## 1 Title

## 2 Functional Anabolic Network Analysis of Human-associated *Lactobacillus* Strains

## 3 Authors

4 Thomas J. Moutinho Jr.[a], Benjamin C. Neubert[a], Matthew L. Jenior[a], Maureen A. Carey[a,b], Gregory L.

5 Medlock[a], Glynis L. Kolling[a], Jason A. Papin[a]#

6 [a] Department of Biomedical Engineering, University of Virginia, Charlottesville, Virginia, USA

7 [b] Division of Infectious Disease and International Health, Department of Medicine, University of Virginia,

8 Charlottesville, Virginia, USA

9 Running Head: Anabolic Network Analysis of Lactobacilli

10

11 #Address correspondence to Jason A. Papin, papin@virginia.edu

## 12 Abstract

13 Members of the *Lactobacillus* genus are frequently utilized in the probiotic industry with many species
14 conferring demonstrated health benefits; however, these effects are largely strain-dependent. We
15 designed a method called PROTEAN (Probabilistic Reconstruction Of constituent Anabolic Networks) to
16 computationally analyze the genomic annotations and predicted metabolic production capabilities of
17 144 strains across 16 species of *Lactobacillus* isolated from human intestinal, oral, and vaginal body
18 sites. Using PROTEAN we conducted a genome-scale metabolic network comparison between strains,
19 revealing that metabolic capabilities differ by isolation site. Notably, PROTEAN does not require a well-
20 curated genome-scale metabolic network reconstruction to provide biological insights. We found that
21 predicted metabolic capabilities of lactobacilli isolated from the vaginal microbiota cluster separately
22 from intestinal and oral isolates, and we also uncovered an overlap in the predicted metabolic
23 production capabilities of intestinal and oral isolates. Using machine learning, we determined the most
24 informative metabolic products driving the difference between predicted metabolic capabilities of
25 intestinal, oral, and vaginal isolates. Notably, intestinal and oral isolates were predicted to have a higher
26 likelihood of producing D-alanine, D/L-serine, and L-proline, while the vaginal isolates were
27 distinguished by a higher predicted likelihood of producing L-arginine, citrulline, and D/L-lactate. We
28 found the distinguishing products to be consistent with published experimental literature. This study
29 showcases a systematic technique, PROTEAN, for comparing the predicted functional metabolic output
30 of microbes using genome-scale metabolic network analysis and computational modeling and provides
31 unique insight into human-associated *Lactobacillus* biology.

## 32 Importance

33 The *Lactobacillus* genus has been shown to be important for human health. Lactobacilli have been
34 isolated from human intestinal, oral, and vaginal sites. Members of the genus contribute significantly to
35 the maintenance of vaginal health by providing colonization resistance to invading pathogens. A wide
36 variety of clinical studies have indicated that *Lactobacillus*-based probiotics confer health benefits for

37    several gut- and immune-associated diseases. Microbes interact with the human body in several ways,
38    including the production of metabolites that influence physiology or other surrounding microbes. We
39    have conducted a strain-level genome-scale metabolic network reconstruction analysis of human-
40    associated *Lactobacillus* strains, revealing that predicted metabolic capabilities differ when comparing
41    intestinal/oral isolate to vaginal isolates. The technique we present here allows for direct interpretation
42    of discriminating features between the experimental groups.

## Introduction

44    *Lactobacillus* is a diverse genus of bacteria with many member strains associated with the human body.
45    Lactobacilli are Gram-positive, lactic acid-producing bacteria typically with a low GC content (1,2). They
46    are known for their production of lactic acid, being facultative anaerobes, and are capable of being
47    metabolically active in a large variety of conditions (3). There is evidence that human-associated
48    lactobacilli colonize mucosal surfaces of the intestinal tract (4), vagina (5–12), and oral cavity (13,14).
49    While strains of *Lactobacillus* have been isolated from all three of these body sites, it remains unknown
50    which are permanent members of the resident microbiota (autochthonous) opposed to transient
51    members (allochthonous). Transient intestinal lactobacilli are either resident members of the oral
52    microbiota or have been ingested, most commonly from unpasteurized fermented foods (4,15).

53    Lactobacilli have been used for a broad range of applications primarily associated with human intestinal
54    probiotics and industrial production of useful metabolites. *Lactobacillus*-based probiotics have been
55    shown to confer health benefits in clinical studies for a variety of conditions including prevention of
56    antibiotic associated diarrhea (16), *Clostridium difficile*-associated diarrhea (17), constipation (18),
57    irritable bowel syndrome (19), and eczema/atopic dermatitis (20). Probiotics are controversial, likely due
58    to claims made by currently marketed probiotics that lack FDA approval for the treatment of specific
59    diseases (21,22). The primary benefits associated with lactobacilli-based probiotics may be a function of
60    their presence in the gut, production of metabolites, and modulation of the immune system (23,24).
61    Metabolism plays a key role in all three of these general mechanisms; therefore, a better understanding
62    of their metabolic capabilities will help to elucidate the mechanisms contributing to probiotic effects
63    (25).

64    In recent years, there has been an explosion of genomic and metagenomic sequencing of human-
65    associated microbiota, which provides a unique opportunity to apply genome-scale metabolic network
66    reconstructions (GENREs) to enhance our current understanding of human-associated lactobacilli
67    metabolism utilizing *in silico* techniques (25). Systems biology has the potential to advance design,
68    selection, and delivery of *Lactobacillus*-based probiotics (26,27). GENREs are a powerful computational
69    tool for mathematically modeling the metabolic processes within a cell at a systems-level, including all
70    known metabolic reactions, metabolites, and metabolic genes in an organism (28). GENREs are created
71    by referencing an annotated genome against biochemical databases, then integrating experimental data
72    when available (29). There are several examples of *Lactobacillus*-specific comparative genomics studies
73    (30–35); however, GENREs allow for a more functional perspective than genomics data alone because of
74    the quantitative accounting for interactions between components in the network (25,36). Simulations
75    with GENREs can accurately predict microbial growth yields and the metabolic pathways utilized for the
76    production of metabolites during exponential growth of a microbe (37). A variety of analytical
77    approaches can be applied to interrogate emergent properties of a GENRE. Flux Balance Analysis (FBA)
78    and related methods have proven highly successful in the analysis of metabolic networks (38). FBA is a

79    mathematical technique for analyzing the flow of metabolites through a GENRE; it can be used to
80    identify a set of reaction fluxes that maximize growth in a specified media condition among other
81    applications (28,39,40). Metabolic network reconstructions and FBA provide a mechanistic look into
82    cellular metabolism and are increasingly used to study biochemical processes of single bacterial species
83    as well as communities of organisms (41).

84    GENREs enable the computational prediction of metabolic capabilities of microbes, both catabolic and
85    anabolic. Additionally, GENREs are capable of contextualizing large 'omic datasets (i.e. genomics,
86    transcriptomics, and metabolomics) with known biochemistry and biological network architectures for
87    improved understanding of the experimental data (42). An important recent finding demonstrated that
88    metabolomics data alone can be used to differentiate between bacterial cultures at the strain level (43).
89    We developed a computational method using GENREs to predict the metabolic products that a strain is
90    likely able to produce. We used predicted production capabilities to then differentiate between
91    different human-associated *Lactobacillus* strains. Just as metabolomics data can be used to differentiate
92    bacterial strains, predicted production capabilities can be used for the same comparisons. We assessed
93    the metabolic potential across a broad set of *Lactobacillus* species, consisting of 144 strains, which have
94    all been isolated from three human-related body sites: intestinal, oral, and vaginal. We found that
95    intestinal and oral isolates have a great deal of overlap in their metabolic functionality, while vaginal
96    isolates have more unique metabolic production capabilities. These analyses can facilitate additional
97    experimental interrogation of this important genus of bacteria.

## Results and Discussion

98    

**Annotated metabolic genes associated with known metabolic functions are sufficiently represented**
**among human-associated lactobacilli**

99    
100   

101   In this study we predict the metabolic production capabilities of 144 lactobacilli strains. We utilized the
102   PATRIC Cross-Genus Protein Families (PGfams) (4) for an initial genomic analysis. PGfams are
103   comparable clusters of proteins that likely have similar functions. These clusters are intended to be used
104   for cross-genus comparison due to their slightly relaxed clustering criteria. However, PGfams allow for
105   the comparison of the large number of strains analyzed in this study. Lactobacilli consist of a broad
106   range of species and thus using the PGfams was appropriate for an initial genomic comparison in this
107   study. We first filtered the PGfams to only include metabolic gene families associated with known
108   metabolic functions (see Methods). The distribution of total metabolic PGfams associated with each
109   genome ranges from 340 to 580 and has a median value of 515 (Figure 1A). Across these 144 strains we
110   found that they share 116 core metabolic PGfams, spanning a variety of cellular functions including, but
111   not limited to, carbohydrate, nucleotide, and amino acid metabolism (Table S1). The pan set of
112   metabolic PGfams, which represents the total set of unique PGfams, expanded to over 1500 after
113   considering all strains utilized within this study (Figure 1B). The *Lactobacillus* strains we studied
114   consisted of 16 species and were isolated from intestinal, oral, and vaginal human body sites (Figure 1C).

