexity communities may trend towards low butyrate production, we found it useful to
216 consider the full community as a reference frame, representing a potential compositional attractor
217 state, when visualizing the relationship between community composition and butyrate production.
218 Consistent with this notion, the model predicted the full community to have butyrate production
219 similar to the average of all communities, with other communities diverging from this average
220 behavior with increasing distance in composition (Euclidean distance between endpoint
221 abundances) from the full community (**Figure 4a**). The landscape of communities was partitioned
222 into two large clusters based on the presence or absence of the prevalent sulfate-reducing
223 Proteobacteria DP⁴⁷. Corroborating these results, DP had the strongest negative complementarity
224 interaction with AC and CC in the regression model of M2 as well as a significant negative impact
225 on butyrate production in the single-species dropout consortia (**Figure 3d, Figure S3a**). This
226 inferred complementarity effect and predicted shift in the butyrate production landscape suggests
227 that the presence of DP may substantially alter the metabolic activities shaping butyrate
228 production.

229 We evaluated the capability of our model to guide broader exploration of the functional
230 landscape and identify infrequent communities that deviate from the typical behavior by designing
231 28 low- and 54 high-butyrate communities each with 11-17 members and containing all 5 butyrate-
232 producers. In addition, we randomly selected 82 communities with the same complexity
233 constraints to evaluate whether our model-guided design procedure could elucidate a set of
234 communities that spanned a broader range of metabolic states, scored by the variance in butyrate
235 concentration (**Figure 4a**). The 82 designed communities exhibited a higher variance in mean
236 butyrate production than the 82 random communities (designed communities s.d.=11 mM,
237 random communities s.d.=8 mM, Levene test, p=0.043), demonstrating a major advantage of the
238 model-guided approach for designing communities to broadly explore regions of the functional
239 landscape (**Figure 4c**). Consistent with our model predictions, communities containing DP
240 exhibited lower butyrate compared to communities excluding DP in both the designed and the
241 randomly chosen communities (**Figure 4c**). While the model predicted the rank order of butyrate

242 concentrations in these communities moderately well (Spearman $\rho=0.67$, $p=3*10^{-25}$), some of
243 the highest butyrate production communities were underpredicted by the model (**Figure 4c**).

244 *Selection effects dominate in high complexity communities lacking AC*

245 In microbial consortia, the contributions of individual members to a given function can be broadly
246 distributed, wherein key driver species can exhibit a substantially larger contribution to
247 community-level functions than the other members^{6,11,14}. In our system, the 24-member
248 community lacking AC exhibited 1.9 ± 1.0 mM butyrate (mean \pm s.d., $n=8$), substantially lower than
249 any observed complex (>10 species) community containing AC and qualitatively consistent with
250 our model which predicted no butyrate (**Figure 4b**, **Figure S3b**). Therefore, AC was a driver of
251 butyrate production in complex communities. There are large interpersonal differences in gut
252 microbiota composition due to environmental factors and host-microbe interactions⁴⁸. Thus, some
253 species such as AC may not be present in certain individuals⁴⁹. To evaluate the capability of
254 our model to steer systems from low to high butyrate producing states independent of the
255 presence of particular species, we designed high butyrate producing complex communities
256 lacking the driver species AC.

257 To do so, we used model M2 to simulate all communities containing the four butyrate
258 producers excluding AC (1,048,575 communities) to forecast species abundance and butyrate
259 production. Similar to the 5 butyrate-producer case, we used the 24-member community lacking
260 AC as a reference frame for quantifying deviations from a potential compositional attractor. While
261 most communities were predicted to have low butyrate production, butyrate production increased
262 with the distance from the 24-member community (**Figure 4b**). To explore this design space and
263 evaluate whether our model could identify communities with low or high butyrate activity, we
264 experimentally assembled 84 communities containing 11-19 members that were predicted to
265 display a broad range of butyrate production capabilities (**Figure 4b**). Mirroring our model
266 prediction, distance from the 24-member community in species composition was positively
267 correlated with butyrate production (Spearman $\rho=0.56$, $p=3*10^{-8}$) as well as butyrate producer
268 abundance (Spearman $\rho=0.85$, $p=1*10^{-24}$) (**Figure 4b**, inset). However, the model substantially
269 overpredicted butyrate production (**Figure 4d**). Therefore, we sought to continue the DTL
270 paradigm by training our model on the new data.

271 *Updated model predicts butyrate production in high complexity communities*

272 The discrepancies between our model predictions and experimental measurements in complex
273 communities were either due to missing information about certain pairwise interactions (i.e. poor
274 parameter estimates due to unobservable interactions) or higher-order interactions that could not
275 be captured by our pairwise model (i.e. model structure fails to represent system behaviors). To
276 distinguish between these possibilities, we updated our model by training on a subset of high-
277 complexity communities: the random set of communities containing all butyrate producers (82
278 communities) and a randomly sampled half of the communities lacking AC (42 communities)
279 (**Figure 4a-d**). Notably, the updated model M3 predicted the measurements of high-complexity
280 communities with high accuracy, demonstrating that the pairwise model structures could explain
281 the quantitative trends in the data when provided with sufficient information (**Figure 4e**). The
282 predictive capability of the model required information from complex communities, supporting
283 recent theoretical work suggesting that the typical pairwise community experimental design may
284 not be the most efficient for building predictive models of complex systems²⁶.

285 We next examined the changes in the inferred parameters between models M2 and M3
286 to provide insights into key microbial interactions impacting complex community behaviors. The
287 major changes in the updated gLV M3 model were new values for all previously unobserved
288 pairwise interactions as well as modification of previously observed interaction parameters
289 (**Figure 4f**, **Figure S2f**). Negative interactions (<-0.05 hr⁻¹ (OD₆₀₀ Species j)⁻¹) dominated the

290 network, representing 49.8% of the interspecies interaction parameters. By contrast, only 1.7%
291 of interactions were strong positive ($>0.05 \text{ hr}^{-1} (\text{OD}_{600} \text{ Species } j)^{-1}$), consistent with previous
292 observations of the prevalence of negative interactions in microbial communities^{6,50}. Notably, 70%
293 of the previously observed interspecies interaction parameters fell within the 60% confidence
294 interval of the posterior distribution for model M2, demonstrating that our M2 model was accurate
295 but lacked sufficient information to be highly confident in the estimated parameter values.

296 In the updated gLV model, species which secreted lactate in monoculture tended to have
297 a positive impact on the growth of AC (**Figure 4f**). Although DP is also a lactate consumer⁴⁷, it
298 did not tend to benefit from monospecies lactate producers. This result highlights the benefits of
299 using data from multiple levels of community complexity for training gLV models as these
300 interactions were not captured by models M1 and M2, trained only on lower-order community
301 contexts. To understand how inter-species interactions vary across chemical composition
302 contexts, we compared the inferred inter-species interaction coefficients in the M3 gLV model to
303 those from a previous study that used a gLV model to study a 12-member subset of our
304 community (PC, BV, BO, BT, BU, DP, CA, EL, FP, CH, BH, and ER) in a different (rich) media⁶
305 and found that 27 parameters with magnitude $>0.1 \text{ hr}^{-1} (\text{OD}_{600} \text{ species } j)^{-1}$ shared a sign and only
306 5 had opposite sign (**Figure 4f**). The high percentage (84%) of qualitatively consistent interaction
307 coefficients inferred based on measurements in two different environmental contexts provides
308 confidence in using parameterized gLV models as prior information to forecast system behaviors
309 in new environments.

310 The regression model from M3 identified three interactions driving complementarity effects
311 in the 5 butyrate producer communities including *Eggerthella lenta* (EL)-AC, DP-AC, and *Dorea*
312 *formicigenerans* (DF)-RI (**Figure 4g**). In the absence of AC, substantial complementarity effects
313 were not detected (**Figure 4g**), consistent with the absence of strong complementarity effects in
314 lower-complexity communities lacking AC. In our system, AC has the unique capability to
315 transform lactate to butyrate in addition to production of butyrate from sugars (**Figure 2b**),
316 suggesting that metabolic flexibility may be a key determinant of complementarity effects. In sum,
317 our modeling framework representing pairwise interactions could accurately predict community
318 composition and butyrate production in complex communities and could be used to decipher key
319 microbial interactions impacting metabolic outputs.

320 *Mechanistic insights identified from inferred interaction networks*

321 We sought to analyze the patterns in our inferred interactions to identify mechanistic hypotheses
322 about the potential ecological and molecular factors driving butyrate production. The low butyrate
323 productivity of specific communities could stem from a global reduction in metabolic activities for
324 the conversion of sugars to organic acid fermentation products. However, the amount of total
325 carbon in acetate, lactate, and propionate was inversely proportional to the amount of carbon in
326 butyrate in complex communities (**Figure S4**), indicating that metabolic tradeoffs dictated the
327 production of specific organic acids. Therefore, we considered how interactions identified by our
328 model could influence such tradeoffs.

329 We analyzed the inferred interaction networks to provide generalizable insights into
330 metabolic processes impacting butyrate production in our system. We first considered the largest
331 negative complementarity effect in our system between AC and DP (**Figure 4g**). While these two
332 species have previously been shown to compete for lactate *in vitro*⁵¹, excess lactate was present
333 in communities containing both DP and AC, suggesting that competition over limited lactate was
334 not a major determinant of the negative complementarity effect (**Figure 5a**). In addition, a large
335 negative complementarity effect was observed in the 3-5 member communities between DP and
336 CC, which does not utilize lactate for butyrate production (**Figure 3d**).

337 Since some *Desulfovibrio* species have the capability to use butyrate as an energy
338 source⁵², we tested whether decreased butyrate in the presence of DP could be due to butyrate

339 consumption. To investigate this hypothesis, we grew DP in media supplemented with different
340 concentrations of sodium butyrate ranging between 0 and 100 mM and measured the butyrate
341 concentration after 48 hours of incubation. The presence of DP did not alter the concentration of
342 butyrate in any condition, suggesting that decreased butyrate due to consumption or degradation
343 was not a major factor contributing to the negative complementarity effects associated with DP
344 (**Figure S5**). One unique metabolic characteristic of DP in our system is the capability to reduce
345 sulfate to hydrogen sulfide (H₂S). Therefore, we hypothesized that H₂S may contribute the
346 negative impact of DP on butyrate production (**Figure 4a**). To test this hypothesis, each butyrate
347 producer was grown in media supplemented with a range of sulfide concentrations. Notably, the
348 butyrate production per unit biomass decreased with increasing sulfide concentration for all
349 butyrate producers (**Figure 5b**). These data suggest that the levels of H₂S produced by the host
350 and constituent members of gut microbiota could shape butyrate production in the human gut
351 microbiome.

