

1 **Title:** Boolean Implication Analysis Unveils Candidate Universal Relationships in Microbiome
2 Data

3

4 **Authors:** Daniella Vo^{1†}; Shayal Charisma Singh^{2†}; Sara Safa^{3†}; Debashis Sahoo^{3,4,5}

5

6

7 **Affiliations:**

8 ¹Department of Bioinformatics and Systems Biology, University of California San Diego, La Jolla, CA
9 92093-083.

10 ²Hacıoğlu Data Science Institute, University of California San Diego, La Jolla, CA 92093-083.

11 ³Department of Computer Science & Engineering, Jacob's School of Engineering, University of
12 California San Diego, La Jolla, CA 92093-083.

13 ⁴Department of Pediatrics, University of California San Diego, La Jolla, CA 92093-083.

14 ⁵Moore's Cancer Center, University of California San Diego, La Jolla, CA 92093-083.

15

16 † Equal contribution

17

18

19

20

21

22 **Corresponding author:**

23 **Debashis Sahoo, Ph.D.;** Assistant Professor, Department of Pediatrics, University of California San
24 Diego; 9500 Gilman Drive, MC 0730, Leichtag Building 132; La Jolla, CA 92093-0831. **Phone:** 858-
25 246-1803; **Fax:** 858-246-0019; **Email:** dsahoo@ucsd.edu

26 **Abstract**

27 **Background:** Microbiomes consist of bacteria, viruses, and other microorganisms, and are
28 responsible for many different functions in both organisms and the environment. Some previous
29 analyses of microbiomes focus on the relationships between specific microbiomes and a particular
30 disease. These typically use correlation which is fundamentally symmetric with respect to pairs of
31 microbes. Correlation focuses on the symmetry of the data distribution, and asymmetric data is
32 often discarded as having a weak correlation. With all the data available on the microbiome, there
33 is a need for a method that comprehensively studies microbiomes and how they are related to each
34 other.

35
36 **Results:** We collect publicly available datasets from human, environment, and animal samples to
37 determine both symmetric and asymmetric Boolean relationships between a pair of microbes. We
38 then find relationships that are potentially invariants, meaning they will hold in any microbe
39 community. In other words, if we determine there is a relationship between two microbes, we
40 expect the relationship to hold in almost all context. We discovered that certain pairs of microbes
41 always exhibit the same relationship in almost all the datasets we studied, thus making them good
42 candidates for universal relationships. Our results confirm known biological properties and seem
43 promising in terms of disease diagnosis.

44
45 **Conclusions:** Since the relationships are likely universal, we expect that they will hold in a clinical
46 setting as well as in the general population. Strong universal relationships may provide insight on
47 prognostic, predictive, or therapeutic properties of a clinically relevant disease. These new

48 analyses may improve disease diagnosis and drug development in terms of accuracy and
49 efficiency.

50

51 **Keywords:** Boolean analysis, microbiome, invariants, systems biology, microbes interaction

52

53 **Background**

54 In recent years, microbiome research has progressed rapidly and there is an abundance of data
55 publicly available. It is important to understand these microbes that are found in organisms and
56 the environment, along with the relationships between them. There is also a growing interest to
57 find the connection between microbiomes and the healthy or diseased states of an organism.
58 Therefore, it is vital to efficiently analyze publicly available microbiome data, which will highlight
59 the microorganisms and their interactions with the host.

60

61 Current methods of analysis fundamentally use Pearson's correlation coefficient to determine
62 relationships within a microbiome, such as in co-occurrence networks (1). Correlation only
63 identifies relationships whose distribution is symmetric or linear. Therefore, this method may only
64 work for certain data in which there is a linear relationship between the two variables, but not for
65 asymmetric data. In some cases, past studies disregard asymmetric data by concluding a weak
66 correlation due to its nonlinearity, resulting in an incomplete analysis of asymmetric data. In
67 addition, there may be other types of relationships that cannot be identified through the standard
68 methods of linear analysis currently in use.

69

70 Past research used Boolean analysis to find the connection between oral microbiome and HIV-
71 associated periodontists (2) and to find a metabolic network of interactions in the gut microbiome
72 (3). However, these methods tend to analyze smaller datasets, which may have reproducibility
73 issues. Additionally, these studies focus on specific regions, such as the mouth and gut
74 microbiomes.

75

76 Instead, we propose using Boolean analysis, a logical method to determine dependency between
77 two variables in order to study large amounts of microbiome data, comprehensively. This method
78 of Boolean implication analysis was previously used to analyze relationships between genes using
79 publicly available microarray datasets for humans, mice, and fruit flies. Our research adapts the
80 method of analysis using Boolean implications, which was successfully used to discover markers
81 of blood stem cells (4), progenitors and a branch point in B-cell and T-cell differentiation (5), and
82 has also been applied to study colon (6, 7), bladder (8) and prostate cancer (9). Since this method
83 of Boolean analysis was previously used on gene expression data, we want to demonstrate the
84 universality of this method by analyzing pairwise microbe relationships.