A

B

C

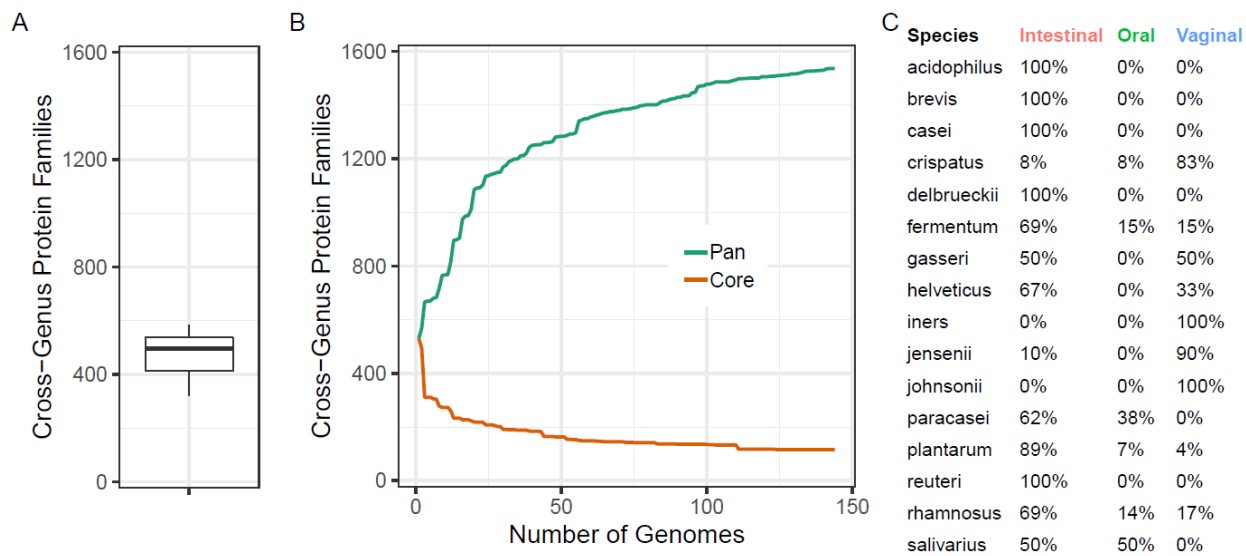| Species | Intestinal | Oral | Vaginal |
|---|---|---|---|
| acidophilus | 100% | 0% | 0% |
| brevis | 100% | 0% | 0% |
| casei | 100% | 0% | 0% |
| crispatus | 8% | 8% | 83% |
| delbrueckii | 100% | 0% | 0% |
| fermentum | 69% | 15% | 15% |
| gasseri | 50% | 0% | 50% |
| helveticus | 67% | 0% | 33% |
| iners | 0% | 0% | 100% |
| jensenii | 10% | 0% | 90% |
| johnsonii | 0% | 0% | 100% |
| paracasei | 62% | 38% | 0% |
| plantarum | 89% | 7% | 4% |
| reuteri | 100% | 0% | 0% |
| rhamnosus | 69% | 14% | 17% |
| salivarius | 50% | 50% | 0% |



**Figure 1: Known metabolic annotations are extensively sampled across the 16 *Lactobacillus* species included in this study.** The genomic features used for this analysis are PATRIC Cross-Genera Protein families (PGfams), a standardized set of features across the PATRIC Database (4). (A) The number of metabolic PGfams for each genome are shown here, with the median value indicated by the middle line in the boxplot. (B) For the 144 strains from 16 species of *Lactobacillus*, we found that there are 116 protein families in the core set of metabolic PGfams, while the pan set of PGfams expands to over 1500 families. The nearly plateau shape of the curve for the pan set of PGfams curve indicates that this sampling represents a large portion of the genetic diversity among the 16 species included in the study. (C) This table shows the complete list of species used in this study and indicates the percentage of strains that were isolated from each human body site. Each strain in this study is a member from one of the 16 species and isolated from one of three human-associated body sites; intestinal, oral, or vaginal (Table S2).

**Probabilistic Reconstruction Of constituent Anabolic Networks (PROTEAN)**

We developed PROTEAN to predict the metabolic production capabilities of microbes based on genomic data alone. PROTEAN generates constituent metabolic production networks with maximum parsimony and probability to predict the production of a given metabolite with a defined set of input metabolites. PROTEAN is a combination of well-validated methods, including Parsimonious Enzyme Usage Flux Balance Analysis (pFBA) (37), likelihood-based gap filling (44), fastGapFill (45), and CarveMe (46). The algorithm uses the ModelSEED biochemical reaction database, a large set of known metabolic reactions, for constituent network generation (47). First, reaction likelihoods are calculated for each reaction in the ModelSEED database using Probannopy (48) (Figure 2). Reaction likelihoods correspond to the probability that a given reaction is catalyzed by an enzyme that is encoded for by the genome. We modified pFBA to utilize reaction likelihoods for weighted minimization of flux through each reaction, while still maintaining near-optimal flux through the objective function. Standard pFBA assumes that metabolism is optimized to minimize enzymatic turnover and thus the method is driven by a minimization of the total flux through the metabolic network (37). Weighted pFBA allows for the reconstruction of constituent anabolic networks while accounting for maximum genomic probability and resource parsimony (see Methods). The constituent anabolic networks output by PROTEAN consist of flux-carrying reactions required for the production of a certain metabolite with preferential flux through reactions that have higher reaction likelihoods. A constituent network represents a theoretically optimal

145    biosynthetic network while accounting for the greatest genomic evidence for production of a given
146    metabolite in a set media condition (Table S4). We represent the information from each constituent
147    network using a single summary metric referred to as the Production Likelihood by calculating the
148    average of all likelihoods of reactions that carry flux. The average of all reaction likelihoods in a
149    metabolic pathway has been previously shown to be a valuable metric for making comparisons between
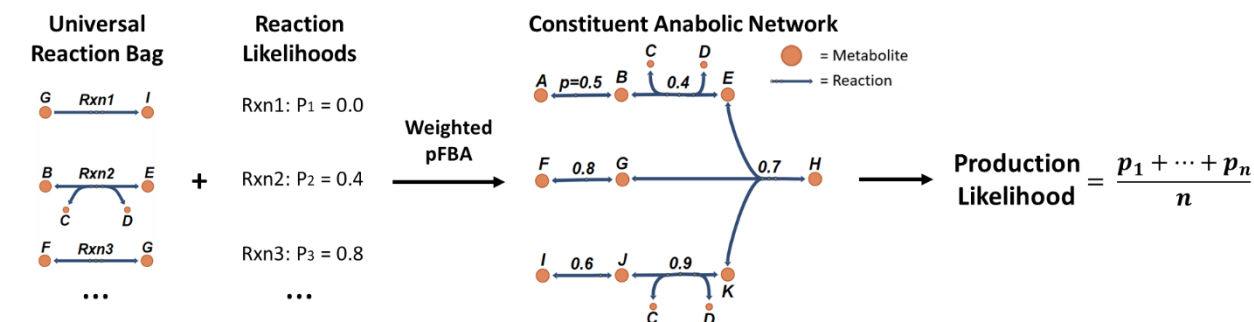150    networks (44).

151



152    **Figure 2: PROTEAN is an approach for quantifying the likelihood that a given metabolic network, derived**
153    **exclusively from genomic evidence, is capable of synthesizing a particular metabolite.** A modified version of
154    Parsimonious Enzyme Usage FBA (weighted pFBA) was performed on a standardized set of reactions to generate
155    constituent anabolic networks for each genome. Reaction likelihoods were used to weight the minimization of flux
156    through each reaction in the network. Therefore, reactions with a greater likelihood were more likely to be
157    included in the resulting constituent anabolic network. Each constituent network has a set of input metabolites
158    representing the media condition (Table S4) and a demand reaction for a certain metabolic product. The resulting
159    constituent network is the set of reactions that requires flux to produce the metabolic product in the given media
160    condition. The production likelihood metric is an average of all the reaction likelihoods associated with the
161    reactions included in the constituent network. This metric is used as a summary statistic that allows for the
162    comparison of constituent networks across different metabolic products and strains, where a higher production
163    likelihood corresponds with greater genetic evidence for that particular constituent anabolic network.

164    **The Scaled Production Likelihood metric facilitates comparison of anabolic capabilities between**
165    **species and strains**

166    Predicted constituent anabolic networks were generated for a set of 50 biologically-relevant metabolic
167    products for each of the 144 *Lactobacillus* strains. The 50 metabolites were selected based on known
168    *Lactobacillus* biology (see Methods). For each metabolic product, we generated a constituent anabolic
169    network (Table S3) across all strains. For each genome we scaled the Production Likelihoods metric by
170    calculating the corresponding z-score. The standard deviation for the z-score calculation was across all
171    metabolic products for each strain. This metric allows for a relative comparison of production
172    capabilities across strains that does not rely on well-curated metabolic network reconstructions. The
173    resulting Scaled Production Likelihood (SPL) is a metric indicating likelihood that a genome encodes for
174    the cellular machinery required to produce a metabolite, given a specific media condition, relative to all
175    of the other SPLs for the metabolic products per strain. For visualization, these data were grouped by
176    species and summarized using the median of the SPLs across all of the strains within each species (Figure
177    3).