352 We next investigated the factors that contribute to strong positive complementarity
353 interactions influencing AC from EL or DF in complex communities with all butyrate producers
354 (**Figure 4f**). Butyrate concentration exhibited a strong negative correlation with lactate
355 concentration in complex communities (**Figure 5a**). Based on this correlation, we hypothesized
356 that communities with higher butyrate concentration than expected based on monoculture
357 butyrate yield (i.e. total positive butyrate complementarity) would exhibit lower lactate
358 concentration than expected based on monoculture lactate yield (i.e. total negative lactate
359 complementarity) (**Methods**). Our results demonstrated a negative correlation between butyrate
360 and lactate complementarity in communities with AC, but not in communities without AC (**Figure**
361 **5a**, inset). These results suggest that the majority of excess butyrate that was not predicted based
362 on monospecies butyrate yield was attributed to conversion of lactate to butyrate by AC,
363 suggesting that this metabolic mode for butyrate production was driving the inferred
364 complementarity effects. Thus, we next considered potential environmental factors that could
365 inhibit the conversion of lactate to butyrate.

366 Previous studies have shown that the environmental pH has a major impact on organic
367 acid production by gut microbiota^{53–56}. For example, in batch cultures of fecal inocula,
368 supplemented lactate was converted entirely to butyrate, propionate, and acetate at pH 5.9 and
369 6.4, but not at pH 5.2. This abrupt metabolic shift at low environmental pH was attributed to
370 inhibition of lactate consumption by AC and *E. hallii* (a closely related lactate-consuming butyrate
371 producer in the clostridial cluster XIVa)⁵⁵. Consistent with these results, butyrate concentration
372 and pH were positively correlated in complex communities with AC (Spearman rho=0.73,
373 p=1*10⁻⁵⁷) or without AC (Spearman rho=0.29, p=1*10⁻⁴), though the correlation was much
374 stronger in communities with AC (**Figure 5c**). A positive correlation between butyrate and pH
375 could be attributed to reduced acidification of the media on a per carbon basis because one
376 butyrate molecule is produced from (or as an alternative to) two acetate molecules (**Figure 2b**).
377 However, we postulate that in the presence of AC, a different mechanism drives the substantially
378 stronger correlation between pH and butyrate, wherein high butyrate production was enabled by
379 an environmental pH maintained above 5.9 (below which lactate conversion to butyrate by AC
380 was inhibited⁵⁵) (**Figure 5c**).

381 The abundance of EL and DF were both positively correlated with pH (**Figure S6**) and had
382 positive complementarity effects in the regression model (**Figure 4g**), consistent with the potential
383 role of pH in mediating positive complementarity effects. Further, EL had a unique environmental
384 impact by increasing the pH in monoculture, suggesting that this mechanism could contribute to
385 the inferred positive complementarity effect towards butyrate production (**Figure S6**). The
386 environmental pH for monospecies did not forecast the correlations between species abundance
387 and pH in complex communities. For example, *Dorea longicatena* (DL) and *Dorea*
388 *formicigenerans* (DF) strongly acidify the media in monoculture but are positively correlated with
389 pH in complex communities (**Figure S6**), highlighting a challenging problem in relating species

390 composition to broad functions such as environmental pH modifications across community
391 contexts.

392 In sum, we postulate that transformation of lactate into butyrate by AC was controlled by
393 a combination of pH modification and resource competition (**Figure 5d**). Based on this
394 mechanism, AC can switch between low and high butyrate producing states depending on the
395 environmental pH and availability of sugars. In communities containing pH buffering species such
396 as EL that maintain the pH above the threshold, the butyrate yield per biomass is dependent on
397 the strength of competition for limited pools of sugars (high growth, low butyrate yield state). After
398 sugars have been depleted, AC switches to a low growth and high butyrate yield metabolic state
399 that transforms lactate into butyrate. The timing of the AC metabolic switch depends on the
400 strength of resource competition in the community. In low pH environments, transformation of
401 lactate to butyrate is inhibited and thus AC competes for limited sugars, resulting in butyrate
402 production that is proportional to growth (i.e. no complementarity effects). Corroborating this
403 notion, lactate-utilizing butyrate producers, including AC, have been shown to prefer glucose over
404 lactate and produce ~5x more butyrate per unit biomass when grown on lactate versus glucose³⁹.
405 Consistent with the proposed mechanism, the abundance of AC was negatively correlated with
406 butyrate in conditions with an endpoint pH > 6 (**Figure S7**). Above this pH threshold, there exists
407 a tradeoff between the biomass of AC and butyrate production depending on the proportion of AC
408 biomass derived from sugars (i.e. high biomass, low butyrate) or lactate (i.e. low biomass, high
409 butyrate) (**Figure 5d**). In sum, the proposed mechanism indicates that in a pH buffered
410 community, resource competition over energy rich nutrients could enhance butyrate production
411 by AC by triggering a shift in metabolism from a low to high butyrate producing state. Further, this
412 hypothesis may explain why positive butyrate complementarity effects from pH-buffering species
413 were not captured by the M1 and M2 models trained on lower-order communities as there were
414 fewer species and thus a reduced strength of resource competition. This analysis highlights that
415 an interpretable statistical model that maps community composition to function can provide key
416 biological insights into ecological and molecular mechanisms driving community functions and
417 illuminates key metabolic modes of ecological drivers of community functions.

418 **DISCUSSION**

419 We demonstrated that community-level functions can be designed using a modeling framework
420 that predicts community assembly (selection effects) and then maps community composition to
421 function (complementarity effects). Our results showed that the capability for butyrate production
422 can vary over a broad range (0-20 mM or 10-60 mM butyrate in the absence and presence of AC,
423 respectively) by manipulating the presence/absence of diverse non-butyrate producing species,
424 highlighting the critical role of microbial interactions in community-level functions. We used a DTL
425 cycle to develop a predictive model of butyrate production by synthetic human gut microbiome
426 communities which enabled the identification of key microbial interactions and insights into
427 potential molecular mechanisms driving butyrate production. Our results demonstrated that
428 accurate prediction of community function in complex multi-member consortia (i.e. >10 species)
429 required measurements of communities at similar levels of complexity. Thus, the predictive
430 capability of computational models of microbial communities could be improved by choosing
431 communities that span the range of complexities of interest, rather than implementing a standard
432 procedure of characterizing pairwise communities^{6,14,42}. Consistent with this proposed
433 experimental design approach, recent theoretical work has demonstrated a similar perspective²⁶.

434 While our approach lacks a host-interaction component, the mechanistic nature of insights
435 derived from our model will enable future work to adapt our pipeline to predict community-level
436 functions in the mammalian gut environment. For instance, DP has been previously associated
437 with IBD⁵⁷, attributed to its H₂S activity inhibiting oxidation of short chain fatty acids by the host
438 via short-chain acyl-CoA dehydrogenase⁵⁸. However, an additional mechanism through which

439 hydrogen sulfide producers could contribute to IBD is by inhibiting microbial production of the anti-
440 inflammatory metabolite butyrate via the analogous bacterial short-chain acyl-CoA
441 dehydrogenase. Indeed, a previous study demonstrated that cecal contents of gnotobiotic mice
442 colonized with an 8-member community plus DP contained less propionate and elevated 3-
443 hydroxybutyrate (upstream intermediate of butyrate production) compared to the 8-member
444 community alone. In this study, the butyrate concentration did not vary between conditions, which
445 could have been masked by host butyrate consumption as the concentration was very low for
446 both communities (<1 mM)⁴⁷. This could be explained by H₂S inhibition of bacterial short chain
447 acyl-CoA dehydrogenases in butyrate and propionate metabolic pathways, observed as
448 accumulation of 3-hydroxybutyrate in the former case and decreased propionate production in
449 the latter. Additionally, this mechanistic insight could explain associations between colitis and
450 other sulfur-reducing bacteria, such as *Bilophila wadsworthia*⁵⁹, which has been shown to be
451 associated with reduced expression of microbial butyrate synthesis pathways in a mouse model
452 of colitis⁶⁰.

453 A major strategy for microbiome modulation involves administration of non-resident
454 species predicted to perform a target beneficial function⁶¹, including butyrate-producing bacteria³².
455 Due to the plasticity of microbial metabolism, our results demonstrate that it is important to
456 consider both how the resident community will enable growth of supplemented butyrate-producing
457 bacteria as well as promote the desired metabolic states. Indeed, our results showed that in the
458 presence of AC, the abundance of functional strains may not correlate with community-level
459 metabolic functions due to complementarity effects that modify microbial metabolic modes.

460 More broadly, our work provides a foundation for implementing model-guided procedures
461 to design community properties and guide development of ecological and mechanistic hypotheses
462 for a wide range of applications. Simple modifications can be made to this framework to
463 accommodate different observed system behaviors. For instance, we modeled our system using
464 two models incorporating only pairwise interaction terms. While this provided a high level of
465 interpretability, it has a limited flexibility for studying higher-order interactions, which may play a
466 critical role in shaping microbiome properties. Additionally, we focused on a predicting single
467 function, whereas designing communities for multifunctionality may be desirable in many cases.
468 Both of these limitations could be addressed by modifying our approach using alternative growth
469 and function models, leveraging advances in machine learning⁶² or integrating information from
470 genome-scale metabolic models¹³.

471 In this work, we constructed models of community dynamics and function in a single
472 media. The gut microbiome is exposed to a wide range of dietary substrates and the temporal
473 changes in resource availability can dramatically shape community composition⁶³. Our approach
474 could be adapted to represent the molecular environment as a design variable to allow
475 simultaneous exploration of the community and chemical composition design spaces to better
476 understand how the molecular environment shapes microbial community functions. In sum, our
477 methods provide a flexible foundation to explore design strategies for building microbial
478 communities with target functions from the bottom-up and to understand molecular and ecological
479 mechanisms influencing community-level functions.

480 **METHODS**

481 *Strain Maintenance and Culturing*

482 All anaerobic culturing was carried out in an anaerobic chamber with an atmosphere of $2.5 \pm 0.5\%$
483 H₂, $15 \pm 1\%$ CO₂ and balance N₂. All prepared media and materials were placed in the chamber at
484 least overnight before use to equilibrate with the chamber atmosphere. The strains used in this
485 work were obtained from the sources listed in **Table S2** and permanent stocks of each were
486 stored in 25% glycerol at -80°C . Batches of single-use glycerol stocks were produced for each

487 strain by first growing a culture from the permanent stock in anaerobic basal broth (ABB) media
488 (HiMedia or Oxoid) to stationary phase, mixing the culture in an equal volume of 50% glycerol,
489 and aliquoting 400 μL into Matrix Tubes (ThermoFisher) for storage at -80°C . Quality control for
490 each batch of single-use glycerol stocks included (1) plating a sample of the aliquoted mixture
491 onto LB media (Sigma-Aldrich) for incubation at 37°C in ambient air to detect aerobic
492 contaminants and (2) Illumina sequencing of 16S rDNA isolated from pellets of the aliquoted
493 mixture to verify the identity of the organism. For each experiment, precultures of each species
494 were prepared by thawing a single-use glycerol stock and combining the inoculation volume and
495 media listed in **Table S2** to a total volume of 5 mL (multiple tubes inoculated if more preculture
496 volume needed) for stationary incubation at 37°C for the preculture incubation time listed in **Table**
497 **S2**. All experiments were performed in a chemically defined medium (DM38), the composition of
498 which is provided in **Table S3**.