85

86 In our process, the frequency of the microbes is first classified as either ‘low’ or ‘high’ using a
87 threshold that is derived for each microbe species. An example of a Boolean implication rule is “if
88 there is a high number of microbe A, then there will almost always be a low number of microbe
89 B”, or $A \text{ high} \rightarrow (\text{implies}) B \text{ low}$. There are six possible Boolean implication relationships: four
90 asymmetric ($A \text{ low} \rightarrow B \text{ low}$, $A \text{ low} \rightarrow B \text{ high}$, $A \text{ high} \rightarrow B \text{ low}$, and $A \text{ high} \rightarrow B \text{ high}$) and two
91 symmetric relationships ($A \text{ equivalent } B$, and $A \text{ opposite } B$).

92

93 Although our particular method of Boolean implication analysis has not been widely used in the
94 area of microbiology, a group of researchers did attempt to use the method to perform analysis on
95 microbiome data (10). In their research, they analyzed environmental data from Visualization and
96 Analysis of Microbial Population Structures (VAMPS)
97 (<http://vamps.mbl.edu/diversity/diversity.php>), and focused specifically on marine microbes.
98 Their research differs from ours in the aspect of diversity, as we want to determine relationships
99 that are present in a broader range of microbiomes, including humans, animals, and the
100 environment. While this research demonstrates the impact Boolean implication analysis will have
101 on microbiome analysis, we want to incorporate more as well as a greater variety of samples. By
102 having a wider array of samples, we want to find microbe relationships that not only exist in the
103 environment, but also in humans and animals.

104
105 Since Boolean implication analysis is able to capture relationships that are often overlooked in the
106 existing methods of analysis, we aim to uncover candidate invariants between pairs of microbes
107 that are likely applicable to every microbe community. For example, if we find a recurring Boolean
108 relationship between two microbe species in our large, diverse datasets, we expect this relationship
109 to be a promising candidate invariant. This research includes more diverse datasets and a novel
110 mathematical model compared to past studies, which helps produce stronger universal candidates.

111
112 This project's goal is to identify Boolean relationships that are comprehensive and likely universal.
113 Through Boolean implication analysis, we will be able to determine candidate universal invariants
114 within diverse microbe communities. If there are general rules that any microbe community
115 follows, it would make the process of identifying diseases easier because we expect these

116 relationships to hold in the general population. These candidate universal invariants provide a basis
117 which scientists could use to determine which microbes are associated with diseases. Therapeutic
118 use of microbes depends on their reproducibility in the general population, which makes our
119 approach more suitable for discovering appropriate microbes.

120

121 **Results**

122 Many universal Boolean relationships were uncovered using the proposed method of Boolean
123 implication analysis. Out of the 29,743,824 relationships discovered, some of the stronger ones
124 are presented in Figures 2 and 3, and are described below.

125

126 **Boolean implication relationships are conserved across environments and species**

127 The high \rightarrow low Boolean implication shows a high frequency of *Akkermansia muciniphila*
128 (Operational Taxonomic Unit (OTU) ID 4306262) implying a low frequency of Stramenopiles
129 (OTU ID 4350498) (Figure 2c). *A. muciniphila* is a human gut bacteria linked to preventing
130 obesity, diabetes, and inflammation (11) and Stramenopiles is found in aquatic environments,
131 mostly made up of algae (12). Since these two microbes are rarely found in similar environments,
132 it makes sense that when one microbe's frequency is high, the other is low. This logic is consistent
133 with the high \rightarrow low Boolean relationship found.

134

135 The graph in Figure 2e displays a strong low \rightarrow low relationship, showing that when
136 *Polynucleobacter* (OTU ID 145533) is low, *Candidatus Xiphinematobacter* (OTU ID 786420) is
137 also low. This relationship is confirmed in other studies which found *Polynucleobacter* makes up
138 a large portion of freshwater bacterioplankton (13) and *Candidatus Xiphinematobacter* is a known

139 nutrient supplier to nematodes, which are abundant in freshwater environments (14). It is
140 presumable that a low frequency count of *Polynucleobacter* indicates an environment that does
141 not contain freshwater; therefore, it is unlikely that the frequency of *Candidatus*
142 *Xiphinematobacter* is high, further confirming this low → low Boolean implication.