178



179

180 **Figure 3: Predicted metabolic production capabilities with the Scaled Production Likelihood (SPL) metric align**
181 **poorly with phylogeny.** There is a single production likelihood for each genome associated with each metabolite. A
182 median SPL can be calculated for a species that allows for more general comparisons across species, illustrated
183 here by the distribution for one species (*L. rhamnosus*) and one metabolite (adenine). There are 50 metabolites
184 used as features to allow for the comparison of predicted production capabilities across the lactobacilli analyzed.

185 The strains were grouped by species and clustered based on median SPLs. We found that across the 16
186 species, D- and L-lactate both have high median SPLs, as we would expect with lactobacilli. Additionally,
187 fumarate and GABA have particularly low SPLs across all species. We were able to find several
188 publications indicating GABA can be produced by select lactobacilli in specific environments (49,50).
189 However, we were unable to find publications discussing the production of fumarate by lactobacilli.
190 Additionally, we found that the dendrogram from clustering based on predicted metabolic production
191 capabilities does not qualitatively align well with published phylogenetic trees generated using the 16S
192 rRNA gene (34). The misalignment to established phylogenetic trees indicates that phylogeny is a poor
193 indicator of metabolic production capabilities. It is likely that evolution of metabolic production
194 capabilities is driven independently from classical genes used for phylogenetic comparisons, such as the
195 16S rRNA gene. Therefore, we need more precise computational tools to better understand the
196 phenotypic differences between microbial species when interrogating metabolism. Perhaps
197 phylogenetic analysis would be augmented with the consideration of metabolic genes in addition to the
198 16S rRNA gene.

199 **Intestinal and oral *Lactobacillus* strains have different metabolic capabilities compared to vaginal**
200 **strains**

201 We performed principle coordinate analysis (PCoA) on the SPLs for each species and determined that
202 the *Lactobacillus* strains cluster significantly by both species (Figure 4A) and isolation site (Figure 4B)
203 (PERMANOVA; $P < 0.001$). The vaginal isolates differ from both the oral and gut cluster (Figure 4B).
204 Substantial overlap was found between oral and gut isolates, specifically within *L. gasseri, L. rhamnosus,*
205 and *L. salivarius*, likely due to the consistent transmission of orally colonized microbes to the intestines
206 (15). It has been hypothesized that many of the lactobacilli isolated from the gut are actually transient
207 strains that are colonized in the oral cavity (51). Our data supports this hypothesis by showing that oral

208 isolates are metabolically similar to a portion of the intestinal isolates. However, there are lactobacilli,
209 such as *L. reuteri*, which likely colonize the human intestines (52). Five of the 16 species in this study are
210 only represented by strains isolated from the intestines; although this result is influenced by sampling
211 bias in the PATRIC Database, it provides support that our data contains species that are only found in
212 the intestines. The vaginal isolates cluster separately from the intestinal/oral isolates along the primary
213 coordinate that accounts for 78% of the variation in these data. The vaginal microbiota is frequently
214 dominated by several *Lactobacillus* species, such as *L. iners*, *L. crispatus*, and *L. jensenii* (53–55). This
215 separation of vaginal isolates from intestinal/oral isolates indicates that these two main clusters have
216 differences in their metabolic production capabilities. This result is to be expected because the
217 intestinal/oral nutrient environment is drastically different from the vaginal environment and the
218 dominant species appear to have metabolic capabilities that reflect this difference.



219

220 **Figure 4: The Scaled Production Likelihood metric distinguishes metabolic functionality among species.** (A) We
221 found that *Lactobacillus* strains cluster significantly by species (PERMANOVA; P < 0.001). (B) Additionally, they
222 cluster significantly by isolation site (PERMANOVA; P < 0.001). Both plots are PCoA using the Bray-Curtis distance
223 metric of the SPLs for each isolate. Points in both panels are identical, but displayed with different color schemes.

224 In addition to distinguishing isolates by body site, the SPL metric is capable of defining collections of
225 functional components that drive differences between groups. Using standard genomic analyses,
226 differences between groups are typically defined by the differential gene content. Genes are intrinsically
227 part of a larger network of metabolism where absence of specific functionality related to a gene may be
228 compensated for within the system. Since our approach is based on Production Likelihoods of specific
229 metabolites, it functions within a more complex metabolic framework compared to the analysis of
230 genomic data without the network context. Using machine learning, we were able to identify the set of
231 metabolites for which each group of strains is more likely to encode the cellular machinery required for
232 production. We conducted a machine learning feature selection to determine the metabolites that are
233 most likely to be produced by each group of strains, intestinal/oral strains and vaginal strains. We
234 grouped the intestinal and oral strains together due to their inherent similarity (Figure 4B) and the
235 observed transmission of oral strains to the intestines (15,51). We generated two separate area under
236 the curve random forest (AUCRF) models to determine the metabolites that were more likely to be
237 produced by each of the groups. Two models were necessary to enrich for the most discriminatory
238 metabolites that were more likely to be produced in each of the groups, rather than simply identifying
239 the metabolites that best classify the samples based on isolation site regardless of being more or less

240   likely to be produced (See methods). The first model was generated to select the metabolites that are
241   most likely to be produced by the intestinal and oral isolates compared to the vaginal isolates, while
242   maximally discriminating the groups. The eight metabolites selected accurately classify greater than 90%
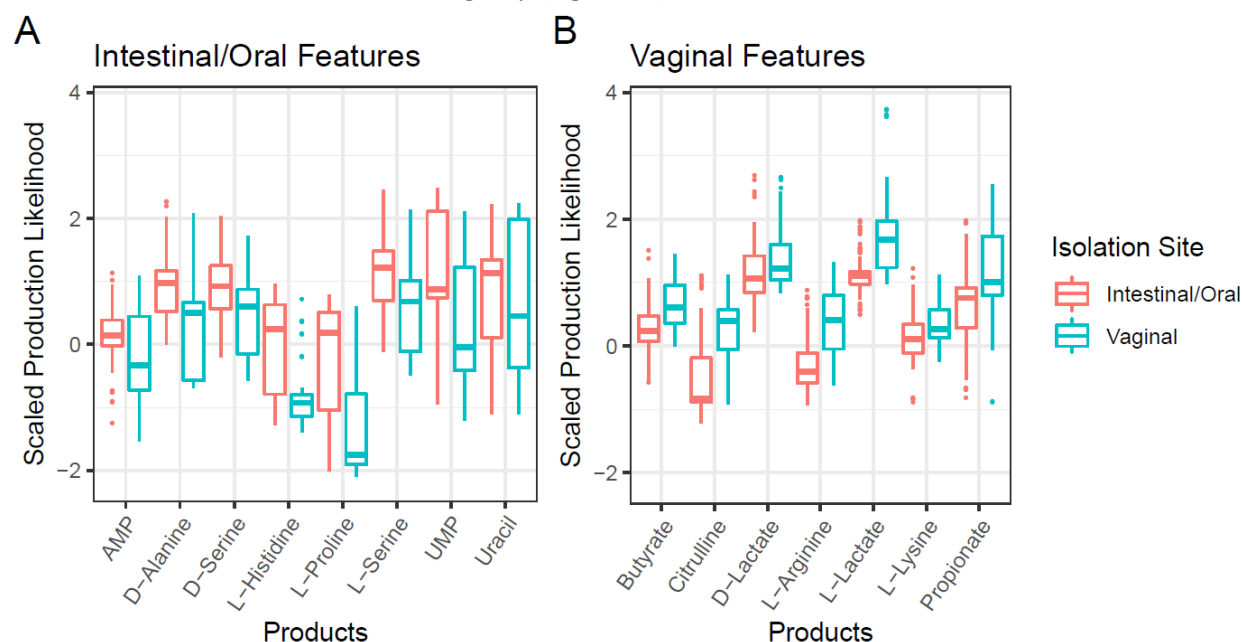243   of isolates to the correct group (Figure 5A). The second model was generated to select the metabolites
244   that are most likely to be produced by the vaginal isolates compared to the intestinal and oral isolates,
245   while maximally discriminating the groups. The seven metabolites selected accurately classify greater
246   than 90% of the isolates to the correct group (Figure 5B).



247

248   **Figure 5: Machine learning of the SPL scores identifies metabolites that discriminate *Lactobacillus* strains.**
249   Machine learning feature selection identified the metabolites that are both most likely to be produced by each
250   group and capable of classifying the strains into two groups, intestinal/oral and vaginal, with greater than 90%
251   accuracy. (A) There are eight metabolites that are more likely to be produced by the intestinal/oral isolates
252   compared to the vaginal isolates. (B) There are seven metabolites that are more likely to be produced by the
253   vaginal isolates compared to intestinal/oral isolates. Both models are more than 90% accurate in predicting the
254   membership to which the given isolate belongs using the SPLs of the metabolites listed.

255   Using SPLs as an input for AUCRF feature selection, we identified the metabolites that are most likely to
256   be produced by the strains associated with the two isolate groups, intestinal/oral and vaginal. The
257   selected metabolite products may contribute to how the strains interact with the mucosal tissues in
258   each site. We hypothesize that these metabolites are related to key phenotypic differences between the
259   two isolate groups. Four of the selected metabolites that are likely produced by intestinal/oral strains,
260   D-alanine, D/L-serine, and L-proline (Figure 5A), have all been previously identified to have an impact on
261   the human intestinal epithelium (23,24,56–58). Additionally, four of the selected metabolites that are
262   likely produced by vaginal strains, L-arginine, citrulline, and D/L-lactate (Figure 5B), have been previously
263   identified to have an impact on the human vaginal microbiome (59–62). The metabolites for which we
264   have not found existing experimental evidence for are likely worth focusing on in future experimental
265   studies.