499 *Monoculture Dynamic Growth Quantification*

500 Each species' preculture was diluted to an OD_{600} of 0.0066 (Tecan F200 Plate Reader, 200 μL in
501 96-Well Microplate) in DM38 and aliquoted into 3 replicates of 1 mL each in a 96 Deep Well
502 (96DW) plate and covered with a semi-permeable membrane (Diversified Biotech) for stationary
503 incubation at 37°C . At each time point, samples were mixed and OD_{600} was measured by diluting
504 an aliquot of each sample into phosphate-buffered saline (PBS) into the linear range of the plate
505 reader.

506 *Community Culturing Experiments and Sample Collection*

507 To produce all desired community combinations, each species' preculture was diluted to an OD_{600}
508 of 0.0066 in DM38. Community combinations were arrayed in 96DW plates by pipetting equal
509 volumes of each species' diluted preculture into the appropriate wells using a Tecan Evo Liquid
510 Handling Robot inside an anaerobic chamber. Each 96DW plate was covered with a semi-
511 permeable membrane and incubated at 37°C . After 48 hours, 96DW plates were removed from
512 the incubator and samples were mixed. Cell density was measured by pipetting 200 μL of each
513 sample into one microplate and diluting 20 μL of each sample into 180 μL of PBS in another
514 microplate and measuring the OD_{600} of both plates (Tecan F200 Plate Reader). We selected the
515 value that was within the linear range of the instrument for each sample. 200 μL of each sample
516 was transferred to a new 96DW plate and pelleted by centrifugation at $2400\times g$ for 10 minutes. A
517 supernatant volume of 180 μL was removed from each sample and transferred to a 96-well
518 microplate for storage at -20°C and subsequent metabolite quantification by high performance
519 liquid chromatography (HPLC). Cell pellets were stored at -80°C for subsequent genomic DNA
520 extraction and 16S rDNA library preparation for Illumina sequencing. In some experiments, 20 μL
521 of each supernatant was used to quantify pH using a phenol Red assay⁶⁴. Phenol red solution
522 was diluted to 0.005% weight per volume in 0.9% w/v NaCl. Bacterial supernatant (20 μL) was
523 added to 180 μL of phenol red solution, and absorbance was measured at 560 nm (Tecan Spark
524 Plate Reader). A standard curve was produced by fitting the Henderson-Hasselbach equation to
525 fresh media with a pH ranging between 3 to 11 measured using a standard electro-chemical pH
526 probe (Mettler-Toledo). We used the following equation to map the pH values to the absorbance
527 measurements.

$$528 \quad \text{pH} = \text{p}K_a + b \cdot \log_{10} \left(\frac{A - A_{\min}}{A_{\max} - A} \right)$$

529
530 The parameters b and $\text{p}K_a$ were determined using a linear regression between pH and the log
531 term for the standards in the linear range of absorbance (pH between 5.2 and 11) with A_{\max}
532 representing the absorbance of the pH 11 standard, A_{\min} denoting the absorbance of the pH 3
533 standard and A representing the absorbance of each condition.

534 *Sulfide Titration Experiment*

535 Each species' preculture was diluted to an OD₆₀₀ of 0.0066 in DM38. FP cultures were
536 supplemented with 1 g/L bacto yeast extract (BD) and 33 mM sodium acetate (Sigma Aldrich).
537 Different volumes of a concentrated solution of sodium sulfide (Alfa Aesar) were added to the
538 cultures to achieve the desired concentration and the cultures were incubated in capped 1.6 mL
539 microfuge tubes for 24 hours at which point the OD₆₀₀ was measured (Tecan F200 Plate Reader,
540 200 uL in 96-Well Microplate) and supernatants were collected for organic acid quantification via
541 HPLC. Sulfide concentrations in the initial cultures were measured via the Cline assay⁶⁵ to
542 account for degradation of the sulfide stock during experimental setup. Briefly, 14.8 uL of Cline
543 reagent was added to 185.2 uL of culture supernatant and incubated in a sealed 96-Well
544 Microplate for 2 hours before diluting in 1% zinc acetate (Fisher) to the linear range of absorbance
545 measurement at 667 nm (Tecan Spark Plate Reader). A standard curve was prepared similarly
546 using sodium sulfide fixed in 1% zinc acetate. Cline reagent was prepared by dissolving 1.6 g
547 N,N-dimethyl-p-phenylenediamine sulfate (Acros Organics) and 2.4 g FeCl₃ (Fisher) in 100 mL
548 50% v/v HCl (Fisher) in water.

549 *HPLC Quantification of Organic Acids*

550 Supernatant samples were thawed in a room temperature water bath before addition of 2 μL of
551 H₂SO₄ to precipitate any components that might be incompatible with the running buffer. The
552 samples were then centrifuged at 2400xg for 10 minutes and then 150 μL of each sample was
553 filtered through a 0.2 μm filter using a vacuum manifold before transferring 70 μL of each sample
554 to an HPLC vial. HPLC analysis was performed using either a ThermoFisher (Waltham, MA)
555 Ultimate 3000 UHPLC system equipped with a UV detector (210 nm) or a Shimadzu HPLC system
556 equipped with a SPD-20AV UV detector (210 nm). Compounds were separated on a 250 x 4.6
557 mm Rezex® ROA-Organic acid LC column (Phenomenex Torrance, CA) run with a flow rate of
558 0.2 ml min⁻¹ and at a column temperature of 50°C. The samples were held at 4°C prior to injection.
559 Separation was isocratic with a mobile phase of HPLC grade water acidified with 0.015 N H₂SO₄
560 (415 μL L⁻¹). At least two standard sets were run along with each sample set. Standards were
561 100, 20, and 4 mM concentrations of butyrate, succinate, lactate, and acetate, respectively. For
562 most runs, the injection volume for both sample and standard was 25 μl. The resultant data was
563 analyzed using the Thermofisher Chromeleon 7 software package.

564 *Genomic DNA Extraction and Sequencing Library Preparation*

565 Genomic DNA was extracted from cell pellets using a modified version of the Qiagen DNeasy
566 Blood and Tissue Kit protocol. First, pellets in 96DW plates were removed from -80°C and thawed
567 in a room temperature water bath. Each pellet was resuspended by pipette in 180 μL of enzymatic
568 lysis buffer (20 mM Tris-HCl (Invitrogen), 2 mM Sodium EDTA (Sigma-Aldrich), 1.2% Triton X-
569 100 (Sigma-Aldrich), 20 mg/mL Lysozyme from chicken egg white (Sigma-Aldrich)). Plates were
570 then covered with a foil seal and incubated at 37°C for 30 minutes with orbital shaking at 600
571 RPM. Then, 25 μL of 20 mg mL⁻¹ Proteinase K (VWR) and 200 μL of Buffer AL (QIAGEN) were
572 added to each sample before mixing with a pipette. Plates were then covered by a foil seal and
573 incubated at 56°C for 30 minutes with orbital shaking at 600 RPM. Next, 200 μL of 100% ethanol
574 (Koptec) was added to each sample before mixing and samples were transferred to a Nucleic
575 Acid Binding (NAB) plate (Pall) on a vacuum manifold with a 96DW collection plate. Each well in
576 the NAB plate was then washed once with 500 uL Buffer AW1 (QIAGEN) and once with 500 μL
577 of Buffer AW2 (QIAGEN). A vacuum was applied to the Pall NAB plate for an additional 10 minutes
578 to remove any excess ethanol. Samples were then eluted into a clean 96DW plate from each well
579 using 110 μL of Buffer AE (QIAGEN) preheated to 56°C. Genomic DNA samples were stored at
580 -20°C until further processing.

581 Genomic DNA concentrations were measured using a SYBR Green fluorescence assay
582 and then normalized to a concentration of $1 \text{ ng } \mu\text{L}^{-1}$ by diluting in molecular grade water using a
583 Tecan Evo Liquid Handling Robot. First, genomic DNA samples were removed from -20°C and
584 thawed in a room temperature water bath. Then, $1 \mu\text{L}$ of each sample was combined with $95 \mu\text{L}$
585 of SYBR Green (Invitrogen) diluted by a factor of 100 in TE Buffer (Integrated DNA Technologies)
586 in a black 384-well microplate. This process was repeated with two replicates of each DNA
587 standard with concentrations of 0, 0.5, 1, 2, 4, and $6 \text{ ng } \mu\text{L}^{-1}$. Each sample was then measured
588 for fluorescence with an excitation/emission of 485/535 nm using a Tecan Spark plate reader.
589 Concentrations of each sample were calculated using the standard curve and a custom Python
590 script was used to compute the dilution factors and write a worklist for the Tecan Evo Liquid
591 Handling Robot to normalize each sample to $1 \text{ ng } \mu\text{L}^{-1}$ in molecular grade water. Samples with
592 DNA concentration less than $1 \text{ ng } \mu\text{L}^{-1}$ were not diluted. Diluted genomic DNA samples were
593 stored at -20°C until further processing.

594 Amplicon libraries were generated from diluted genomic DNA samples by PCR
595 amplification of the V3-V4 of the 16S rRNA gene using custom dual-indexed primers (**Table S3**)
596 for multiplexed next generation amplicon sequencing on Illumina platforms (Method adapted from
597 Venturelli et al. *Mol. Sys. Bio.*, 2018). Primers were arrayed in skirted 96 well PCR plates (VWR)
598 using an acoustic liquid handling robot (Labcyte Echo 550) such that each well received a different
599 combination of one forward and one reverse primer ($0.1 \mu\text{L}$ of each). After liquid evaporated, dry
600 primers were stored at -20°C . Primers were resuspended in $15 \mu\text{L}$ PCR master mix ($0.2 \mu\text{L}$
601 Phusion High Fidelity DNA Polymerase (Thermo Scientific), $0.4 \mu\text{L}$ 10 mM dNTP Solution (New
602 England Biolabs), $4 \mu\text{L}$ 5x Phusion HF Buffer (Thermo Scientific), $4 \mu\text{L}$ 5M Betaine (Sigma-
603 Aldrich), $6.4 \mu\text{L}$ Water) and $5 \mu\text{L}$ of normalized genomic DNA to give a final concentration of 0.05
604 μM of each primer. Primer plates were sealed with Microplate B seals (Bio-Rad) and PCR was
605 performed using a Bio-Rad C1000 Thermal Cycler with the following program: initial denaturation
606 at 98°C (30 s); 25 cycles of denaturation at 98°C (10 s), annealing at 60°C (30 s), extension at
607 72°C (60 s); and final extension at 72°C (10 minutes). $2 \mu\text{L}$ of PCR products from each well were
608 pooled and purified using the DNA Clean & Concentrator (Zymo) and eluted in water. The
609 resulting libraries were sequenced on an Illumina MiSeq using a MiSeq Reagent Kit v3 (600-
610 cycle) to generate 2x300 paired end reads.