143

144 **Boolean relationships confirm some known biological properties**

145 A strong high → high relationship that was found (Figure 2d) is between *Corynebacterium* (OTU
146 ID 1062356) and *Staphylococcus aureus* (OTU ID 4446058). *Corynebacterium* and *S. aureus*
147 species both reside in the nose trail and skin microbiota of humans. *S. aureus* can be pathogenic
148 and cause infections. Studies have shown that *Corynebacterium* spp. and *S. aureus* reside together,
149 meaning they are positively correlated (15). In addition to being positively correlated, the
150 *Corynebacterium* high → *S. aureus* high relationship reveals that it is also possible to have a low
151 frequency of *Corynebacterium* and a high frequency of *S. aureus*.

152

153 An example of a symmetric relationship is shown in Figure 2f, where *Corynebacterium* (OTU ID
154 1062356) is equivalent to *Corynebacterium* (OTU ID 282360). *Corynebacterium* are a family of
155 Gram-positive bacteria with a large number of known species due to interest in the medical field
156 (16). However, the specific species of these *Corynebacterium* are not stated in the GreenGenes
157 database. Further analysis could determine the species for these *Corynebacterium*, which would
158 allow us to confirm this relationship. Although the species are not known, there is a symmetric
159 relationship between the two species in that as the count of one increases, so does the other.

160

161 **Microbes yield different Boolean implications in environmental versus animal samples**

162 In the next two examples in Section 1 of Figure 3, green represents environmental samples and red
163 represents the animal samples, including both humans and animals. Figure 3a shows the
164 relationship *Polynucleobacter* (OTU ID 145533) high \rightarrow *Candidatus Xiphinematobacter* (OTU
165 ID 786420) low. The abundance of green and the absence of red samples suggest that this
166 relationship is only present in the environmental microbiomes, and is not present in animal
167 microbiomes. As stated previously, *Polynucleobacter* makes up a large portion of freshwater
168 bacterioplankton (13), while *Candidatus Xiphinematobacter* tends to be found in soil samples (17),
169 meaning they are both environmental microbes.

170

171 Figure 3b presents the relationship *Polynucleobacter* (OTU ID 3071019) high \rightarrow *Bacteroides*
172 *uniformis* (OTU ID 197072) low. This Boolean relationship suggests that *Polynucleobacter* is
173 mainly present in the environmental microbiome while *Bacteroides uniformis* mostly exist in the
174 animal microbiome. Previous studies have shown that *Bacteroides uniformis* is one of the main
175 bacterial species of the human gut microbiome (18).

176

177 **Different body sites affect the presence of microbes relationships**

178 In Section 2 of Figure 3, both plots show the relationship between *Staphylococcus aureus* (OTU
179 ID 446058) and *Corynebacterium* (OTU ID 1000986) but in different regions of the human body.
180 While Figure 3c suggests that *S. aureus* and *Corynebacterium* have a low \rightarrow low relationship in
181 the skin region, based on Figure 3d, there is no specific Boolean relationship between these two
182 microbes when they are present in human feces. This emphasizes the importance of different
183 factors, like body region, in the way two microbes behave in the presence of each other.

184

185 **Boolean implications using disease-specific microbes is promising in potential diagnosis**

186 Irritable bowel disease (IBD) is a gastrointestinal disorder that currently is difficult to treat, but
187 treatments using the gut microbiome have been proposed (19). The relationship *Actinomyces* (OTU
188 ID 12564) high → *Lachnospiraceae* (OTU ID 4469576) low (Figure 3e) specifically highlights
189 samples from patients that either have Crohn’s Disease (CD), Ulcerative Colitis (UC), or neither
190 (No IBD). Numerous studies show that there is a decreased amount of Clostridiales
191 (*Lachnospiraceae* is in the class Clostridiales) in patients with IBS (19-22), but not much
192 information about IBD. Although there have been studies that show there is a connection between
193 *Actinomyces* and infections (23), studies have not compared the amount of *Actinomyces* to IBD.
194 This indicates that more research needs to be done to determine if Clostridiales and *Actinomyces*
195 are connected to certain types of IBD. There tends to be a higher proportion of *Lachnospiraceae*
196 in patients with CD than UC, and a higher proportion of *Actinomyces* in patients of UC versus CD.
197 The relationship *Streptococcus* (OTU ID 4467992) high → *Lachnospiraceae* (OTU ID 4469576)
198 low (Figure 3f) also highlights differences between samples of IBD patients. A similar trend of a
199 higher proportion of *Lachnospiraceae* in CD patients and a higher proportion of *Streptococcus* in
200 UC patients appears with these microbes. There is no previous research suggesting there is a
201 connection between *Streptococcus* and IBD, but more research needs to be done to see if this
202 relationship has any disease-identifying properties. Relationships like these might have only been
203 discovered using methods taking into account asymmetry like Boolean implication analysis which
204 may be why there are no studies that can confirm these relationships.