266 For intestine-associated lactobacilli in this study, there is a connection between intestinal immune
267 system regulation and D-alanine rich lipotechoic acid, a glycolipid expressed by some lactobacilli, such as
268 *L. plantarum* (23,24). D-alanine rich lipotechoic acid, produced by lactobacilli, has been shown to down-
269 regulate local colonic inflammation in a murine colitis model (23,24). With PROTEAN we identified that
270 intestinal lactobacilli were more likely to produce D-alanine (Figure 5A). It is possible that a positive
271 interaction with the intestinal host immune system would result in an evolutionary advantage by
272 reducing local immune response. Additionally, serine rich serine-threonine peptides have been shown to
273 have a similar regulatory effect on intestinal dendritic cells (56,57). These peptides expressed by *L.*
274 *plantarum* are resistant to intestinal proteolysis and appear to be present in the colon of most healthy
275 individuals (56,57). Similar to D-alanine, the production of D/L-serine would require a robust
276 biosynthesis pathway present in those strains.

277 A final gut-related connection involves the biosynthesis of L-proline (Figure 5A). One of the primary
278 stress responses in *L. acidophilus* to high osmotic pressure results in the accumulation of L-proline in the
279 cell; there is little evidence that this response is a result of L-proline transport into the cell (58). These
280 *Lactobacillus* strains are exposed to a large range of stressors in the gut, including suboptimal osmotic
281 pressures. There is strong evidence that L-proline is used by *L. acidophilus* to tolerate suboptimal
282 osmotic pressures and there is a lack of evidence for L-proline transporters. As such, the biosynthesis of
283 L-proline may be advantageous for growth in the gut.

284 For the enriched metabolic products in vaginal isolates (Figure 5B), there is evidence for an
285 arginine/ornithine antiporter and arginine deiminase in *L. fermentum* (59). These enzymes are part of
286 the arginine deiminase pathway through which there is the production of citrulline which is exported
287 from the cell and contributes to acid tolerance (59). It has also been demonstrated that treatment with
288 probiotics containing arginine deiminase-positive lactobacilli can improve clinical symptoms of vaginosis
289 in parallel with significant declines in polyamine (i.e. arginine, ornithine, and citrulline) levels in the
290 vagina (60,61). The vaginal isolates in this study show enrichment for the cellular machinery required for
291 the production of both citrulline and L-arginine (Figure 5B). The importance of lactate for the adequate
292 maintenance of vaginal health in many individuals is known. The current hypothesis revolves around
293 colonization resistance where vaginal lactobacilli establish an acidic environment by producing lactate
294 (62). The acidic environment is generally inhospitable to invading pathogens as well as other microbes
295 that are otherwise capable of residing in the vaginal environment (62). It has been shown that higher
296 levels of D-lactate over L-lactate present in the vagina, produced by lactobacilli, further decrease the
297 chance of infections in female patients (62). However, both isoforms of lactate remain important in
298 maintaining vaginal health.

## Conclusions

300 Microbial biosynthesis of metabolites has a broad range of applications, from bio-manufacturing to
301 microbiome research (63). There is a wealth of well-curated and accessible knowledge stored in
302 biochemical reaction databases such as ModelSEED (64). Genome-Scale Metabolic Network
303 Reconstructions access this fundamental knowledge while accounting for systems-level interactions.
304 This study represents one such application of GENREs that is a step toward predicting the metabolic
305 production capabilities of understudied organisms. Experimental validation of the production
306 capabilities predicted with PROTEAN will allow for conclusions to be made beyond the statement that a
307 microbe is genetically likely to be able to produce a metabolite. Utilizing PROTEAN data, we found that

308    human-associated lactobacilli strains cluster significantly by species and isolation site. Additionally,
309    many of the metabolic products that drive the clustering of strains by the isolation sites have known
310    physiological function and importance in the respective isolation sites.

311    Future applications of PROTEAN could include optimal strain selection for bio-manufacturing of a certain
312    compound, generating predicted metabolomics data for an organism to generate a prioritized list of
313    conditions that would be most worthwhile to validate experimentally, and predicting the metabolites
314    that are most likely to be produced in a microbiota. Microbes can have a wide range of physiological
315    impacts on human health; these impacts are, in part, a result of the metabolites that are or are not
316    produced by members of a microbiota. One of the core limitations of this study includes the lack of
317    reaction likelihoods for some reactions in the universal reaction bag we used from ModelSEED. The
318    number of reactions we could generate likelihoods for was limited by the Probannopy reaction
319    template. However, this template can be expanded to continue to improve the utility of PROTEAN. With
320    the inclusion of validation data, additional analyses will be possible, such as determining metabolic
321    production pathways lacking proper annotation. By determining the reactions that are most likely
322    required for biosynthesis of a known product, it would be possible to generate additional hypotheses for
323    enzyme annotation experiments. PROTEAN is an algorithm with potential for a wide range of
324    applications in the study and use of microbial metabolic networks.

# Methods

**Constituent Anabolic Network Generation (PROTEAN)**

327    Probabilistic pFBA-based constituent anabolic network generation was accomplished using three Python
328    packages, Cobrapy (65), Mackinac (66), and Probannopy (48). The complete ModelSEED universal
329    reaction bag was downloaded from the github repository and filtered based on the annotation quality
330    score, including all reactions with an 'OK' quality status or better (64). For each reaction in the
331    ModelSEED universal reaction bag, we used Probannopy to generate a reaction likelihood based on the
332    FASTA file for each genome obtained from the PATRIC database (4). The Cobrapy implementation of
333    Parsimonious Enzyme Usage Flux Balance Analysis (pFBA) was altered to allow for each reaction's linear
334    constraint to be set individually based on the reaction likelihood. The linear constraint for each reaction
335    was set to one minus the reaction likelihood (a value between 0 and 1). There were reactions included in
336    the universal reaction bag that were lacking from the Probannopy template model, therefore resulting
337    in several gene-associated reactions lacking reaction likelihood scores. The reactions without likelihoods
338    were left at a full minimization penalty (linear constraint value of 1). We chose to penalize the reactions
339    without likelihoods to bias our results towards the construction of networks for which all reactions had
340    evidence of presence. The linear constraints applied to each reaction based on likelihood acted as a
341    weighting (inclusion penalty) for the minimization step in pFBA, resulting in the reactions with greater
342    likelihood having a lower penalty for carrying flux; therefore, the reactions had a higher likelihood of
343    being included in the constituent anabolic networks.

344    Using PROTEAN, we generated constituent anabolic networks by setting a certain input media condition
345    (Table S4) and constraining flux through the single metabolite objective function (Table S3). We ran our
346    likelihood-weighted pFBA flux minimization across the entire universal reaction bag and isolated the
347    reactions that carried flux to get the desired product. The resulting networks consist of the direct
348    reactions that would be part of a production pathway as might be shown in a typical biosynthesis
349    pathway figure, while also accounting for all of the secondary and energy metabolites that are required

350    for the production of the metabolite in consideration. Additionally, this algorithm is optimizing for three
351    core characteristics in the constituent networks: 1) minimum flux through the network (loosely, the
352    minimum number of reactions), 2) maximum average reaction likelihood across the constituent
353    network, and 3) output flux within 90% of the optimal yield of the metabolic product. We chose to allow
354    flux through any reaction in the universal reaction bag during the generation of the constituent anabolic
355    production pathways rather than simply pulling from a GENRE that was first gapfilled to allow
356    production of biomass. Using the universal reaction bag instead of a gapfilled model was important
357    because the biomass function is difficult to define for understudied organisms and unnecessary for our
358    applications.

**Scaled Production Likelihood Metric**

360    We represent the information from each constituent network using a single summary metric for ease of
361    comparison, simply named the Production Likelihood. This metric is the average of the reaction
362    likelihoods included in the constituent network. The average reaction likelihood for a metabolic pathway
363    has been previously used for making comparisons between networks (44). The Production Likelihoods
364    for all 50 metabolites are scaled for each given genome by calculating the z-score to create the Scaled
365    Production Likelihoods used for the majority of the analysis in this study. The z-score is calculated for
366    each individual strain using the median and standard deviation for the production likelihoods across the
367    50 metabolic products. The Scaled Production Likelihood allows for a ranked comparison of metabolic
368    products across the genome set and corrects for annotation bias by essentially comparing the ranked z-
369    score for each metabolic product.

**Supporting data for pathway generation**

371    The simulated media formulation was based on *in vitro* minimal media growth conditions for *L.
372    plantarum* (Table S4) (67–69). The techniques used in this study do not assume that all species are
373    capable of growth in the given media condition, therefore this media condition simply provides a
374    standard reference for comparison. The product list was developed by identifying metabolites that have
375    been shown to be produced by lactobacilli during *in vitro* growth experiments, in addition to other
376    metabolites that have been shown to be related to human physiology (70–74).