611 *Bioinformatic Analysis for Quantification of Species Abundance*

612 Sequencing data were demultiplexed using Basespace Sequencing Hub's FastQ Generation
613 program. Custom python scripts were used for further data processing (Method adapted from
614 Venturelli et al. *Mol. Sys. Bio.*, 2018)⁶⁶. Paired end reads were merged using PEAR (v0.9.10)⁶⁶
615 after which reads without forward and reverse annealing regions were filtered out. A reference
616 database of the V3-V5 16S rRNA gene sequences was created using consensus sequences from
617 next-generation sequencing data or Sanger sequencing data of monospecies cultures.
618 Sequences were mapped to the reference database using the mothur (v1.40.5)⁶⁷ command
619 classify.seqs (Wang method with a bootstrap cutoff value of 60). Relative abundance was
620 calculated as the read count mapped to each species divided by the total number of reads for
621 each condition. Absolute abundance of each species was calculated by multiplying the relative
622 abundance by the OD_{600} measurement for each sample. Samples were excluded from further
623 analysis if they had $\text{OD}_{600} > 0.1$ and they had less than 1000 total reads or $> 1\%$ of the reads were
624 assigned to a species not expected to be in the community.

625 *Model-Free Quantification of Complementarity*

626 We quantified the contribution of complementarity effects to butyrate and lactate production in
627 each community by calculating the difference between the measured metabolite concentration

628 and the expected metabolite concentration based on monoculture yield according to the following
629 equation:

630

631
$$\text{Complementarity of } M_k = [M_k]_{\text{Community}} - \sum_{i \in \text{Species}} \frac{[M_k]_{\text{Monoculture } i}}{X_{i, \text{Monoculture}}} X_{i, \text{Community}}$$

632

633 The variables M_k represents metabolite k (e.g. butyrate or lactate), $[M_k]_{\text{Community}}$ represents the
634 concentration of metabolite k measured in the community, $[M_k]_{\text{Monoculture } i}$ denotes the
635 concentration of metabolite k in the monoculture of species i , $X_{i, \text{Monoculture}}$ represents the
636 absolute abundance of species i in monoculture, $X_{i, \text{Community}}$ is the absolute abundance of
637 species i in the community, and the summation is across all species in the community.

638 *gLV Models and Training*

639 We used a model with two modules: the gLV model to predict composition of the assembled
640 community and a regression model with interaction terms to predict butyrate production as a
641 function of the predicted community composition (**Figure 1a**). The gLV model is a set of N coupled
642 first-order ordinary differential equations, where N denotes the number of species, of the form:

643

644
$$\frac{1}{X_i} \frac{dX_i}{dt} = r_i + \sum_{j=1}^N a_{ij} X_j$$

645

646 The species X_i is the abundance of species i , r_i is a parameter that represents the basal growth
647 rate of species i , and a_{ij} is a parameter that represents interactions by modifying the growth rate
648 of species i proportional to the abundance of species j . To prevent unbounded growth, a_{ij} is
649 constrained to be negative when $i=j$, representing intra-species competition. This model has
650 previously been used to understand and predict the behavior of complex microbial communities⁶
651 and provides an interpretable model form (e.g. which interspecies interactions are important)
652 without introducing an excessive number of parameters (e.g. complex mechanistic models⁶⁸).

653 We used a Bayesian parameter inference approach to estimate parameters for the gLV
654 model from experimental measurements (adapted from Shin et al., *PLoS Computational Biology*,
655 2019³⁰). Briefly, our method has a prior distribution for each model parameter and then varies the
656 parameters to fit the model to the measured species abundances (mean of biological replicates)
657 while penalizing deviations from the parameter prior distributions. These penalties provide a
658 regularization effect, which is necessary when the model is underdetermined. We used L2
659 regularization because we expected inter-species competition to be prevalent and thus did not
660 expect many interaction parameters to be negligible. After an optimal parameter set is found, this
661 process is repeated hundreds of times after applying random noise to the experimental data
662 proportional to the measured experimental variance to generate an ensemble of parameter sets
663 (i.e. the posterior distribution). This posterior distribution is then used as the prior distribution when
664 updating the model with new data. We adapted a previous implementation of this method in Julia
665 for this work.

666 Before training the model on any data, we assumed a normally distributed prior for each
667 parameter with mean of 0 and standard deviation equal to 1. We then trained the gLV model on
668 time-series measurements of monoculture growth for each species, estimating a posterior
669 distribution for each r_i and a_{ij} parameter (other a_{ij} posterior distributions were equal to the prior
670 distribution). We used this posterior distribution as a prior distribution to update the model with
671 the pairwise community data and generated the gLV module of Model M1, where posterior
672 distributions were estimated from experimental data for r_i , a_{ii} , and a_{ij} where species i and species

673 j co-occurred in the experimental data and the posterior distribution of a_{ij} for unobserved pairs
674 was equal to the prior. We similarly updated the model using the 3 to 5-member community
675 experiments to generate Model M2. Regularization coefficients for each iteration of the model
676 updating process are shown in **Table S4**.

677 The gLV modules of Models M1 and M2 were underdetermined due to pairs of species
678 never being observed in the same community within the training dataset. To generate parameters
679 for these unobserved interactions, we used a matrix imputation approach to estimate the
680 interaction parameters informed by the phylogenetic relatedness of species. First, we sorted the
681 a_{ij} interaction parameter matrix such that the rows and columns occurred in the same order as the
682 phylogenetic tree (**Figure 2a**). Next, we used K-nearest neighbors matrix imputation with $K = 2$ to
683 estimate interaction parameters for species that were not observed in the training data
684 (implemented in Python 3 using the fancyimpute package, <https://pypi.org/project/fancyimpute/>).
685 This process was repeated independently for each parameter set in the posterior distribution.

686 While the parameter optimization portion of this model-training process had previously
687 been found to scale with increasing number of pairwise community datasets³⁰, we found that the
688 optimization problem became intractable when attempting to estimate parameters from complex
689 community data (i.e. >10 species). To address this problem, we used the nonlinear programming
690 solver FMINCON in MATLAB to generate the gLV module of Model M3 by training on all data
691 simultaneously. Using this method, the cost function for the optimization algorithm is computed
692 using an ODE solver to simulate each community and the sum of mean squared errors for the
693 community is computed and added to a L2 regularization term penalizing the magnitude of the
694 parameter vector. To ensure that the model did not sacrifice the goodness of fit to the time-series
695 monospecies data, the mean squared errors for these data were weighted more highly. The
696 resulting optimization function was as follows:

$$697 \varphi = \sum_{k \in \text{Single}} (X_{exp,k} - X_{model,k})^2 + w \sum_{l \in \text{Dynamic}} (X_{exp,l} - X_{model,l})^2 + \lambda \sum_{j \in \text{Params}} \theta_j^2.$$

698 In this equation, single denotes the set of experiments where only the end point community
699 composition was measured, dynamic indicates the set of time-series monospecies
700 measurements, w is the weighting factor the time-series monospecies measurements, and λ
701 represents the regularization coefficient. The FMINCON function identifies a parameter estimate
702 which minimizes the cost function. We provided the median parameter values from Model M2 as
703 an initial guess for the FMINCON function. We repeated this process with various values of λ and
704 w to find a parameter set that simultaneously fits the Dynamic and Single datasets with maximal
705 regularization penalty to prevent overfitting to the data (**Table S4**). We used a procedure based
706 on the one described above for the Julia implementation to generate an ensemble of parameter
707 sets (i.e. the posterior distribution) using FMINCON. Because each iteration of the FMINCON
708 parameter estimation took several hours to complete, we massively parallelized the generation of
709 each of the hundreds of parameter sets in the ensemble using resources from the UW-Madison
710 Center for High Throughput Computing.

713 *Regression Models and Training*

714 We used a regression model to represent a microbial community function with interaction terms:

$$715 B = \sum_{i \in \text{BPB}} \alpha_i \chi_i + \sum_{j \in \text{BPB}} \sum_{k \in \text{ALL}} \beta_{jk} \chi_j \chi_k$$

716 The variable B is the predicted butyrate concentration, α_i are parameters corresponding to each
717 of the variables χ_i (end point abundances and time 0 presence or absence (1 or 0) for each
718 butyrate producer, 10 variables total), and β_{jk} are interaction parameters corresponding to each
719 butyrate producer, 10 variables total), and β_{jk} are interaction parameters corresponding to each
720 butyrate producer, 10 variables total).

721 pair of variables χ_j (end point abundances and time 0 presence or absence (1 or 0) for each
722 butyrate producer, 10 variables total) and χ_k (end point abundances and time 0 presence or
723 absence (1 or 0) for all species, 50 variables total), excluding cases where χ_j and χ_k refer to the
724 same species (450 total parameters). Model fitting was performed using custom scripts written in
725 MATLAB and Python. We used L1 regularization to minimize the number of nonzero parameters.
726 Regularization coefficients were chosen by using 10-fold cross validation and choosing the
727 coefficient value with the lowest median mean-squared error for the test data. For models M1 and
728 M2, ensembles of regression models were generated, one for each possible combination of
729 butyrate producers, where samples containing butyrate producers from outside of each set were
730 excluded. In this case, butyrate production from less productive species (e.g. FP) were small
731 compared to more productive species (e.g. AC, ER, RI, CC) thus reducing the model accuracy
732 for communities lacking the high productivity species. For Model M3, one regression model was
733 generated using all data because all communities of interest contained highly productive butyrate
734 producers.

735 *Model Simulations to Predict New Communities*

736 Custom MATLAB scripts were used to predict community assembly and butyrate production, for
737 many communities as described in the text (e.g. all communities containing all 5 butyrate
738 producers for **Figure 4a**). For each community, the growth dynamics were simulated using each
739 parameter set from the posterior distribution of the gLV model. The resulting community
740 compositions for each simulation were an input to the regression model to predict butyrate
741 concentration. Statistics on the resulting distributions of butyrate concentration and abundance of
742 each species were stored for later plotting. Because of the large number of communities and the
743 large number of parameter sets (i.e. hundreds of simulations per community), we used parallel
744 computing (MATLAB parfor) to complete the simulations in a reasonable timeframe (~4 days for
745 the communities in **Figure 4a**).