205

206 Certain microbes seem to be related to different skin conditions, such as eczema and psoriasis. The
207 relationship *Acinetobacter johnsonii* (OTU ID 4482374) low → *Corynebacterium* (OTU ID

208 361600) low shows how patients with psoriasis tend to have higher counts of both *Acinetobacter*
209 *johnsonii* and *Corynebacterium* than patients with eczema and patients with neither skin condition.
210 In looking at another relationship, *Ruminococcaceae* (OTU ID 4346675) high → *Anaerococcus*
211 (OTU ID 927089) low, it is clear that patients with psoriasis have higher counts of *Anaerococcus*,
212 while patients with eczema have higher counts of *Ruminococcaceae*, with both having minimal
213 amounts of the other microbes. Although there has been research to detect microbial diversity on
214 the skin, experts still cannot agree on a universal method of using microbes to diagnose psoriasis
215 in patients (24). Our method of Boolean implication analysis attempts to provide a mathematical
216 model of identifying candidate universal invariants. This will allow the properties determined
217 within these microbiomes to apply in almost all situations, and hopefully provide treatments for
218 diseases like these that will be universally successful.

219

220 **Discussion**

221 There is potential in using Boolean implication analysis to comprehensively determine microbial
222 relationships, which can then be used to build abstract versions of biological systems.
223 Understanding and simulating biological systems has always been the goal of researchers, but
224 current analysis has not met that objective due to the lack of complexity that symmetric analysis
225 has, and the smaller datasets researchers tend to use. Relationships between microbes and diseases
226 have always been evident, so our research intends to build the foundation of the biological system.
227 Specified research will provide more concrete evidence of the connection between microbes and
228 diseases.

229

230 Our results reveal that Boolean analysis is a promising method for analysis of different
231 microbiomes. After analyzing more than 400 diverse datasets consisting of over 100,000 samples,
232 we uncovered candidate invariants that held in all our datasets, which is consistent with our
233 hypothesis. However, only four of the six Boolean implication relationships were found in the
234 datasets: low \rightarrow low, high \rightarrow low, high \rightarrow high, and equivalent. The other two relationships,
235 opposite and low \rightarrow high, did not appear in the datasets. In both of these cases, the low, low
236 quadrant is considered sparse, meaning there will be a high number of at least one microbe in most
237 of the samples. Due to the diverse nature of microbes, we know that certain microbes are not
238 present in every type of sample. Therefore, it should be rare for the low, low quadrant to be sparse,
239 which is justified by the lack of opposite and low \rightarrow high relationships in our results. Boolean
240 analysis also unveils differences in microbe interactions due to factors such as body site, disease,
241 and environment. Each implication is believed to be a universal candidate because it holds in all
242 the datasets we analyzed. These microbe relationships can be validated in the lab by generating
243 and sequencing additional samples to confirm these relationships.

244

245 One of the limitations in using Boolean analysis in which the data focuses on the stronger
246 relationships, making the analysis less noisy. However, the weaker relationships are lost in the
247 process. Further analysis might prove these weaker relationships have significance, but this
248 method will only focus on the stronger relationships. A second limitation is that we only analyzed
249 datasets downloaded from Qiita that were processed using the GreenGenes database (25).
250 However, the latest version of GreenGenes was published in 2013, which may not include the most
251 up to date information involving microbiome taxonomy. Additionally, since we are not focusing
252 on microbiomes found in a specific region, we are limiting our scope to focus on microbes that are

253 found universally in humans, animals, and the environment. We may be excluding relationships
254 that are compelling due to the fact that they might not be found everywhere. Future studies using
255 this method of Boolean analysis should be done to focus on specific regions, which will give better
256 insight into particular microbiomes, such as the gut microbiome.

257

258 **Conclusions**

259 The lack of comprehensive analysis of microbiome data created a need for more extensive
260 approaches. Boolean implication analysis presents a solution that takes into account both
261 symmetric and asymmetric relationships. Our results show that some biological properties were
262 confirmed by Boolean analysis. For example, it is proven that the *Corynebacterium* and *S. aureus*
263 species reside together and are positively correlated, which is consistent with the high → high
264 relationship found. Our results also show that different microbiomes affect the presence of
265 microbe relationships, on a broader scope such as environmental versus animal samples, or on a
266 smaller scope such as various body sites. Boolean implication analysis is promising in terms of
267 potential disease diagnosis, as well such as IBD. We found that higher frequencies of certain
268 microbes seem to be associated with either CD or UC.

269

270 Future work includes building a Boolean implication network to further analyze how microbe
271 implications are connected to each other. A Boolean implication network with the candidate
272 microbe invariants may help in developing better models for biological systems. Our research also
273 helps determine strong properties of biological systems, and future research on this topic provide
274 novel directions in understanding biological system. Invariants help formulate new theories that
275 may provide more effective diagnostic and therapeutic applications.