**Machine learning feature selection**

378    Discriminating intestinal/oral and vaginal features were selected using area under the ROC curve
379    random forest (AUCRF) using default parameters (75) (see Code). We generated two separate AUCRF
380    models to determine the metabolites that were more likely to be produced by each of the groups,
381    intestinal/oral and vaginal. Two models allowed us to enrich for likely products rather than simply
382    selecting for the metabolites that provide the greatest discrimination between the groups but which
383    may have poor likelihood scores. We conducted the enrichment for likely metabolic products for each
384    model by reducing the feature set down to only metabolites that were more likely to be produced by
385    the group of interest. Likely metabolic products were determined by comparing the median SPLs of each
386    metabolite between the groups. Additionally, the feature sets were reduced to include only metabolites
387    with a median value greater than zero for the group of interest. An AUCRF model was then generated to
388    select the features that provided the greatest discrimination between the two groups.

**Statistical modeling and figure generation**

390 The principle coordinate analysis (PCoA) ordinations were created using the R vegan package (76),
391 implemented with the Bray-Curtis dissimilarity metric. Statistical significance for comparing the PCoA
392 clusters was determined using a PERMANOVA (R Adonis test). A variety of R packages were used for all
393 figure generation (77–81).

394 **Genome Quality and PATRIC Cross Genus Protein Family Data**

395 Genomes used in the study were filtered for quality before being included in the analysis. Strains with
396 greater than 0.2% unknown nucleotide calls in the genome were eliminated. Low quality genome
397 assemblies with greater than 300 contigs were removed. Non-human associated *Lactobacillus* strains
398 from the PATRIC database were used to determine the GC content range for each species (82,83), and
399 significant outliers (plus or minus two percent) were removed to control for sequencing bias (84,85).
400 Only isolates from the three human-associated sites (oral, intestinal, and vaginal) were included in the
401 final dataset.

402 The inclusion of metabolic PATRIC cross genus protein families was conducted by filtering the PGfams
403 for each genome based on the existence of an associated known reaction and Probannopy likelihood
404 greater than 0. Pan and core metabolic PGfam sets were evaluated after the addition of all genomic
405 features from each genome. The pan set of metabolic PGfams was defined as the total number of
406 unique PGfams included in the data set after the above filtering steps. The core set of metabolic PGfams
407 are those that existed within each genome included in this study.

408 **Data and code availability**

409 Genome FASTA files and metadata were downloaded from the PATRIC Database (4). Python and R code
410 is available at: Github.com/Tjmoutinho/Lactobacillus

# 411 References

412 1. de Vos WM. Systems solutions by lactic acid bacteria: from paradigms to practice. Microb Cell
413    Factories. 2011 Aug 30;10(1):S2.

414 2. de Vos WM, Hugenholtz J. Engineering metabolic highways in Lactococci and other lactic acid
415    bacteria. Trends Biotechnol. 2004 Feb 1;22(2):72–9.

416 3. Ljungh Å, Wadström T. Lactobacillus Molecular Biology: From Genomics to Probiotics. Horizon
417    Scientific Press; 2009. 217 p.

418 4. Wattam AR, Abraham D, Dalay O, Disz TL, Driscoll T, Gabbard JL, et al. PATRIC, the bacterial
419    bioinformatics database and analysis resource. Nucleic Acids Res. 2014 Jan 1;42(D1):D581–91.

420 5. OHanlon DE. In vivo versus in vitro metabolomics profiling of vaginal lactobacilli for probiotic use.
421    2013 Jun 4 [cited 2018 Sep 24]; Available from: https://www.omicsonline.org/proceedings/in-vivo-
422    versus-in-vitro-metabolomics-profiling-of-vaginal-lactobacilli-for-probiotic-use-785.html

423 6. O'Hanlon DE, Moench TR, Cone RA. Vaginal pH and Microbicidal Lactic Acid When Lactobacilli
424    Dominate the Microbiota. PLoS ONE [Internet]. 2013 Nov 6 [cited 2018 Sep 24];8(11). Available
425    from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3819307/

426    7.    Tachedjian G, Aldunate M, Bradshaw CS, Cone RA. The role of lactic acid production by probiotic
427          Lactobacillus species in vaginal health. Res Microbiol. 2017 Nov 1;168(9):782–92.

428    8.    Tachedjian G, O'Hanlon DE, Ravel J. The implausible "in vivo" role of hydrogen peroxide as an
429          antimicrobial factor produced by vaginal microbiota. Microbiome [Internet]. 2018 Feb 6 [cited 2018
430          Sep 24];6. Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5801833/

431    9.    Parolin C, Foschi C, Laghi L, Zhu C, Banzola N, Gaspari V, et al. Insights Into Vaginal Bacterial
432          Communities and Metabolic Profiles of Chlamydia trachomatis Infection: Positioning Between
433          Eubiosis and Dysbiosis. Front Microbiol [Internet]. 2018 Mar 28 [cited 2018 Sep 24];9. Available
434          from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5883401/

435    10.   Vitali B, Cruciani F, Picone G, Parolin C, Donders G, Laghi L. Vaginal microbiome and metabolome
436          highlight specific signatures of bacterial vaginosis. Eur J Clin Microbiol Infect Dis. 2015 Dec
437          1;34(12):2367–76.

438    11.   Gosmann C, Anahtar MN, Handley SA, Farcasanu M, Abu-Ali G, Bowman BA, et al. Lactobacillus-
439          Deficient Cervicovaginal Bacterial Communities Are Associated with Increased HIV Acquisition in
440          Young South African Women. Immunity. 2017 Jan 17;46(1):29–37.

441    12.   Ratzke C, Gore J. Modifying and reacting to the environmental pH can drive bacterial interactions.
442          PLOS Biol. 2018 Mar 14;16(3):e2004248.

443    13.   Palmer RJ. Composition and development of oral bacterial communities. Periodontol 2000
444          [Internet]. 2014 Feb [cited 2018 Sep 24];64(1). Available from:
445          https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3876289/

446    14.   Tannock GW. A Special Fondness for Lactobacilli. Appl Environ Microbiol. 2004 Jun;70(6):3189–94.

447    15.   Schmidt TSB, Hayward MR, Coelho LP, Li SS, Costea PI, Voigt AY, et al. Extensive transmission of
448          microbes along the gastrointestinal tract. Nieuwdorp M, editor. eLife. 2019 Feb 12;8:e42693.

449    16.   Szajewska H, Ruszczyński M, Radzikowski A. Probiotics in the prevention of antibiotic-associated
450          diarrhea in children: A meta-analysis of randomized controlled trials. J Pediatr. 2006 Sep
451          1;149(3):367-372.e1.

452    17.   Hempel S, Newberry SJ, Maher AR, Wang Z, Miles JNV, Shanman R, et al. Probiotics for the
453          Prevention and Treatment of Antibiotic-Associated Diarrhea: A Systematic Review and Meta-
454          analysis. JAMA. 2012 May 9;307(18):1959–69.

455    18.   Ford AC, Quigley EMM, Lacy BE, Lembo AJ, Saito YA, Schiller LR, et al. Efficacy of Prebiotics,
456          Probiotics, and Synbiotics in Irritable Bowel Syndrome and Chronic Idiopathic Constipation:
457          Systematic Review and Meta-analysis. Am J Gastroenterol. 2014 Oct;109(10):1547–61.

458    19.   Nikfar S, Rahimi R, Rahimi F, Derakhshani S, Abdollahi M. Efficacy of Probiotics in Irritable Bowel
459          Syndrome: A Meta-Analysis of Randomized, Controlled Trials. Dis Colon Rectum. 2008 Dec
460          1;51(12):1775–80.

461   20. Elazab N, Mendy A, Gasana J, Vieira ER, Quizon A, Forno E. Probiotic Administration in Early Life,
462       Atopy, and Asthma: A Meta-analysis of Clinical Trials. Pediatrics. 2013 Sep 1;132(3):e666–76.

463   21. Berstad A, Raa J, Midtvedt T, Valeur J. Probiotic lactic acid bacteria – the fledgling cuckoos of the
464       gut? Microb Ecol Health Dis [Internet]. 2016 May 26 [cited 2018 Sep 24];27. Available from:
465       https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4884264/

466   22. Suez J, Zmora N, Zilberman-Schapira G, Mor U, Dori-Bachash M, Bashiardes S, et al. Post-Antibiotic
467       Gut Mucosal Microbiome Reconstitution Is Impaired by Probiotics and Improved by Autologous
468       FMT. Cell. 2018 Sep 6;174(6):1406-1423.e16.

469   23. de Vos WM. Lipotechoic acid in lactobacilli: D-alanine makes the difference. Proc Natl Acad Sci.
470       2005;102(31):10763–4.

471   24. Grangette C, Nutten S, Palumbo E, Morath S, Hermann C, Dewulf J, et al. Enhanced
472       antiinflammatory capacity of a Lactobacillus plantarum mutant synthesizing modified teichoic acids.
473       Proc Natl Acad Sci. 2005;102(29):10321–6.

474   25. Branco dos Santos F, de Vos WM, Teusink B. Towards metagenome-scale models for industrial
475       applications—the case of Lactic Acid Bacteria. Curr Opin Biotechnol. 2013 Apr 1;24(2):200–6.

476   26. Le Barz M, Anhê FF, Varin TV, Desjardins Y, Levy E, Roy D, et al. Probiotics as Complementary
477       Treatment for Metabolic Disorders. Diabetes Metab J. 2015 Aug;39(4):291–303.