746 **ACKNOWLEDGEMENTS**

747 We would like to thank Sungho Shin, Jordan Jalving, and Victor Zavala for their advice related to
748 implementing Julia parameter estimation methods. In addition, we are grateful to Mayank
749 Baranwal and Alfred Hero for conversations which inspired the matrix imputation approach for
750 estimating unobserved interaction parameters. We would like to thank Federico Rey for
751 generously taking the time to provide advice that improved the manuscript. Research was
752 sponsored by the National Institutes of Health and was accomplished under Grant Number
753 R35GM124774 and University of Wisconsin-Madison Office of the Chancellor and Vice
754 Chancellor for Research and Graduate Education with funding from the Wisconsin Alumni
755 Research Foundation. S.E.H. was supported by the National Institute of General Medical
756 Sciences of the National Institutes of Health under Award Number T32GM008349. R.L.C. was
757 supported in part by an NHGRI training grant to the Genomic Sciences Training Program (T32
758 HG002760). This research was performed using the computing resources and assistance of the
759 UW-Madison Center for High Throughput Computing (CHTC) in the Department of Computer
760 Sciences. The CHTC is supported by UW-Madison, the Advanced Computing Initiative, the
761 Wisconsin Alumni Research Foundation, the Wisconsin Institutes for Discovery, and the National
762 Science Foundation, and is an active member of the Open Science Grid, which is supported by
763 the National Science Foundation and the U.S. Department of Energy's Office of Science.
764

765 AUTHOR CONTRIBUTIONS

766 O.S.V. and R.L.C conceived the study. R.L.C., J.J.H., S.E.H., and B.M.C. carried out the
767 experiments. R.L.C. implemented computational modeling. R.L.C., S.E.H. and O.S.V. analyzed
768 the data. B.M.C. proposed inhibition of butyrate production by hydrogen sulfide. D.A.N. and
769 D.M.S. designed and implemented metabolite measurements. O.S.V. secured funding. R.L.C.
770 and O.S.V. wrote the paper and all authors provided feedback on the manuscript.

771 CONFLICT OF INTEREST

772 The authors do not have a conflict of interest.

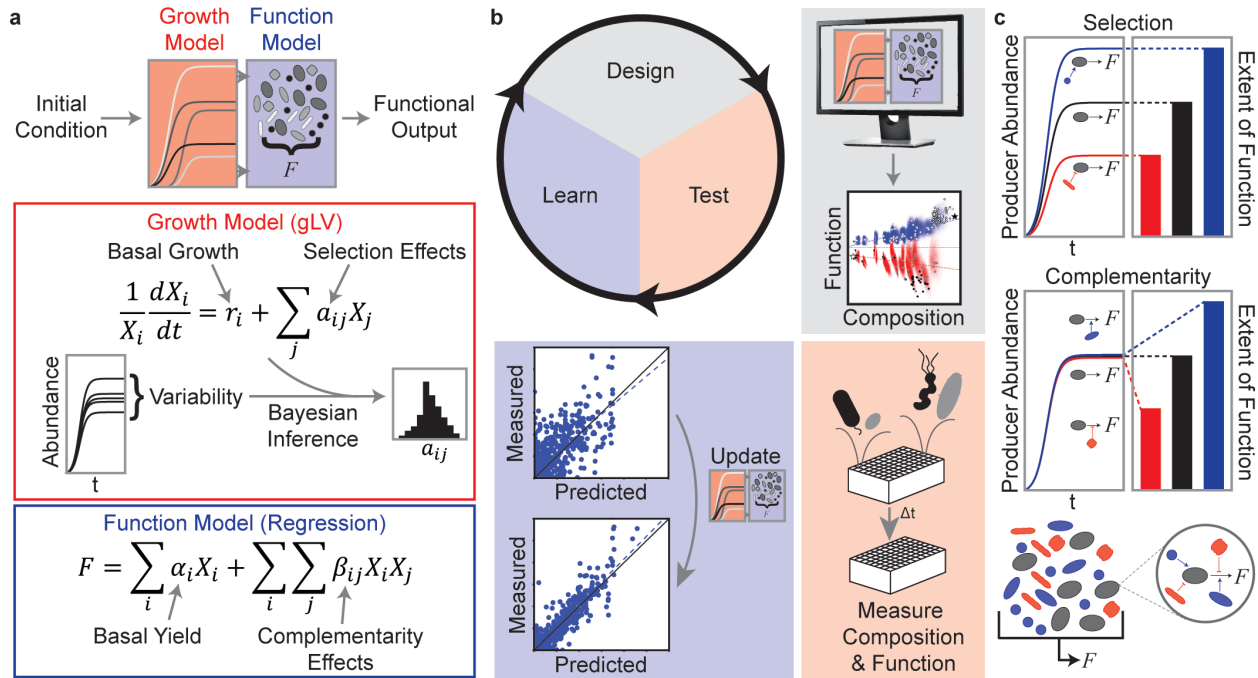
773 REFERENCES

- 774 1. Thompson, L. R. *et al.* A communal catalogue reveals Earth’s multiscale microbial
775 diversity. *Nature* **551**, 457–463 (2017).
- 776 2. Moran, M. A. The global ocean microbiome. *Science* (80-.). **350**, (2015).
- 777 3. Ahmad, F., Ahmad, I. & Khan, M. S. Screening of free-living rhizospheric bacteria for their
778 multiple plant growth promoting activities. *Microbiol. Res.* **163**, 173–181 (2008).
- 779 4. Flint, H. J., Scott, K. P., Duncan, S. H., Louis, P. & Forano, E. Microbial degradation of
780 complex carbohydrates in the gut. *Gut Microbes* 289–306 (2012).
- 781 5. Lawson, C. E. *et al.* Common principles and best practices for engineering microbiomes.
782 *Nat. Rev. Microbiology* (2019). doi:10.1038/s41579-019-0255-9
- 783 6. Venturelli, O. S. *et al.* Deciphering microbial interactions in synthetic human gut
784 microbiome communities. *Mol. Syst. Biol.* **14**, e8157 (2018).
- 785 7. Friedman, J., Higgins, L. M. & Gore, J. Community structure follows simple assembly
786 rules in microbial microcosms. *Nat. Ecol. Evol.* **1**, 0109 (2017).
- 787 8. Ratzke, C. & Gore, J. Modifying and reacting to the environmental pH drives bacterial
788 interactions. *PLoS Biol.* 136838 (2018). doi:10.1101/136838
- 789 9. Abreu, C. I., Friedman, J., Andersen Woltz, V. L. & Gore, J. Mortality causes universal
790 changes in microbial community composition. *Nat. Commun.* **10**, (2019).
- 791 10. Ratzke, C., Barrere, J. & Gore, J. Strength of species interactions determines biodiversity
792 and stability in microbial communities. *Nat. Ecol. Evol.* (2020).
- 793 11. Piccardi, P., Vessman, B. & Mitri, S. Toxicity drives facilitation between 4 bacterial
794 species. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 15979–15984 (2019).
- 795 12. Lax, S., Abreu, C. I. & Gore, J. Higher temperatures generically favour slower-growing
796 bacterial species in multispecies communities. *Nat. Ecol. Evol.* (2020).
- 797 13. Bauer, E., Zimmermann, J., Baldini, F., Thiele, I. & Kaleta, C. BacArena: Individual-based
798 metabolic modeling of heterogeneous microbes in complex communities. *PLoS Comput.*
799 *Biol.* **13**, 1–22 (2017).
- 800 14. Sanchez-Gorostiaga, A., Bajić, D., Osborne, M. L., Poyatos, J. F. & Sanchez, A. High-
801 order interactions distort the functional landscape of microbial consortia. *PLoS Biol.*
802 (2019). doi:10.1101/333534
- 803 15. Medlock, G. L. *et al.* Inferring Metabolic Mechanisms of Interaction within a Defined Gut
804 Microbiota. *Cell Syst.* **7**, 245-257.e7 (2018).
- 805 16. Stein, R. R. *et al.* Computer-guided design of optimal microbial consortia for immune
806 system modulation. *Elife* **7**, 1–17 (2018).
- 807 17. Rosser, E. C. *et al.* Microbiota-Derived Metabolites Suppress Arthritis by Amplifying Aryl-
808 Hydrocarbon Receptor Activation in Regulatory B Cells. *Cell Metab.* **31**, 1–15 (2020).
- 809 18. Kimura, I. *et al.* The gut microbiota suppresses insulin-mediated fat accumulation via the
810 short-chain fatty acid receptor GPR43. *Nat. Commun.* **4**, 1829 (2013).

- 811 19. Li, Z. *et al.* Butyrate reduces appetite and activates brown adipose tissue via the gut-brain
812 neural circuit. *Gut* Published Online First: 03 November 2017 (2017). doi:10.1136/gutjnl-
813 2017-314050
- 814 20. Lin, H. V. *et al.* Butyrate and propionate protect against diet-induced obesity and regulate
815 gut hormones via free fatty acid receptor 3-independent mechanisms. *PLoS One* **7**, 1–9
816 (2012).
- 817 21. Furusawa, Y. *et al.* Commensal microbe-derived butyrate induces the differentiation of
818 colonic regulatory T cells. *Nature* **504**, 446–450 (2013).
- 819 22. Segain, J. P. *et al.* Butyrate inhibits inflammatory responses through NFkB inhibition:
820 implications for Crohn’s disease. *Gut* **47**, 397–403 (2000).
- 821 23. Rivera-Chávez, F. *et al.* Depletion of Butyrate-Producing Clostridia from the Gut
822 Microbiota Drives an Aerobic Luminal Expansion of Salmonella. *Cell Host Microbe* **19**,
823 443–454 (2016).
- 824 24. Khan, S., Maremanda, K. P. & Jena, G. Butyrate, a Short-Chain Fatty Acid and Histone
825 Deacetylases Inhibitor: Nutritional, Physiological, and Pharmacological Aspects in
826 Diabetes. in *Handbook of Nutrition, Diet, and Epigenetics* 1–15 (Springer International
827 Publishing, 2017). doi:10.1007/978-3-319-31143-2_70-1
- 828 25. Fung, K. Y. C., Cosgrove, L., Lockett, T., Head, R. & Topping, D. L. A review of the
829 potential mechanisms for the lowering of colorectal oncogenesis by butyrate. *Br. J. Nutr.*
830 **108**, 820–831 (2012).
- 831 26. Maynard, D. S., Miller, Z. R. & Allesina, S. Predicting coexistence in experimental
832 ecological communities. *Nat. Ecol. Evol.* **4**, 91–100 (2020).
- 833 27. Hector, A. & Loreau, M. Partitioning selection and complementarity in biodiversity
834 experiments. *Nature* **412**, 72–76 (2001).
- 835 28. Bell, T., Newman, J. A., Silverman, B. W., Turner, S. L. & Lilley, A. K. The contribution of
836 species richness and composition to bacterial services. *Nature* **436**, 1157–1160 (2005).
- 837 29. Rivett, D. W. & Bell, T. Abundance determines the functional role of bacterial phylotypes
838 in complex communities. *Nat. Microbiol.* **3**, 767–772 (2018).
- 839 30. Shin, S., Venturelli, O. S. & Zavala, V. M. Scalable nonlinear programming framework for
840 parameter estimation in dynamic biological system models. *PLoS Comput. Biol.* 1–29
841 (2019).
- 842 31. Petrof, E. O. & Khoruts, A. From stool transplants to next-generation microbiota
843 therapeutics. *Gastroenterology* **146**, 1573–1582 (2014).
- 844 32. Glijamse, P. W. *et al.* Treatment with *Anaerobutyricum soehngenii*: a pilot study of safety
845 and dose–response effects on glucose metabolism in human subjects with metabolic
846 syndrome. *npj Biofilms Microbiomes* 1–10 (2020). doi:10.1038/s41522-020-0127-0
- 847 33. Vital, M., Howe, A. & Tiedje, J. Revealing the Bacterial Synthesis Pathways by Analyzing
848 (Meta) Genomic Data. *MBio* **5**, 1–11 (2014).
- 849 34. Lawson, C. E. & Lückner, S. Complete ammonia oxidation: an important control on
850 nitrification in engineered ecosystems? *Curr. Opin. Biotechnol.* **50**, 158–165 (2018).
- 851 35. Natividad, J. M. *et al.* Impaired Aryl Hydrocarbon Receptor Ligand Production by the Gut
852 Microbiota Is a Key Factor in Metabolic Syndrome. *Cell Metab.* **28**, 737-749.e4 (2018).
- 853 36. Arif, I., Batool, M. & Schenk, P. M. Plant Microbiome Engineering: Expected Benefits for
854 Improved Crop Growth and Resilience. *Trends Biotechnol.* 1–12 (2020).
855 doi:10.1016/j.tibtech.2020.04.015
- 856 37. Forster, S. C. *et al.* A human gut bacterial genome and culture collection for precise and
857 efficient metagenomic analysis. *Nat. Biotechnol.* **37**, (2019).
- 858 38. Louis, P. & Flint, H. J. Diversity, metabolism and microbial ecology of butyrate-producing
859 bacteria from the human large intestine. *FEMS Microbiol. Lett.* **294**, 1–8 (2009).
- 860 39. Duncan, S. H., Louis, P. & Flint, H. J. Lactate-Utilizing Bacteria, Isolated from Human
861 Feces, That Produce Butyrate as a Major Fermentation Product. *Appl. Environ. Microbiol.*