276

277 **Materials and Methods**

278 **Data Collection**

279 We extracted pre-processed OTU tables along with the corresponding metadata from Qiita (26), a
280 microbiome database and study management platform. Qiita uses third party plugins including
281 QIIME (Quantitative Insights Into Microbial Ecology) (<http://qiime.org/>) or QIIME 2
282 (<https://qiime2.org/>) to process microbial 16S rRNA sequences of each study, which are
283 contributed by users on this platform. Qiita classifies the microbes using the GreenGenes database,
284 and generates into OTU tables. OTU tables display the frequencies of all the microbes species
285 present in each of the samples. Each study and the corresponding raw data comes from different
286 individuals and institutes, which makes this database comprehensive. The metadata includes
287 information about the samples, such as location and sample identification. For easier analysis of
288 the collected data, we separated the downloaded studies from QIITA and pooled them into 4
289 datasets. These OTU tables and metadata were uploaded onto a web-based tool for analyzing big
290 data called Hegemon (7, 8, 27-29).

291

292 **Boolean Implication Analysis**

293 To classify a relationship, thresholds are first determined for each microbe using the StepMiner
294 algorithm. The StepMiner algorithm (30) is a tool that helps identify step-wise transitions (either
295 step-up or step-down transitions) calculated using sum-of-square errors. Steps are defined as the
296 sharpest change between low microbe frequency count and high microbe frequency count. In order
297 to fit a step function, the StepMiner algorithm computes the average of the values on both sides of
298 the step for all possible step positions. The midpoint of the step position that minimizes the square

299 error is chosen as the threshold for each respective microbe. The step is placed at the largest jump
300 from low values to high values, and the sets the threshold at the point where the step crosses the
301 original data. If the microbe frequencies are evenly distributed from low to high, the mean
302 frequency level tends to be used. Microbe frequencies are further classified as either ‘high’, ‘low’,
303 or ‘intermediate’. If t is the microbe frequency threshold, frequency levels above $t + 0.5$ are ‘high’,
304 levels below $t - 0.5$ are ‘low’, and levels between $t - 0.5$ and $t + 0.5$ are ‘intermediate’. Points in
305 the intermediate region are ignored because these points might appear on either side of the
306 threshold due to noise.

307

308 The OTU tables are uploaded onto Hegemon to visualize the samples on scatter plots, comparing
309 two microbes against each other. All microbe pairs are analyzed to determine if any Boolean
310 relationships are present. In each graph, one microbe species’ frequency (using OTU ID A) is
311 plotted on the x-axis, and another microbe species’ frequency (OTU ID B) is plotted on the y-axis.
312 Each data point represents a sample, and the frequencies are plotted on a log-log scale. From here,
313 Boolean implication analysis is applied to the scatter plots. Using the graphs constructed on
314 Hegemon, microbe pairs are analyzed to determine if a Boolean relationship is present. There are
315 six possible Boolean implications: symmetric (opposite and equivalent) or asymmetric (low \rightarrow
316 low, low \rightarrow high, high \rightarrow low, high \rightarrow high).

317

318 The asymmetric are determined by checking if one of the four quadrants in the scatter plot is
319 significantly sparse compared with other quadrants. If A low \rightarrow B low and A high \rightarrow B high are
320 both sparsely populated, then A is equivalent to B. If A high \rightarrow B low and A low \rightarrow B high are
321 both sparsely populated, then A is opposite to B. The BooleanNet statistics tests (27) determine

322 whether there is a Boolean relationship between A and B. Consider the relationship A low \rightarrow B
 323 high. First, test if the microbe frequencies in the sparse quadrant is significantly less than the
 324 expected frequencies in an independence model. Let a_{00} , a_{01} , a_{10} , and a_{11} represent the quadrants in
 325 which the microbe frequencies of A and B are low and low, low and high, high and low, and high
 326 and high, respectively.

327

$$328 \quad total = a_{00} + a_{01} + a_{10} + a_{11}$$

$$329 \quad \text{number of A low expression values} = nA_{low} = (a_{00} + a_{01})$$

$$330 \quad \text{number of B low expression values} = nB_{low} = (a_{00} + a_{10})$$

$$331 \quad expected = \left(\frac{nA_{low}}{total} \times \frac{nB_{low}}{total} \right) \times total = (a_{00} + a_{01}) \times \frac{(a_{00} + a_{10})}{total}$$

$$332 \quad observed = a_{00}$$

$$333 \quad S \text{ statistic} = \frac{expected - observed}{\sqrt{expected}}$$

334 Second, the observed values in the sparse quadrant are considered erroneous points, so a sparse
 335 quadrant should have a small number of these erroneous points. A maximum likelihood estimate
 336 of the error rate is then computed:

337

$$338 \quad errorrate = \frac{1}{2} \left(\frac{a_{00}}{a_{00}+a_{01}} + \frac{a_{00}}{a_{00}+a_{10}} \right)$$

339

340 If both tests succeed, the low-low quadrant is considered sparse, so the implication A low \rightarrow B
 341 high is true. An implication is considered significant if the S statistic is greater than 3 and the
 342 error rate is less than 0.1. All the microbe relationships that pass both of these tests are now
 343 considered candidate invariants.