478   27. Saulnier DM, Santos F, Roos S, Mistretta T-A, Spinler JK, Molenaar D, et al. Exploring Metabolic
479       Pathway Reconstruction and Genome-Wide Expression Profiling in Lactobacillus reuteri to Define
480       Functional Probiotic Features. PLOS ONE. 2011 Apr 29;6(4):e18783.

481   28. Lewis NE, Nagarajan H, Palsson BO. Constraining the metabolic genotype–phenotype relationship
482       using a phylogeny of in silico methods. Nat Rev Microbiol. 2012 Apr;10(4):291–305.

483   29. Haggart CR, Bartell JA, Saucerman JJ, Papin JA. Whole-genome metabolic network reconstruction
484       and constraint-based modeling. Methods Enzymol. 2011;500:411–33.

485   30. Kant R, Blom J, Palva A, Siezen RJ, de Vos WM. Comparative genomics of Lactobacillus. Microb
486       Biotechnol. 2011 May;4(3):323–32.

487   31. Drissi F, Merhej V, Angelakis E, El Kaoutari A, Carrière F, Henrissat B, et al. Comparative genomics
488       analysis of Lactobacillus species associated with weight gain or weight protection. Nutr Diabetes.
489       2014 Feb;4(2):e109.

490   32. France MT, Mendes-Soares H, Forney LJ. Genomic Comparisons of Lactobacillus crispatus and
491       Lactobacillus iners Reveal Potential Ecological Drivers of Community Composition in the Vagina.
492       Appl Env Microbiol. 2016 Dec 15;82(24):7063–73.

493   33. Morita H, Toh H, Fukuda S, Horikawa H, Oshima K, Suzuki T, et al. Comparative Genome Analysis of
494       Lactobacillus reuteri and Lactobacillus fermentum Reveal a Genomic Island for Reuterin and
495       Cobalamin Production. DNA Res. 2008 Jun 1;15(3):151–61.

34. Zhang Z-G, Ye Z-Q, Yu L, Shi P. Phylogenomic reconstruction of lactic acid bacteria: an update. BMC Evol Biol. 2011 Jan 1;11:1.

35. Kleerebezem M, Vos WM de. Lactic acid bacteria: life after genomics. Microb Biotechnol. 2011 May 1;4(3):318–22.

36. Rau MH, Zeidan AA. Constraint-based modeling in microbial food biotechnology. Biochem Soc Trans. 2018 Mar 27;BST20170268.

37. Lewis NE, Hixson KK, Conrad TM, Lerman JA, Charusanti P, Polpitiya AD, et al. Omic data from evolved E. coli are consistent with computed optimal growth from genome-scale models. Mol Syst Biol. 2010 Jul 27;6:390.

38. Feist AM, Palsson BO. The biomass objective function. Curr Opin Microbiol. 2010 Jun;13(3):344–9.

39. Altafini C, Facchetti G. Metabolic Adaptation Processes That Converge to Optimal Biomass Flux Distributions. PLoS Comput Biol. 2015 Sep 4;11(9):e1004434.

40. Orth JD, Thiele I, Palsson BØ. What is flux balance analysis? Nat Biotechnol. 2010 Mar;28(3):245–8.

41. Pinto F, Medina DA, Pérez-Correa JR, Garrido D. Modeling metabolic interactions in a consortium of the infant gut microbiome. Front Microbiol. 2017;8:2507.

42. Schmidt BJ, Ebrahim A, Metz TO, Adkins JN, Palsson BØ, Hyduke DR. GIM3E: condition-specific models of cellular metabolism developed from metabolomics and expression data. Bioinformatics. 2013 Nov 15;29(22):2900–8.

43. Li H, Zhu J. Targeted metabolic profiling rapidly differentiates Escherichia coli and Staphylococcus aureus at species and strain level. Rapid Commun Mass Spectrom. 2017;31(19):1669–76.

44. Benedict MN, Mundy MB, Henry CS, Chia N, Price ND. Likelihood-Based Gene Annotations for Gap Filling and Quality Assessment in Genome-Scale Metabolic Models. PLOS Comput Biol. 2014 Oct 16;10(10):e1003882.

45. Thiele I, Vlassis N, Fleming RMT. fastGapFill: efficient gap filling in metabolic networks. Bioinformatics. 2014 Sep 1;30(17):2529–31.

46. Machado D, Andrejev S, Tramontano M, Patil KR. Fast automated reconstruction of genome-scale metabolic models for microbial species and communities. Nucleic Acids Res. 2018 Sep 6;46(15):7542–53.

47. Devoid S, Overbeek R, DeJongh M, Vonstein V, Best AaronA, Henry C. Automated Genome Annotation and Metabolic Model Reconstruction in the SEED and Model SEED. In: Alper HS, editor. Systems Metabolic Engineering [Internet]. Humana Press; 2013 [cited 2017 Apr 6]. p. 17–45. (Methods in Molecular Biology). Available from: http://dx.doi.org/10.1007/978-1-62703-299-5_2

48. King B, Farrah T, Richards MA, Mundy M, Simeonidis E, Price ND. ProbAnnoWeb and ProbAnnoPy: probabilistic annotation and gap-filling of metabolic reconstructions. Bioinformatics. 2018 May 1;34(9):1594–6.

531   49. Komatsuzaki N, Shima J, Kawamoto S, Momose H, Kimura T. Production of γ-aminobutyric acid
532       (GABA) by Lactobacillus paracasei isolated from traditional fermented foods. Food Microbiol.
533       2005;22(6):497–504.

534   50. Li H, Qiu T, Huang G, Cao Y. Production of gamma-aminobutyric acid by Lactobacillus brevis NCL912
535       using fed-batch fermentation. Microb Cell Factories. 2010;9(1):85.

536   51. Walter J. Ecological Role of Lactobacilli in the Gastrointestinal Tract: Implications for Fundamental
537       and Biomedical Research. Appl Env Microbiol. 2008 Aug 15;74(16):4985–96.

538   52. Valeur N, Engel P, Carbajal N, Connolly E, Ladefoged K. Colonization and immunomodulation by
539       Lactobacillus reuteri ATCC 55730 in the human gastrointestinal tract. Appl Environ Microbiol.
540       2004;70(2):1176–81.

541   53. Romero R, Hassan SS, Gajer P, Tarca AL, Fadrosh DW, Nikita L, et al. The composition and stability of
542       the vaginal microbiota of normal pregnant women is different from that of non-pregnant women.
543       Microbiome. 2014 Feb 3;2(1):4.

544   54. Gajer P, Brotman RM, Bai G, Sakamoto J, Schütte UME, Zhong X, et al. Temporal Dynamics of the
545       Human Vaginal Microbiota. Sci Transl Med. 2012 May 2;4(132):132ra52-132ra52.

546   55. Ravel J, Gajer P, Abdo Z, Schneider GM, Koenig SSK, McCulle SL, et al. Vaginal microbiome of
547       reproductive-age women. Proc Natl Acad Sci. 2011 Mar 15;108(Supplement 1):4680–7.

548   56. Al-Hassi HO, Mann ER, Sanchez B, English NR, Peake STC, Landy J, et al. Altered human gut dendritic
549       cell properties in ulcerative colitis are reversed by Lactobacillus plantarum extracellular encrypted
550       peptide STp. Mol Nutr Food Res. 2014;58(5):1132–43.

551   57. Bernardo D, Sánchez B, Al-Hassi HO, Mann ER, Urdaci MC, Knight SC, et al. Microbiota/Host
552       Crosstalk Biomarkers: Regulatory Response of Human Intestinal Dendritic Cells Exposed to
553       Lactobacillus Extracellular Encrypted Peptide. PLOS ONE. 2012 May 14;7(5):e36262.

554   58. Jewell JB, Kashket ER. Osmotically regulated transport of proline by Lactobacillus acidophilus IFO
555       3532. Appl Env Microbiol. 1991 Oct 1;57(10):2829–33.

556   59. Vrancken G, Rimaux T, Weckx S, De Vuyst L, Leroy F. Environmental pH determines citrulline and
557       ornithine release through the arginine deiminase pathway in Lactobacillus fermentum IMDO
558       130101. Int J Food Microbiol. 2009 Nov 15;135(3):216–22.

559   60. Famularo G, Pieluigi M, Coccia R, Mastroiacovo P, Simone CD. Microecology, bacterial vaginosis and
560       probiotics: perspectives for bacteriotherapy. Med Hypotheses. 2001 Apr 1;56(4):421–30.

561   61. Rousseau V, Lepargneur JP, Roques C, Remaud-Simeon M, Paul F. Prebiotic effects of
562       oligosaccharides on selected vaginal lactobacilli and pathogenic microorganisms. Anaerobe. 2005
563       Jun 1;11(3):145–53.

564   62. Witkin SS, Mendes-Soares H, Linhares IM, Jayaram A, Ledger WJ, Forney LJ. Influence of Vaginal
565       Bacteria and d- and l-Lactic Acid Isomers on Vaginal Extracellular Matrix Metalloproteinase Inducer:

566     Implications for Protection against Upper Genital Tract Infections. mBio. 2013 Aug 30;4(4):e00460-
567     13.

568   63. LeBlanc JG, Milani C, de Giori GS, Sesma F, van Sinderen D, Ventura M. Bacteria as vitamin suppliers
569     to their host: a gut microbiota perspective. Curr Opin Biotechnol. 2013 Apr 1;24(2):160–8.

