1  **Species level resolution of female bladder microbiota from marker gene**
2  **surveys**

3  Carter Hoffman, MS[1]; Nazema Y Siddiqui, MD, MHSc[2]; Ian Fields, MD[3]; W. Thomas Gregory,
4  MD[3]; Holly Simon, PhD[4]; Michael A. Mooney, PhD[1]; Alan J. Wolfe, PhD[5]; Lisa Karstens,
5  PhD[1,3]*
6

## 7  Abstract

8  The human bladder contains bacteria in the absence of infection. Interest in studying these
9  bacteria and their association with bladder conditions is increasing, but the chosen experimental
10 method can limit the resolution of the taxonomy that can be assigned to the bacteria found in the
11 bladder. 16S rRNA gene sequencing is commonly used to identify bacteria, but is typically
12 restricted to genus-level identification. Our primary aim was to determine if accurate species-
13 level identification of bladder bacteria is possible using 16S rRNA gene sequencing. We
14 evaluated the ability of different classification schemes, each consisting of combinations of a 16S
15 rRNA gene variable region, a reference database, and a taxonomic classification algorithm to
16 correctly classify bladder bacteria. We show that species-level identification is possible, and that
17 the reference database chosen is the most important component, followed by the 16S variable
18 region sequenced.
19
20 **Importance**
21 Species-level information may deepen our understanding of associations between bladder
22 microbiota and bladder conditions, such as lower urinary tract symptoms and urinary tract
23 infections. The capability to identify bacterial species depends on large databases of sequences,
24 algorithms that leverage statistics and available computer hardware, and knowledge of bacterial
25 genetics and classification. Taken together, this is a daunting body of knowledge to become
26 familiar with before the simple question of bacterial identity can be answered. Our results show
27 the choice of taxonomic database and variable region of the 16S rRNA gene sequence makes
28 species level identification possible. We also show this improvement can be achieved through
29 the more careful application of existing methods and use of existing resources.
30
31

32 Author Affiliations:
33     1. Department of Medical Informatics and Clinical Epidemiology, Oregon Health &
34        Science University, Portland, OR USA
35     2. Department of Obstetrics and Gynecology, Duke University, Durham, NC USA
36     3. Department of Obstetrics and Gynecology Oregon Health & Science University,
37        Portland, OR USA
38     4. AnimalBiome, Oakland, CA USA
39     5. Department of Microbiology & Immunology, Loyola University Chicago, Maywood, IL,
40        USA

41 *Corresponding author: karstens@ohsu.edu

1

42  **Introduction**

43  The human body provides a wide range of habitats, supporting a variety of microorganisms that
44  include bacteria, archaea, viruses and fungi, collectively known as the human microbiome(1).
45  Recent evidence from sequence-based and enhanced culturing techniques have revealed a
46  population of microbes (bacteria, fungi and viruses) that exist in the bladder, even in the absence
47  of infection(2–7). The discovery of the bladder microbiota (also known as the bladder urobiome)
48  has led researchers to question how these microbes influence the health of the host. Studies have
49  shown that altered bladder urobiome diversity is associated with urgency urinary incontinence
50  (UUI)(4,8), urinary tract infection after instrumentation of the urinary tract(9,10), and is
51  predictive of response to a common UUI drug (11). These studies collectively provide evidence
52  that the bladder urobiome, while previously overlooked, is clinically relevant and warrants
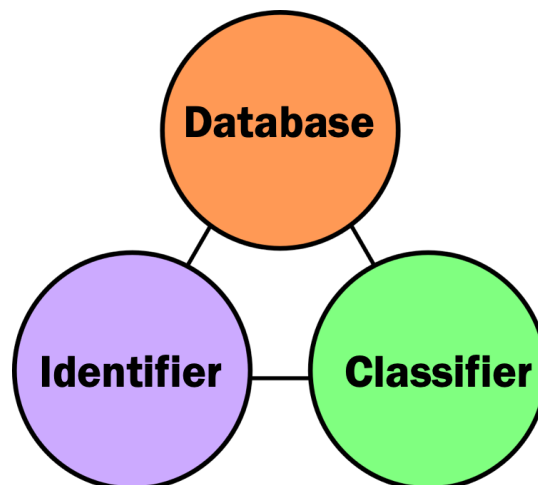53  further investigation.