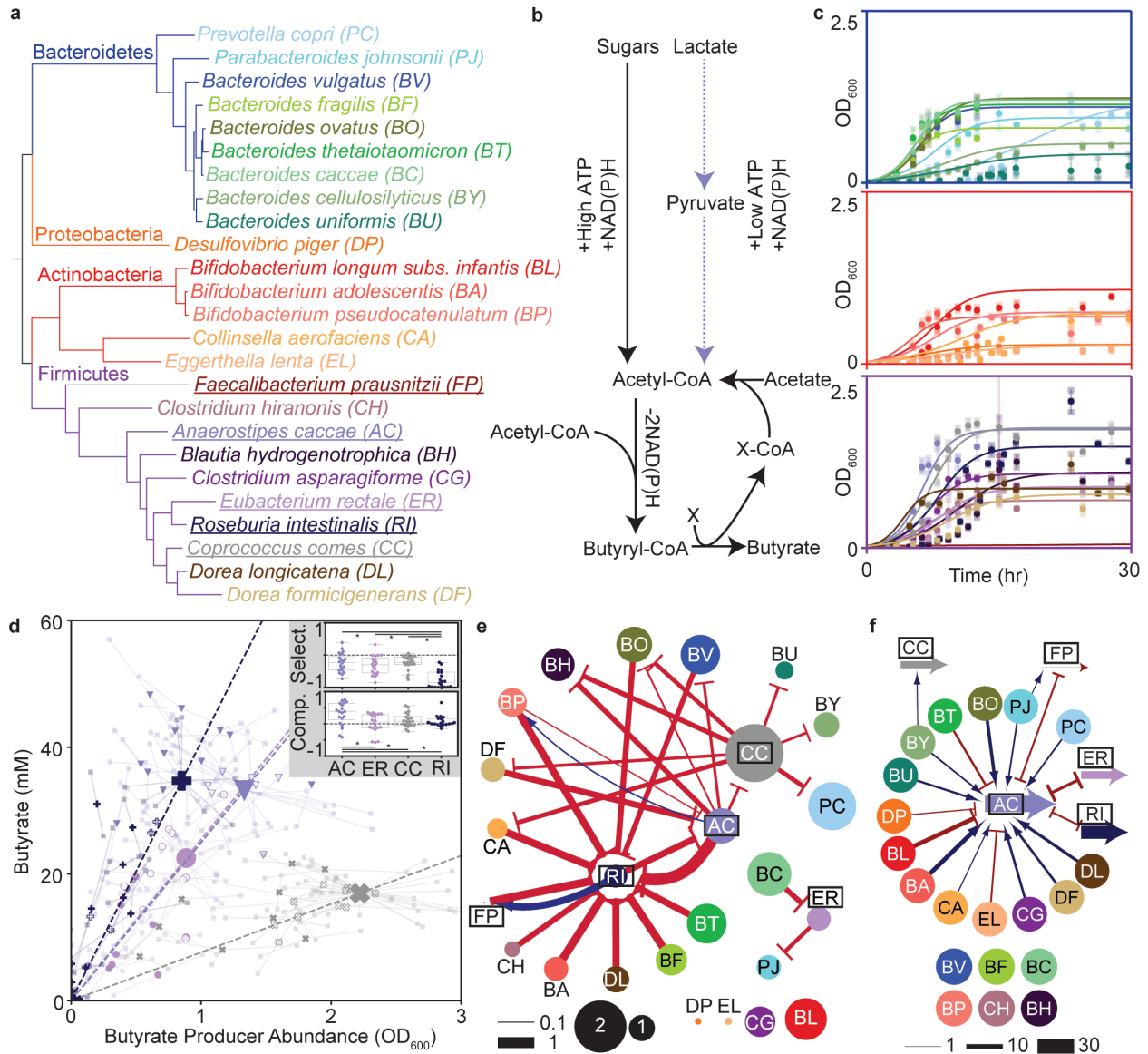
- 862 **70**, 5810–5817 (2004).
- 863 40. Tramontano, M., Andrejev, S., Pruteanu, M. & Klünemann, M. Nutritional preferences of
864 human gut bacteria reveal their metabolic idiosyncrasies. *Nat. Microbiol.* 4–6 (2018).
865 doi:10.1038/s41564-018-0123-9
- 866 41. Fort, H. Making quantitative predictions on the yield of a species immersed in a
867 multispecies community: The focal species method. *Ecol. Modell.* **430**, 109108 (2020).
- 868 42. Kong, W., Meldgin, D. R., Collins, J. J. & Lu, T. Designing microbial consortia with
869 defined social interactions. *Nat. Chem. Biol.* **14**, 821–829 (2018).
- 870 43. Goldford, J. E. *et al.* Emergent simplicity in microbial community assembly. *Science (80-*
871 *)*. **361**, 469–474 (2018).
- 872 44. Zmora, N. *et al.* Personalized Gut Mucosal Colonization Resistance to Empiric Probiotics
873 Is Associated with Unique Host and Microbiome Features. *Cell* **174**, 1388-1405.e21
874 (2018).
- 875 45. Shade, A. *et al.* Fundamentals of microbial community resistance and resilience. *Front.*
876 *Microbiol.* **3**, 1–19 (2012).
- 877 46. Bittleston, L. S., Gralka, M., Leventhal, G. E., Mizrahi, I. & Cordero, O. X. Context-
878 dependent dynamics lead to the assembly of functionally distinct microbial communities.
879 *Nat. Commun.* **11**, 1440 (2020).
- 880 47. Rey, F. E. *et al.* Metabolic niche of a prominent sulfate-reducing human gut bacterium.
881 *Proc. Natl. Acad. Sci. U. S. A.* **110**, 13582–7 (2013).
- 882 48. Flint, H. J., Duncan, S. H., Scott, K. P. & Louis, P. Interactions and competition within the
883 microbial community of the human colon: Links between diet and health. *Environ.*
884 *Microbiol.* **9**, 1101–1111 (2007).
- 885 49. Vital, M., Karch, A. & Pieper, D. H. Colonic Butyrate-Producing Communities in Humans:
886 an Overview Using Omics Data. *mSystems* **2**, 1–18 (2017).
- 887 50. Foster, K. R. & Bell, T. Competition, not cooperation, dominates interactions among
888 culturable microbial species. *Curr. Biol.* **22**, 1845–1850 (2012).
- 889 51. Marquet, P., Duncan, S. H., Chassard, C., Bernalier-Donadille, A. & Flint, H. J. Lactate
890 has the potential to promote hydrogen sulphide formation in the human colon. *FEMS*
891 *Microbiol. Lett.* **299**, 128–134 (2009).
- 892 52. Struchtemeyer, C. G., Duncan, K. E. & Mcinerney, M. J. Evidence for syntrophic butyrate
893 metabolism under sulfate-reducing conditions in a hydrocarbon-contaminated aquifer.
894 *FEMS Microbiol. Ecol.* **76**, 289–300 (2011).
- 895 53. Ilhan, Z. E., Marcus, A. K., Kang, D.-W., Rittmann, B. E. & Krajmalnik-Brown, R. pH-
896 Mediated Microbial and Metabolic Interactions in Fecal Enrichment Cultures. *mSphere* **2**,
897 1–12 (2017).
- 898 54. Walker, A. W., Duncan, S. H., Carol McWilliam Leitch, E., Child, M. W. & Flint, H. J. pH
899 and peptide supply can radically alter bacterial populations and short-chain fatty acid
900 ratios within microbial communities from the human colon. *Appl. Environ. Microbiol.* **71**,
901 3692–3700 (2005).
- 902 55. Belenguer, A. *et al.* Impact of pH on lactate formation and utilization by human fecal
903 microbial communities. *Appl. Environ. Microbiol.* **73**, 6526–6533 (2007).
- 904 56. Reichardt, N. *et al.* Specific substrate-driven changes in human faecal microbiota
905 composition contrast with functional redundancy in short-chain fatty acid production.
906 *ISME J.* **12**, 610–622 (2018).
- 907 57. Loubinoux, J., Bronowicki, J. P., Pereira, I. A. C., Mougénel, J. L. & Le Faou, A. E.
908 Sulfate-reducing bacteria in human feces and their association with inflammatory bowel
909 diseases. *FEMS Microbiol. Ecol.* **40**, 107–112 (2002).
- 910 58. Babidge, W., Millard, S. & Roediger, W. Sulfides impair short chain fatty acid β -oxidation
911 at acyl-CoA dehydrogenase level in colonocytes: Implications for ulcerative colitis. *Mol.*
912 *Cell. Biochem.* **181**, 117–124 (1998).