344

345 **Analysis of metadata**

346 We used a Hegemon function that calculates the differential analysis for specific factors in the
347 metadata using t-tests in R software framework (R version 3.4.4 - 2018-03-15). We then selected
348 the OTU IDs that had higher mean differential values and higher $-\log_{10}(p)$ values. After the list of
349 potential IDs were generated, we visually confirmed on Hegemon whether certain factors such as
350 environment versus animal or body site affects the microbe counts.

351

352 **List of Abbreviations**

353 CD: Crohn's Disease

354 IBD: Irritable Bowel Disease

355 QIIME: Quantitative Insights Into Microbial Ecology

356 OTU: Operational Taxonomic Units

357 UC: Ulcerative Colitis

358 VAMPS: Visualization and Analysis of Microbial Population Structures

359

360 **Declarations**

361 **Ethics approval and consent to participant**

362 Publicly available data were generated after appropriate ethics approval and consent to participate.

363

364 **Consent for publication**

365 Not Applicable.

366

367 **Availability of data and materials**

368 All dataset are publicly available from <https://qiita.ucsd.edu>, under the Study, and View Studies
369 tab. We selected a subset of all the available studies. The specific list of studies used for this
370 research is included as an additional file.

371

372 **Competing interests**

373 The authors declare that they have no competing interests.

374

375 **Funding**

376 DS is supported by NIH UG3TR003355, UG3TR002968, R01AI155696, R01GM138385 and
377 R00CA151673, Hyundai Hope On Wheels, Padres Pedal the Cause/Rady Children's Hospital
378 Translational PEDIATRIC Cancer Research Award (Padres Pedal the Cause/RADY #PTC2017),
379 2017, Padres Pedal the Cause/C3 Collaborative Translational Cancer Research Award (San Diego
380 NCI Cancer Centers Council [C3] #PTC2017).

381

382 **Authors' contributions**

383 DV, SCS, SS, DS conceived and designed the analysis, collected the data, contributed analysis
384 tools, performed the analysis, wrote the paper. DS supervised, acquired funding, reviewed, and
385 edited the paper.

386

387 **Acknowledgements**

388 The authors thank Dr. Christine Alvarado for her guidance and the members of UCSD's Boolean
389 Lab for reviewing our work and providing constructive criticism.

390

391 **References**

- 392 1. Bauer E, Thiele I. From Network Analysis to Functional Metabolic Modeling of the Human
393 Gut Microbiota. *mSystems*. 2018;3(3).
- 394 2. Noguera-Julian M, Guillén Y, Peterson J, Reznik D, Harris EV, Joseph SJ, et al. Oral
395 microbiome in HIV-associated periodontitis. *Medicine (Baltimore)*. 2017;96(12):e5821.
- 396 3. Steinway SN, Biggs MB, Loughran TP, Papin JA, Albert R. Inference of Network
397 Dynamics and Metabolic Interactions in the Gut Microbiome. *PLoS Comput Biol*.
398 2015;11(5):e1004338.
- 399 4. Chen JY, Miyanishi M, Wang SK, Yamazaki S, Sinha R, Kao KS, et al. Hoxb5 marks long-
400 term haematopoietic stem cells and reveals a homogenous perivascular niche. *Nature*.
401 2016;530(7589):223-7.
- 402 5. Inlay MA, Bhattacharya D, Sahoo D, Serwold T, Seita J, Karsunky H, et al. Ly6d marks
403 the earliest stage of B-cell specification and identifies the branchpoint between B-cell and T-cell
404 development. *Genes Dev*. 2009;23(20):2376-81.
- 405 6. Dalerba P, Sahoo D, Clarke MF. CDX2 as a Prognostic Biomarker in Colon Cancer. *N*
406 *Engl J Med*. 2016;374(22):2184.
- 407 7. Dalerba P, Sahoo D, Paik S, Guo X, Yothers G, Song N, et al. CDX2 as a Prognostic
408 Biomarker in Stage II and Stage III Colon Cancer. *N Engl J Med*. 2016;374(3):211-22.