570   64. Henry CS, DeJongh M, Best AA, Frybarger PM, Linsay B, Stevens RL. High-throughput generation,
571     optimization and analysis of genome-scale metabolic models. Nat Biotechnol. 2010 Sep;28(9):977–
572     82.

573   65. Ebrahim A, Lerman JA, Palsson BO, Hyduke DR. COBRApy: COnstraints-Based Reconstruction and
574     Analysis for Python. BMC Syst Biol. 2013 Aug 8;7(1):74.

575   66. Mundy M, Mendes-Soares H, Chia N. Mackinac: a bridge between ModelSEED and COBRApy to
576     generate and analyze genome-scale metabolic models. Bioinformatics. 2017 Aug 1;33(15):2416–8.

577   67. Wegkamp A, Teusink B, De Vos W m., Smid E j. Development of a minimal growth medium for
578     Lactobacillus plantarum. Lett Appl Microbiol. 2010 Jan 1;50(1):57–64.

579   68. Ricciardi A, Ianniello RG, Parente E, Zotta T. Modified chemically defined medium for enhanced
580     respiratory growth of Lactobacillus casei and Lactobacillus plantarum groups. J Appl Microbiol. 2015
581     Sep 1;119(3):776–85.

582   69. Elli M, Zink R, Rytz A, Reniero R, Morelli L. Iron requirement of Lactobacillus spp. in completely
583     chemically defined growth media. J Appl Microbiol. 2000 Apr 1;88(4):695–703.

584   70. Rowland I, Gibson G, Heinken A, Scott K, Swann J, Thiele I, et al. Gut microbiota functions:
585     metabolism of nutrients and other food components. Eur J Nutr. 2018 Feb 1;57(1):1–24.

586   71. Neis EPJG, Dejong CHC, Rensen SS. The Role of Microbial Amino Acid Metabolism in Host
587     Metabolism. Nutrients. 2015 Apr 16;7(4):2930–46.

588   72. Wu G. Intestinal Mucosal Amino Acid Catabolism. J Nutr. 1998 Aug 1;128(8):1249–52.

589   73. Rooj AK, Kimura Y, Buddington RK. Metabolites produced by probiotic Lactobacilli rapidly increase
590     glucose uptake by Caco-2 cells. BMC Microbiol. 2010 Jan 20;10(1):16.

591   74. Belzer C, Chia LW, Aalvink S, Chamlagain B, Piironen V, Knol J, et al. Microbial Metabolic Networks at
592     the Mucus Layer Lead to Diet-Independent Butyrate and Vitamin B12 Production by Intestinal
593     Symbionts. mBio. 2017 Nov 8;8(5):e00770-17.

594   75. Urrea V, Calle M. AUCRF: variable selection with random forest and the area under the curve. R
595     Package Version 11. 2012;

596   76. Oksanen J, Blanchet FG, Kindt R, Legendre P, Minchin PR, O'hara R, et al. vegan: Community ecology
597     package. R Package Version. 2011;117–8.

598   77. Ihaka R, Gentleman R. R: A Language for Data Analysis and Graphics. J Comput Graph Stat. 1996 Sep
599     1;5(3):299–314.

600    78. Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer; 2016. 266 p.

601    79. Wickham H. tidyr: Easily Tidy Data with spread () and gather () Functions. Version 06 0. 2016;

602    80. Wickham H, Francois R, Henry L, Müller K. dplyr: A grammar of data manipulation. R Package
603         Version 04. 2015;3.

604    81. Neuwirth E, Brewer RC. ColorBrewer palettes. R Package Version. 2014;1–1.

605    82. Haywood-Farmer E, Otto SP. The Evolution of Genomic Base Composition in Bacteria. Evolution.
606         2003;57(8):1783–92.

607    83. Bentley SD, Parkhill J. Comparative Genomic Structure of Prokaryotes. Annu Rev Genet.
608         2004;38(1):771–91.

609    84. Benjamini Y, Speed TP. Summarizing and correcting the GC content bias in high-throughput
610         sequencing. Nucleic Acids Res. 2012 May;40(10):e72.

611    85. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies.
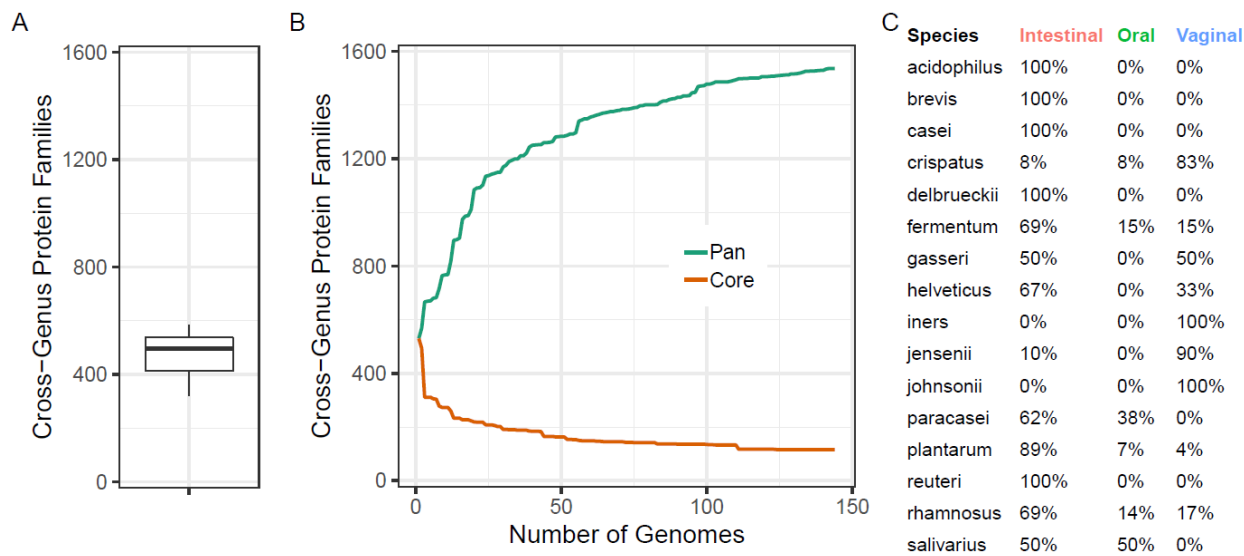612         Bioinformatics. 2013 Apr 15;29(8):1072–5.

613

614

| Species | Intestinal | Oral | Vaginal |
|---|---|---|---|
| acidophilus | 100% | 0% | 0% |
| brevis | 100% | 0% | 0% |
| casei | 100% | 0% | 0% |
| crispatus | 8% | 8% | 83% |
| delbrueckii | 100% | 0% | 0% |
| fermentum | 69% | 15% | 15% |
| gasseri | 50% | 0% | 50% |
| helveticus | 67% | 0% | 33% |
| iners | 0% | 0% | 100% |
| jensenii | 10% | 0% | 90% |
| johnsonii | 0% | 0% | 100% |
| paracasei | 62% | 38% | 0% |
| plantarum | 89% | 7% | 4% |
| reuteri | 100% | 0% | 0% |
| rhamnosus | 69% | 14% | 17% |
| salivarius | 50% | 50% | 0% |

**Figure 1: Known metabolic annotations are extensively sampled across the 16 *Lactobacillus* species included in this study.** The genomic features used for this analysis are PATRIC Cross-Genera Protein families (PGfams), a standardized set of features across the PATRIC Database (4). (A) The number of metabolic PGfams for each genome are shown here, with the median value indicated by the middle line in the boxplot. (B) For the 144 strains from 16 species of *Lactobacillus*, we found that there are 116 protein families in the core set of metabolic PGfams, while the pan set of PGfams expands to over 1500 families. The nearly plateau shape of the curve for the pan set of PGfams curve indicates that this sampling represents a large portion of the genetic diversity among the 16 species included in the study. (C) This table shows the complete list of species used in this study and indicates the percentage of strains that were isolated from each human body site. Each strain in this study is a member from one of the 16 species and isolated from one of three human-associated body sites; intestinal, oral, or vaginal (Table S2).
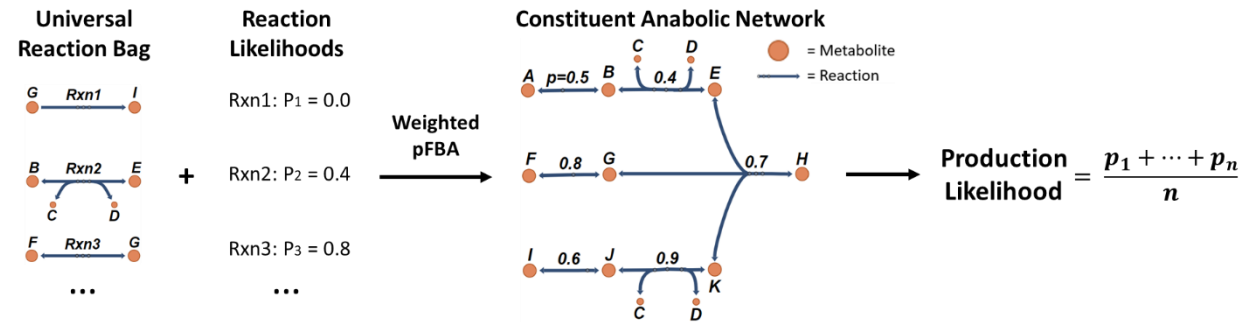