- 913 59. Devkota, S. *et al.* Dietary-fat-induced taurocholic acid promotes pathobiont expansion
914 and colitis in *Il10^{-/-}* mice. *Nature* **487**, 104–108 (2012).
- 915 60. Natividad, J. M. *et al.* *Bilophila wadsworthia* aggravates high fat diet induced metabolic
916 dysfunctions in mice. *Nat. Commun.* **9**, 2802 (2018).
- 917 61. Veiga, P., Suez, J., Derrien, M. & Elinav, E. Moving from probiotics to precision
918 probiotics. *Nat. Microbiol.* (2020). doi:10.1038/s41564-020-0721-1
- 919 62. Costello, Z. & Martin, H. G. A machine learning approach to predict metabolic pathway
920 dynamics from time-series multiomics data. *npj Syst. Biol. Appl.* **4**, 1–14 (2018).
- 921 63. Patnode, M. L. *et al.* Interspecies Competition Impacts Targeted Manipulation of Human
922 Gut Bacteria by Fiber-Derived Glycans. *Cell* **179**, 59-73.e13 (2019).
- 923 64. Silverstein, T. P. Fitting imidazole 1H NMR titration data to the Henderson-Hasselbalch
924 equation. *J. Chem. Educ.* **89**, 1474–1475 (2012).
- 925 65. Cline, J. D. Spectrophotometric Determination of Hydrogen Sulfide in Natural Waters.
926 *Limnol. Oceanogr.* **14**, 454–458 (1969).
- 927 66. Zhang, J., Kobert, K., Flouri, T. & Stamatakis, A. PEAR: A fast and accurate Illumina
928 Paired-End reAd mergeR. *Bioinformatics* **30**, 614–620 (2014).
- 929 67. Schloss, P. D. *et al.* Introducing mothur: Open-source, platform-independent, community-
930 supported software for describing and comparing microbial communities. *Appl. Environ.*
931 *Microbiol.* **75**, 7537–7541 (2009).
- 932 68. Momeni, B., Xie, L. & Shou, W. Lotka-Volterra pairwise modeling fails to capture diverse
933 pairwise microbial interactions. *Elife* **6**, e25051 (2017).
- 934
- 935



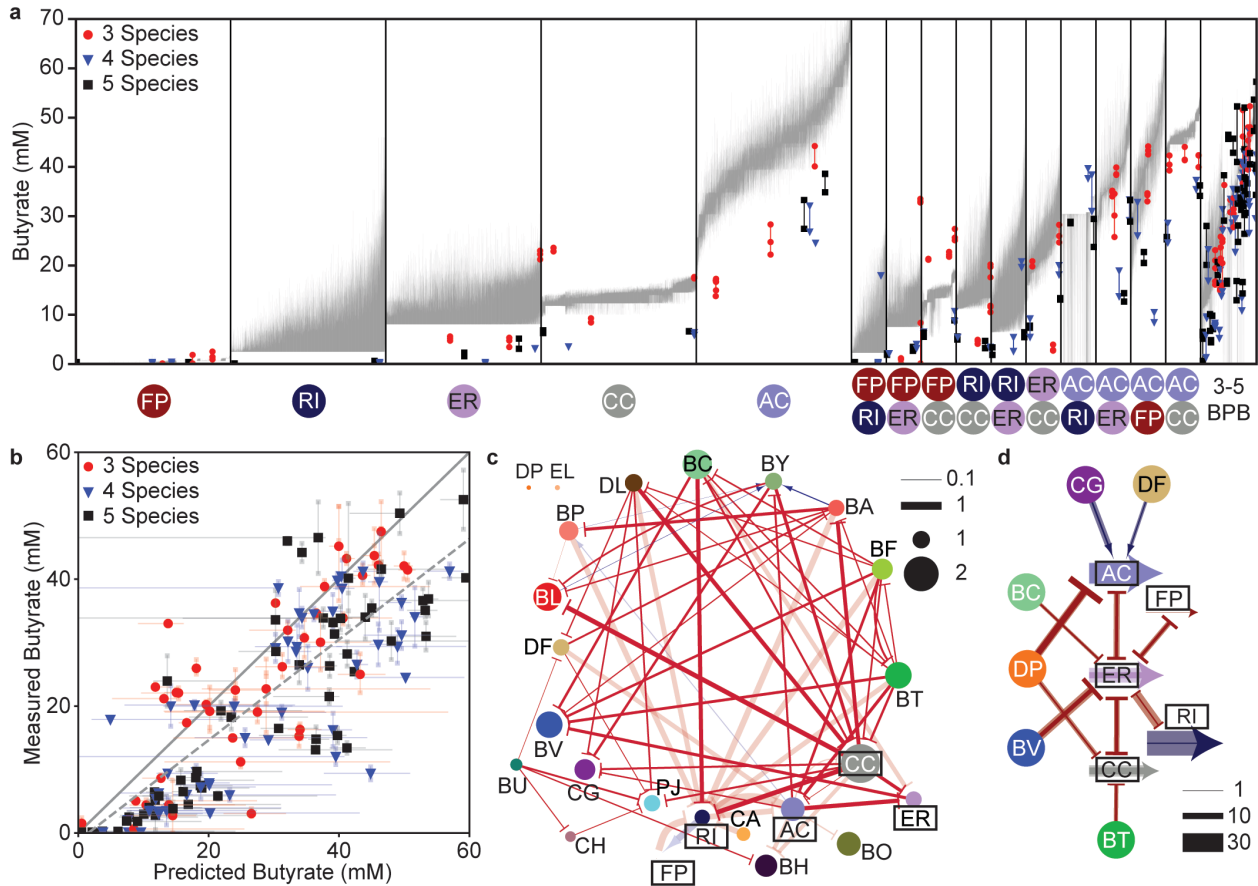
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954

Figure 1. Iterative modeling framework to predict microbial community assembly and function. (a) Two-stage modeling framework for predicting community assembly and function. The generalized Lotka-Volterra model (gLV) represents community dynamics. The inter-species interaction terms represent the selection effects described in c. A Bayesian Inference approach was used to determine parameter uncertainties due to biological and technical variability. A linear regression model with interactions represents the complementarity effects as described in c. Combining these two models enables prediction of a probability distribution of the functional activity from initial species concentrations. (b) Model-guided iterative experimental approach for developing a model to predict community assembly and butyrate production. First, we use our model to explore the design space of possible experiments (i.e. different initial conditions of species presence/absence) and design communities that span the range of expected functional outputs. Next, we use high-throughput experimental methods to measure species abundance and functional outputs. Finally, we evaluate the model's capability to accurately predict the experimental data and train the model on new data for the next iteration. (c) Inter-species interactions that impact the functional output of an organism can be driven by selection (top) or complementarity (bottom) effects. In this model, the total functional output of the communities is determined by a combination of these effects.



955
 956 **Figure 2. Characterizing interaction types in two-species communities.** (a) Phylogenetic tree
 957 of the synthetic human gut microbiome community composed of 25 highly prevalent and diverse
 958 species. Branch color indicates phylum and underlined species denote butyrate producers. (b)
 959 Metabolic pathways for the transformation of sugars, acetate, and lactate into butyrate.
 960 Conversion of sugars or lactate to acetyl-CoA generates ATP and NAD(P)H, with higher ATP
 961 production per NAD(P)H from sugars. NAD(P)H is oxidized through conversion of acetyl-CoA to
 962 butyryl-CoA. Many substrates (X) can be used to exchange CoA between acetate and/or butyrate.
 963 In our system, *Anaerostipes caccae* has the unique capability to utilize the lactate conversion
 964 pathway (purple dashed arrows). (c) Monospecies growth responses over time. Transparent
 965 symbols indicate biological replicates connected to the corresponding mean (solid symbols) by
 966 transparent lines. Solid lines represent the generalized Lotka-Volterra (gLV) model fit to the data.
 967 Each plot shows the growth curves for species within the Bacteroidetes (top),
 968 Actinobacteria/Proteobacteria (middle) or Firmicutes (bottom) phylum. (d) Scatter plot of butyrate
 969 producer abundance and butyrate concentration for all pairwise communities containing at least
 970 one butyrate producer. Solid symbols indicate the mean of biological replicates of a community.
 971 Large symbols indicate butyrate producer monoculture. Smaller symbols indicate two-species

972 communities, with closed symbols denoting significant differences in butyrate concentration
973 and/or butyrate-producer abundance from the monoculture ($p < 0.05$, t-test, unequal variance).
974 Transparent squares indicate biological replicates and are connected to the corresponding mean
975 with lines. Dashed lines indicate the predicted butyrate concentration assuming a constant
976 butyrate yield based on monoculture data. Inset: distribution of selection and complementarity
977 effects normalized by monoculture butyrate concentration for two-species communities. Asterisks
978 indicate significant difference in the mean across butyrate producers ($p < 0.05$, t-test, unequal
979 variance) **(e)** Network representation of the inferred gLV inter-species interaction network based
980 on data from **b** and **c**. Nodes size represents the abundance of each species in monoculture
981 (OD_{600}) at 48 hr and edges indicate interaction parameters with widths proportional to magnitude
982 (units of $hr^{-1} OD_{600}^{-1}$) and color indicating sign (red negative, blue positive). Only edges with $>95\%$
983 confidence in sign are shown. **(f)** Network representation of regression model trained on data
984 from **b** and **c**. Butyrate producer arrows denote monoculture butyrate production, nodes indicate
985 non-butyrate producers, and edges represent modification of butyrate production in two-species
986 communities. Edges connecting two butyrate producer arrows appear as bidirectional arrows
987 since the directionality of the effect cannot be inferred. Edge widths are proportional to butyrate
988 production (units of mM Butyrate). Only interactions with magnitude greater than 2 mM are shown.
989