- 409 8. Volkmer JP, Sahoo D, Chin RK, Ho PL, Tang C, Kurtova AV, et al. Three differentiation
410 states risk-stratify bladder cancer into distinct subtypes. *Proc Natl Acad Sci U S A*.
411 2012;109(6):2078-83.
- 412 9. Sahoo D, Wei W, Auman H, Hurtado-Coll A, Carroll PR, Fazli L, et al. Boolean analysis
413 identifies CD38 as a biomarker of aggressive localized prostate cancer. *Oncotarget*.
414 2018;9(5):6550-61.
- 415 10. Zhu C, Jiang R, Chen T. Constructing a Boolean implication network to study the
416 interactions between environmental factors and OTUs. *Quantitative Biology*. 2014;2(4):127-41.
- 417 11. Depommier C, Everard A, Druart C, Plovier H, Van Hul M, Vieira-Silva S, et al.
418 Supplementation with *Akkermansia muciniphila* in overweight and obese human volunteers: a
419 proof-of-concept exploratory study. *Nat Med*. 2019;25(7):1096-103.
- 420 12. Derelle R, López-García P, Timpano H, Moreira D. A Phylogenomic Framework to Study
421 the Diversity and Evolution of Stramenopiles (=Heterokonts). *Mol Biol Evol*. 2016;33(11):2890-
422 8.
- 423 13. Hahn MW, Jezberová J, Koll U, Saueressig-Beck T, Schmidt J. Complete ecological
424 isolation and cryptic diversity in *Polynucleobacter* bacteria not resolved by 16S rRNA gene
425 sequences. *ISME J*. 2016;10(7):1642-55.
- 426 14. Brown AM, Howe DK, Wasala SK, Peetz AB, Zasada IA, Denver DR. Comparative
427 Genomics of a Plant-Parasitic Nematode Endosymbiont Suggest a Role in Nutritional Symbiosis.
428 *Genome Biol Evol*. 2015;7(9):2727-46.

- 429 15. Ramsey MM, Freire MO, Gabriliska RA, Rumbaugh KP, Lemon KP. *Staphylococcus*
430 *aureus* Shifts toward Commensalism in Response to *Corynebacterium* Species. *Front Microbiol.*
431 2016;7:1230.
- 432 16. Liebl W. *Corynebacterium* taxonomy. *Handbook of Corynebacterium glutamicum* CRC
433 Press, Boca Raton, FL. 2005:9-34.
- 434 17. Lazarova S, Peneva V, Kumari S. Morphological and molecular characterisation, and
435 phylogenetic position of *X. browni* sp. n., *X. penevi* sp. n. and two known species of *Xiphinema*
436 *americanum*-group (Nematoda, Longidoridae). *Zookeys*. 2016(574):1-42.
- 437 18. Renouf M, Hendrich S. *Bacteroides uniformis* is a putative bacterial species associated
438 with the degradation of the isoflavone genistein in human feces. *J Nutr*. 2011;141(6):1120-6.
- 439 19. Pittayanon R, Lau JT, Yuan Y, Leontiadis GI, Tse F, Surette M, et al. Gut Microbiota in
440 Patients With Irritable Bowel Syndrome-A Systematic Review. *Gastroenterology*.
441 2019;157(1):97-108.
- 442 20. Jalanka-Tuovinen J, Salojärvi J, Salonen A, Immonen O, Garsed K, Kelly FM, et al. Faecal
443 microbiota composition and host-microbe cross-talk following gastroenteritis and in postinfectious
444 irritable bowel syndrome. *Gut*. 2014;63(11):1737-45.
- 445 21. Carroll IM, Ringel-Kulka T, Siddle JP, Ringel Y. Alterations in composition and diversity
446 of the intestinal microbiota in patients with diarrhea-predominant irritable bowel syndrome.
447 *Neurogastroenterol Motil*. 2012;24(6):521-30, e248.

- 448 22. Li G, Yang M, Jin Y, Li Y, Qian W, Xiong H, et al. Involvement of shared mucosal-
449 associated microbiota in the duodenum and rectum in diarrhea-predominant irritable bowel
450 syndrome. *J Gastroenterol Hepatol*. 2018;33(6):1220-6.
- 451 23. Könönen E, Wade WG. Actinomyces and related organisms in human infections. *Clin*
452 *Microbiol Rev*. 2015;28(2):419-42.
- 453 24. Langan EA, Griffiths CEM, Solbach W, Knobloch JK, Zillikens D, Thaçi D. The role of
454 the microbiome in psoriasis: moving from disease description to treatment selection? *Br J*
455 *Dermatol*. 2018;178(5):1020-7.
- 456 25. DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, et al. Greengenes, a
457 chimera-checked 16S rRNA gene database and workbench compatible with ARB. *Appl Environ*
458 *Microbiol*. 2006;72(7):5069-72.
- 459 26. Gonzalez A, Navas-Molina JA, Kosciulek T, McDonald D, Vázquez-Baeza Y, Ackermann
460 G, et al. Qiita: rapid, web-enabled microbiome meta-analysis. *Nat Methods*. 2018;15(10):796-8.
- 461 27. Sahoo D, Dill DL, Gentles AJ, Tibshirani R, Plevritis SK. Boolean implication networks
462 derived from large scale, whole genome microarray datasets. *Genome Biol*. 2008;9(10):R157.
- 463 28. Sahoo D, Seita J, Bhattacharya D, Inlay MA, Weissman IL, Plevritis SK, et al. MiDReG:
464 a method of mining developmentally regulated genes using Boolean implications. *Proc Natl Acad*
465 *Sci U S A*. 2010;107(13):5732-7.