54  To study the relationships between the bacteria found in the human bladder and health of the
55  host, it is necessary to accurately identify bacteria in a rapid and large-scale manner. Reliable
56  methods of determining the bacterial identity of an unknown bacterium include Matrix Assisted
57  Laser Desorption/Ionization-time of flight (MALDI-TOF) analysis or whole genome sequencing
58  (WGS) of purified colonies; both techniques permit species-level identification of bacteria(12).
59  However, culturing specific bacterial species is also time consuming and laborious. This
60  limitation has been circumvented by adopting culture-independent methods of sequencing DNA
61  directly from an environmental sample, such as shotgun metagenomic sequencing and targeted
62  amplicon sequencing, the latter most commonly involving the 16S rRNA gene(13). These
63  culture-independent sequencing methods are an attractive strategy because they can more
64  accurately reveal microbiota diversity by identifying bacteria that are difficult to grow in culture.

65  Targeted amplicon sequencing is currently the most practical method for identifying bladder
66  bacteria in a large-scale manner. When performing targeted amplicon sequencing, DNA is first
67  extracted from all cells in a sample, including host and bacterial cells. Next, the polymerase
68  chain reaction (PCR) is used to amplify a small segment of the bacterial genome. This segment is
69  then sequenced in a high-throughput manner. Finally, bioinformatics are used to process the
70  resulting sequences and identify the taxonomy of the bacteria. Algorithms compare the short
71  DNA sequences recovered from a sample to known sequences held in a reference database until
72  the closest match is found. In general, longer or more unique strings of sequenced DNA can be
73  used to identify bacteria at a higher level of precision, though sequence length is often limited by
74  the sequencing technology. The 16S rRNA gene is commonly used in amplicon sequencing
75  studies due to its universal presence in bacteria. The 16S rRNA gene conveniently contains
76  multiple "variable regions" with unique strings of sequence that can be used for bacterial
77  identification. A common target is the 4th variable region (V4), as this region has good
78  phylogenetic resolution down to the genus level for many bacteria(14).

79  When identifying bacteria using targeted amplicon sequencing there are three important
80  components (**Figure 1**). These components are: 1) the identifier, or DNA sequence of the
81  unknown bacterium; 2) a database of DNA sequences annotated with taxonomic information;
82  and 3) a classifier, which is the algorithm that compares the unknown sequence to those in a
83  database until the closest match is found. These components work together as a *classification*
84  *scheme*. One common classification scheme uses the V4 region from the 16S rRNA gene

2

85 sequence(15) as the identifier, the Silva database(16), and the Naïve Bayes algorithm(17). A
86 limitation of this particular classification scheme, and many others commonly used, is that the
87 phylogenetic resolution is usually constrained to the genus level. Recently, several new
88 approaches to sequence processing and taxonomy assignment have become available, which may
89 improve resolution to the species level (e.g. amplicon sequence variant algorithms such as
90 DADA2(18), and taxonomic classifiers such as Bayesian Lowest Common Ancestor
91 (BLCA)(19)).

92 Our primary aim was to determine if species-level identification of bladder bacteria is possible
93 from 16S rRNA gene sequencing studies. To achieve this aim, we used a representative sample
94 of bacteria found in the human female bladder and published by Thomas-White and
95 colleagues(20). This dataset includes bacteria found in the human female bladder that have been
96 cultured, isolated, subjected to whole genome sequencing, and identified using full length
97 sequences of 40 protein-encoding genes.. We used these known DNA sequences to determine
98 which classification schemes would be most useful for future targeted amplicon sequencing
99 studies. We evaluated several variable regions (i.e. potential identifiers), reference databases, and
100 taxonomic classification algorithms for their ability to accurately identify bladder bacteria at the
101 species level.