990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

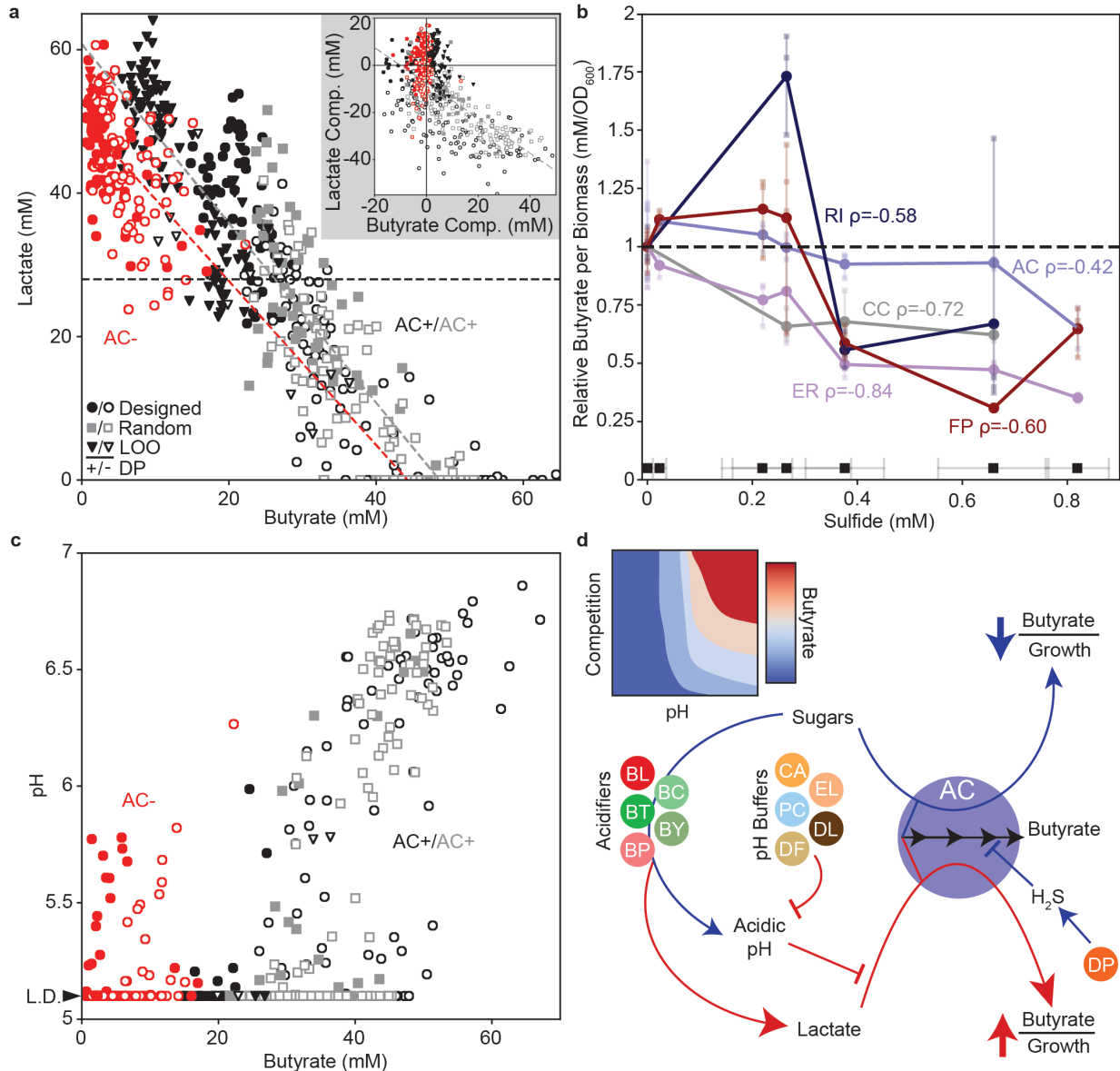
1014

Figure 3. Model-guided investigation of low complexity synthetic human gut communities.

(a) Predicted (grey bars) and measured (data points) butyrate concentrations for all 3-5 member communities containing at least one butyrate producer. Vertical black lines separate groups of communities based on the identities of the combination of butyrate producers specified on the x-axis. Communities with all combinations of 3-5 butyrate producers are included in the final bin for simplicity. Within each bin, communities are sorted in rank order of increasing median predicted butyrate production with a vertical grey bar indicating the 60 percent confidence interval of predictions for each community. Data points represent biological replicates of a selected subset of communities with replicates of a community connected by lines and the data point type indicating the number of species in each community (156 communities total). **(b)** Scatter plot of predicted and measured butyrate concentrations for communities in **a**. Transparent datapoints represent biological replicates of a community and are connected to the corresponding mean measurement (solid datapoints) by transparent lines. Prediction error bars indicate the 60 percent confidence interval of predicted butyrate. Solid grey line indicates $x=y$ and dashed line indicates the linear regression between the mean measurement and median prediction ($y=0.79x-1.2$, Pearson $r=0.83$, $p=6 \times 10^{-41}$). Data point type indicates the number of species in each community. **(c)** Network representation of generalized Lotka-Volterra model updated with data from **b**. Node size represents the abundance of each species in monoculture (OD_{600}), edge widths denote the magnitude of the inter-species interaction coefficients (units of $hr^{-1} OD_{600}^{-1}$) and color of the edges corresponds to the sign (red negative, blue positive). Faint interaction edges indicate interactions that did not change from the model trained on monospecies and pairwise communities (<2-fold change in magnitude of parameter mean). Only edges with >95% confidence in sign are shown. **(d)** Network representation of contributions of updated regression model to butyrate production in communities from **b**. Butyrate producer arrows indicate contribution to butyrate production

1015 independent of interactions. Nodes indicate non-butyrate producing species, and edges indicate
1016 modification of butyrate production in communities from **b**. Edges connecting two butyrate
1017 producers are bidirectional because it is not possible to discern which organism is producing the
1018 butyrate. For each butyrate producer, solid edge widths are proportional to the mean contribution
1019 and faint edge widths are proportional to the maximum contribution of the interaction across
1020 communities where those species were present (units of mM butyrate). Only interactions with
1021 maximum contribution >5 mM and with at least 4 communities including that interaction are
1022 shown.
1023

1025 **Figure 4. Model-guided exploration of butyrate production landscape. (a)** Scatter plot of
1026 Euclidean distance in community absolute abundance from predicted full 25-member community
1027 versus predicted butyrate concentration for all possible communities. Histograms indicate the
1028 distribution of communities across the given axis. Communities are colored according to the
1029 presence (red) or absence (blue) of *Desulfovibrio piger* (DP). Blue and red dashed lines indicate
1030 the linear regression of communities with (red, $y=-1.7x+25.5$, $r=-0.26$) or without (blue,
1031 $y=3.1x+27.8$, $r=0.72$) DP. The white star indicates the full 25-member community and black star
1032 indicates the community of all butyrate producers. Large data points indicate communities chosen
1033 for experimental validation. Black triangles indicate leave-one-out communities, black circles
1034 indicate designed communities, and grey squares indicate random communities, with
1035 open/closed symbols indicating absence/presence of DP. **(b)** Scatter plot of Euclidean distance
1036 in community composition from predicted 24-member community excluding *Anaerostipes caccae*
1037 (AC) versus predicted butyrate concentration for all possible communities. Histograms indicate
1038 the distribution of communities across the given axis. Grey dashed line indicates the mean
1039 predicted butyrate concentration across all communities. Blue dashed line indicates the linear
1040 regression of all communities ($y=8.3x-1.4$, $r=0.50$). The white star indicates the full 24-member
1041 community and the black star indicates the 4 butyrate-producer community. Large data points
1042 indicate communities chosen for experimental validation. Inset: mean experimental
1043 measurements of butyrate concentration (black) and total abundance of butyrate producers (red)
1044 versus the distance from the full 24-member community. The grey and red dashed lines represent
1045 the mean butyrate concentration and total butyrate producer abundances across measured
1046 communities, respectively. **(c)** Scatter plot of predicted versus measured butyrate concentration
1047 for communities in **a**. Transparent data points indicate biological replicates and are connected to
1048 the corresponding mean values by transparent lines. Data points denote the median with error
1049 bars spanning the 60% confidence interval. Solid line indicates $x=y$. Dashed line indicates linear
1050 regression of median prediction versus mean measurement ($y=1.2x+3.0$, $r=0.59$, $p=1*10^{-18}$).
1051 Legend indicates statistically significant differences in measured butyrate between populations of
1052 communities (Kruskal-Wallis test). **(d)** Scatter plot of predicted versus measured butyrate for
1053 communities in **b**. Transparent data points indicate biological replicates and are connected to the
1054 corresponding mean by transparent lines. Data points represent the median with error bars
1055 spanning the 60% confidence interval. Solid grey line indicates $x=y$. Dashed line indicates the
1056 linear regression of the median versus mean measurement ($y=0.1x+2.6$, $r=0.44$, $p=2*10^{-5}$). **(e)**
1057 Scatter plot of predicted versus measured butyrate for complex communities using model M3.
1058 Transparent data points indicate biological replicates and are connected to the corresponding
1059 mean by transparent lines. Solid grey line indicates $x=y$. Dashed line indicates linear regression
1060 of prediction versus mean measurement ($y=0.99x-4.6$, $r=0.86$, $p=1*10^{-44}$). **(f)** Heat-map of the
1061 median value of the inter-species interaction coefficients (a_{ij}) for the M3 gLV model. Interactions
1062 impacting AC and DP are annotated with L+ or L- if species j produced or consumed >10 mM
1063 lactate in monoculture. Inter-species interactions included in the model community C1 from
1064 (Venturelli et al., Mol. Sys. Bio., 2018)⁶ are annotated with C+ or C- if interactions from both
1065 models had magnitudes greater than 0.05 hr^{-1} and had the same or opposite sign, respectively.
1066 **(g)** Network representation of updated M3 regression model. Butyrate producer arrows indicate
1067 contribution to butyrate production independent of inter-species interactions. Nodes indicate non-
1068 butyrate producing species, and edges indicate modification of butyrate production (blue,
1069 increased; red, decreased) in all complex communities (>10 species). Solid edge widths are
1070 proportional to the mean contribution and faint edge widths are proportional to the maximum
1071 contribution of the interaction across communities where those species were present (units of mM
1072 butyrate). Only interactions with maximum contribution >5 mM and with at least 4 communities
1073 including that interaction are shown.
1074
1075



1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091

Figure 5. Model-guided identification of molecular mechanisms impacting butyrate production. (a) Scatter plot of butyrate concentration versus lactate concentration for complex communities (>10 species). Each data point indicates a biological replicate of a community. Grey dashed line indicates the linear regression for communities containing AC ($y = -1.3x + 61$, $r = -0.91$, $p = 5 \times 10^{-182}$), red dashed line indicates the linear regression for communities lacking AC ($y = -1.1x + 51$, $r = -0.56$, $p = 8 \times 10^{-16}$) and black horizontal dashed line indicates initial concentration of lactate in the media (28 mM). Inset: butyrate complementarity versus lactate complementarity. Grey dashed line indicates the linear regression for communities containing AC ($y = -0.75x - 7.5$, $r = -0.63$, $p = 4 \times 10^{-53}$). Pearson correlation for communities lacking AC was not statistically significant ($p = 0.12$). (b) Butyrate concentration per unit biomass as a function of sulfide concentration. Butyrate yield per biomass was normalized to the no sulfide condition. Circles indicate the mean of biological replicates, with individual replicates shown as transparent squares. Black squares indicate the mean measured sulfide concentration for each treatment level with error bars indicating the standard deviation of at least 10 technical replicates. Species labels are accompanied by statistically significant Spearman correlation coefficients (ρ) between all

1092 biological replicates of that species and mean sulfide concentration for each level ($p < 0.05$, AC
1093 $p = 0.02$; CC $p = 0.002$; ER $p = 3 \times 10^{-8}$; RI $p = 0.02$; FP $p = 0.0008$). **(c)** Scatter plot of butyrate
1094 concentration versus pH for complex communities. Each data point indicates a biological replicate
1095 of a community. **(d)** Schematic representing proposed driving mechanisms impacting butyrate
1096 production by AC in complex communities. Red edges denote processes that negatively impact
1097 butyrate production and blue edges represent processes that enhance butyrate production. The
1098 abundance of species that acidify the environment were positively correlated with lactate
1099 concentration and negatively correlated with pH in complex communities. The abundance of pH
1100 buffering species were positively correlated with pH in complex communities. Note that species
1101 contributions to these processes are expected to be context-dependent. Inset: proposed
1102 qualitative butyrate landscape as a function of the strength of resource competition for sugars and
1103 environmental pH.