**Figure 2: PROTEAN is an approach for quantifying the likelihood that a given metabolic network, derived exclusively from genomic evidence, is capable of synthesizing a particular metabolite.** A modified version of Parsimonious Enzyme Usage FBA (weighted pFBA) was performed on a standardized set of reactions to generate constituent anabolic networks for each genome. Reaction likelihoods were used to weight the minimization of flux through each reaction in the network. Therefore, reactions with a greater likelihood were more likely to be included in the resulting constituent anabolic network. Each constituent network has a set of input metabolites representing the media condition (Table S4) and a demand reaction for a certain metabolic product. The resulting constituent network is the set of reactions that requires flux to produce the metabolic product in the given media condition. The production likelihood metric is an average of all the reaction likelihoods associated with the reactions included in the constituent network. This metric is used as a summary statistic that allows for the

comparison of constituent networks across different metabolic products and strains, where a higher production likelihood corresponds with greater genetic evidence for that particular constituent anabolic network.
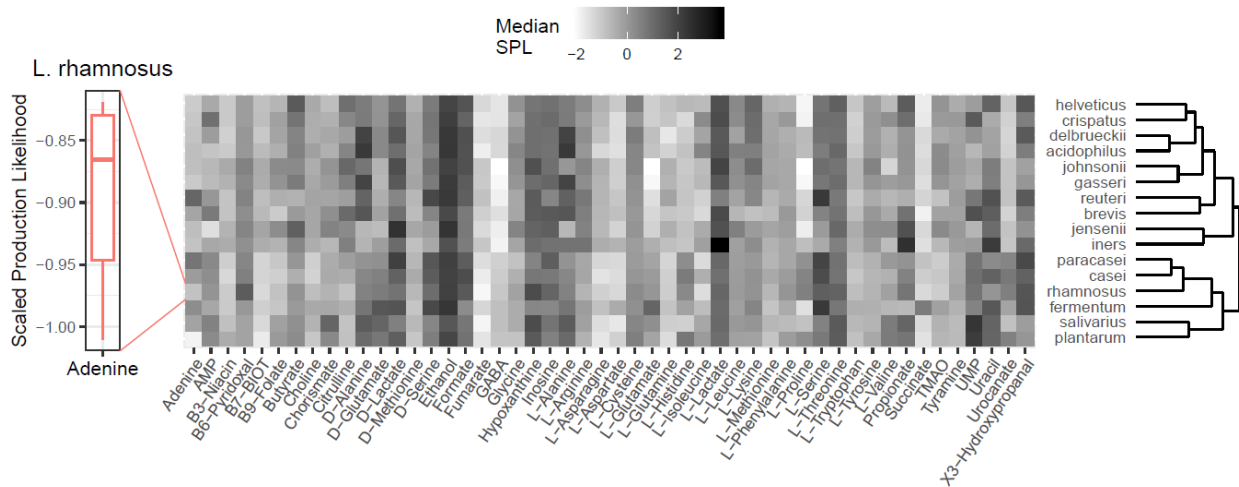


**Figure 3: Predicted metabolic production capabilities with the Scaled Production Likelihood (SPL) metric align poorly with phylogeny.** There is a single production likelihood for each genome associated with each metabolite. A median SPL can be calculated for a species that allows for more general comparisons across species, illustrated here by the distribution for one species (*L. rhamnosus*) and one metabolite (adenine). There are 50 metabolites used as features to allow for the comparison of predicted production capabilities across the lactobacilli analyzed.
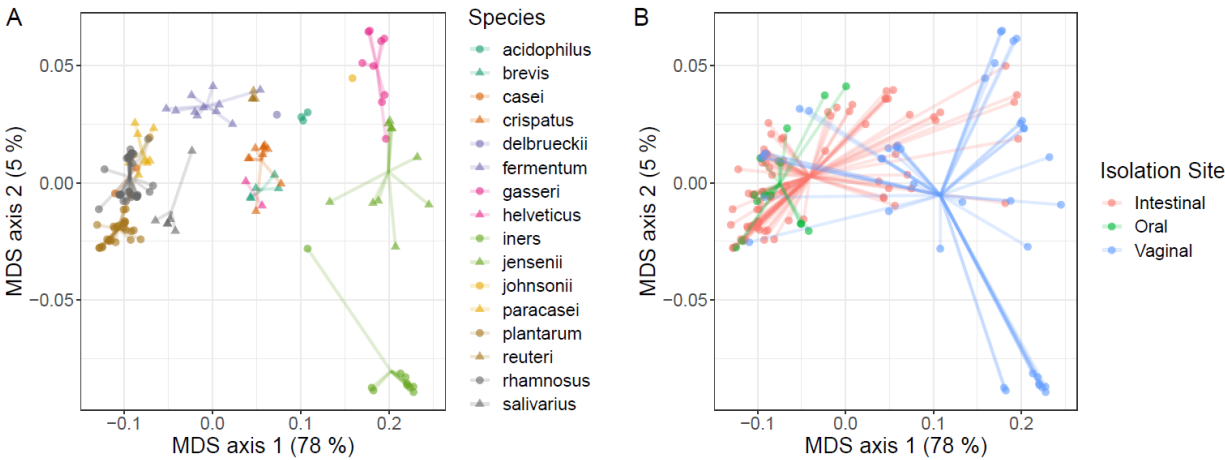


**Figure 4: The Scaled Production Likelihood metric distinguishes metabolic functionality among species.** (A) We found that *Lactobacillus* strains cluster significantly by species (PERMANOVA; P < 0.001). (B) Additionally, they cluster significantly by isolation site (PERMANOVA; P < 0.001). Both plots are PCoA using the Bray-Curtis distance metric of the SPLs for each isolate. Points in both panels are identical, but displayed with different color schemes.
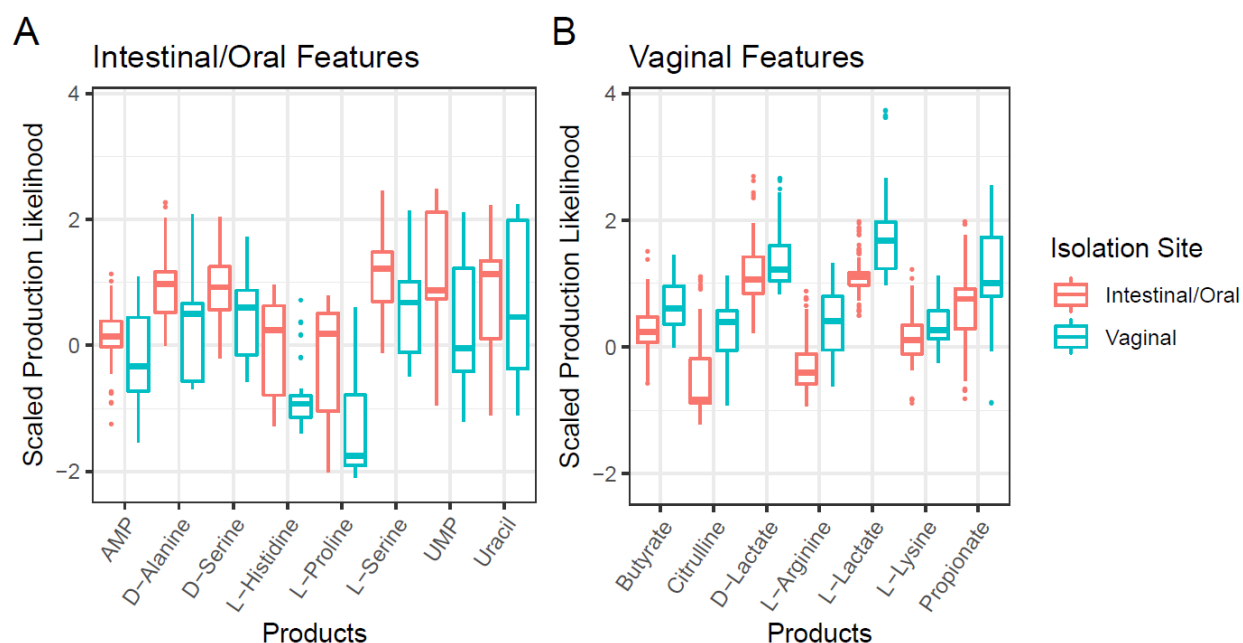
651

**Figure 5: Machine learning of the SPL scores identifies metabolites that discriminate *Lactobacillus* strains.**
Machine learning feature selection identified the metabolites that are both most likely to be produced by each group and capable of classifying the strains into two groups, intestinal/oral and vaginal, with greater than 90% accuracy. (A) There are eight metabolites that are more likely to be produced by the intestinal/oral isolates compared to the vaginal isolates. (B) There are seven metabolites that are more likely to be produced by the vaginal isolates compared to intestinal/oral isolates. Both models are more than 90% accurate in predicting the membership to which the given isolate belongs using the SPLs of the metabolites listed.