466 29. Dalerba P, Kalisky T, Sahoo D, Rajendran PS, Rothenberg ME, Leyrat AA, et al. Single-
467 cell dissection of transcriptional heterogeneity in human colon tumors. Nat Biotechnol.
468 2011;29(12):1120-7.

469 30. Sahoo D, Dill DL, Tibshirani R, Plevritis SK. Extracting binary signals from microarray
470 time-course data. Nucleic Acids Res. 2007;35(11):3705-12.

471

472

473 **Figure Legends:**

474 **Figure 1. Study design**

475 The overview of the research process: (a) OTU tables were collected from publicly available
476 microbiome datasets. (b) Tables were uploaded to Hegemon and all possible microbe pairs were
477 plotted (Using 4 plots as examples - The number of total plots is larger). (c) Boolean analysis was
478 performed on all the plots. (d) The plots that passed the BooleanNet statistics tests were marked
479 as candidate invariants for further analysis and validation. (e) Any of the determined candidates
480 that can be validated in other datasets, represents a likely universal invariant (a rule between two
481 microbes that holds between them, any time the pair are present together in any environment).

482 *Note that this is just an example. A universal invariant can be any of the 6 possible Boolean
483 relationship.

484

485 **Figure 2. Boolean implication relationships represent diversity in microbiome data**

486 All four types of Boolean relationships found in our dataset. (a) Describes the StepMiner algorithm
487 that creates thresholds for each microbe. Among all samples, the frequencies of a particular
488 microbe are sorted and a step function is fitted where the sharpest change between low microbe

489 count and high microbe count takes place. The midpoint of the step position that minimizes the
490 square error is chosen as the threshold for each respective microbe. (b) Depicts a log-log plot of
491 the number of each type of relationship and the corresponding number of microbes that exhibit
492 that specific relationship. The remaining diagrams display each of the four types of relationships
493 found. Each point in the scatter plot corresponds to a sample, where the two axes represent the
494 frequency counts of each microbes. (c) high → low (d) high → high (e) low → low (f) equivalent
495

496 **Figure 3. Boolean implication reveals strong patterns in diverse biological and**
497 **environmental conditions**

498 Analysis of scatter plots with various experimental conditions using metadata files that provided
499 additional information about the samples. *Section 1:* green represents environmental samples
500 (plants, water, soil, etc.) and red indicates animal samples (humans, animals). (a) *Polynucleobacter*
501 (145533) low → *Candidatus Xiphinematobacter* (786420) low; this relationship is only present in
502 environmental microbiomes due to the lack of red samples in the plot. (b) *Polynucleobacter*
503 (3071019) high → *Bacteroides uniformis* (197072) low; this relationship suggests
504 *Polynucleobacter* is mainly present in the environment, and *Bacteroides uniformis* is mainly
505 present in animals. *Section 2:* (c) and (d) have the same microbes on the axes *Staphylococcus*
506 *aureus* (446058) and *Corynebacterium* (1000986), but different regions of the body plotted: skin
507 (dark blue) and feces (green). (c) shows the relationship *S. aureus* low → *Corynebacterium* low
508 holds for the skin region. (d) shows the relationship using fecal samples, and there is not a clear
509 relationships that can be determined from this. *Section 3:* Pink represents Crohn's Disease (CD),
510 teal represents Ulcerative Colitis (UC), and light gray represents neither disease. (e) The
511 relationship *Actinomyces* (12574) high → *Lachnospiraceae* (4469576) low is shown, with higher

512 counts of *Lachnospiraceae* in CD, and higher counts of *Actinomyces* in UC. (f) shows the
513 relationship *Streptococcus* (4467992) high → *Lachnospiraceae* (4469576) low, with higher counts
514 of *Lachnospiraceae* in CD, while higher counts of *Streptococcus* in UC. *Section 4*: Magenta
515 represents eczema, blue represents psoriasis, and beige represents neither skin condition. (g) The
516 relationship *Acinetobacter johnsonii* (4482374) low → *Corynebacterium* (361600) low is shown.
517 *Patients with psoriasis tend to have higher counts of Corynebacterium* than patients with eczema.
518 (h) The relationship *Ruminococcaceae* (4346675) high → *Anaerococcus* (927089) low is shown.
519 Patients with psoriasis tend to have higher counts of *Anaerococcus*, while patients with eczema
520 tend to have higher counts of *Ruminococcaceae*.