102

103 **Figure 1. Model of the components that make up a classification scheme to assign taxonomy to**
104 **unknown sequences.** A classification scheme consists of an identifier, a database, and a classifier. The identifiers
105 used in this study are subsequences of the 16S rRNA gene, computationally generated using published primers as
106 coordinates on the gene sequence. These targeted amplicons are the V3, V4, and V6 variable regions of the gene, or
107 span the V1-V3, V2-V3, V3-V4, and V4-V6 variable regions. The databases used in this study are the Greengenes,
108 Silva, and NCBI 16S. The classifiers used in this study are the Naive Bayes and Bayesian Lowest Common Ancestor
109 (BLCA) algorithms. One example of a classification scheme is the V4 region identifier, Silva database, and Naive
110 Bayes classifier. Another example classification scheme is the V6 region identifier, Greengenes database, and BLCA
111 classifier. These two examples are distinct from each other and can have different outcomes when assigning
112 taxonomy.

113 **Results**

114 **Representation of bladder bacteria in 16S rRNA gene sequence databases.** The Thomas-
115 White genome sequencing dataset consists of 149 bladder bacterial isolates, representing 78
116 bacterial species from 36 genera(20). There are several databases available for bacterial
117 identification using the 16S rRNA gene(21,22). Of these, Greengenes (v.13_5)(23), Silva (v.

118   132)(24), and NCBI 16S Microbial (v. August 2019) were evaluated due to their widespread use
119   in amplicon sequencing studies and availability of species-level annotation. At the genus level,
120   all but 1 genus (*Globicatella*) were present in the Greengenes database and all genera were
121   present in the Silva and NCBI 16S Microbial databases. At the species level, all 78 bladder
122   bacterial species were present in the Silva and NCBI 16S Microbial databases, whereas only 21
123   species were present in the Greengenes database.

124   **Information contained in variable regions differs for bladder bacterial species.** Sequencing
125   studies frequently focus on a small segment of the 16S rRNA gene that can be rapidly sequenced
126   in a high throughput manner using short read sequencing technology, such as the Illumina HiSeq
127   or MiSeq. To evaluate the performance of different variable regions as identifiers, amplicons
128   were computationally generated from the Thomas-White genome sequencing dataset for the V1-
129   V3, V2-V3, V3-V4, V4-V6, V3, V4, and V6 variable regions using published primers (see
130   **Methods**). These computational amplicons (**Figure 2A**) were used to determine how well the
131   currently available classification schemes can distinguish bladder bacterial species. To assess
132   different classification schemes, we tested multiple permutations of the variable regions listed
133   above with different databases (i.e. Greengenes, Silva, or NCBI 16S Microbial) and different
134   classifiers (i.e. Naive Bayes or BLCA, see **Figure 1**).

135   To quantify the amount of information contained across variable regions of the 16S rRNA gene
136   among commonly identified bladder bacteria, we performed a sliding window analysis on a
137   multiple sequence alignment (MSA) of all genomes from the Thomas-White dataset. We
138   calculated entropy as a measure of information content along the MSA (**Figure 2B**). As
139   expected, the defined variable regions contained regions of high entropy, suggesting variability
140   across species, whereas variable regions were flanked by conserved regions with low entropy
141   containing sequences that are similar among species. The V1 and V2 regions contained the
142   highest entropy, while V7 and V8 contained the lowest.

143   **Evaluation of classification scheme performance.** To evaluate the ability of
144   currently available resources to identify bladder species, we calculated the recall, precision and
145   F-measure for each classification scheme implemented (see **Methods**). Briefly, each resulting
146   taxonomic classification was evaluated as a true match, true non-match, false match or false non-
147   match based on whether the taxonomic classification was correctly assigned or not. Recall refers
148   to the proportion of matches that the classification scheme correctly identified out of all possible
149   matches. Precision refers to the proportion of matches that the classification scheme called
150   correctly out of all classified matches. The F-measure is the equally weighted harmonic mean of
151   recall and precision.

152   In general, the classification schemes that use the NCBI 16S Microbial database perform the best
153   (**Figure 3**), with high recall and precision (range 60.3%-91.0% for both classifiers). Those using
154   the Silva database show reduced precision and recall (range 23.1%-70.5% for both measures).
155   Because the Greengenes database is missing many of the bacterial species found in the bladder, it
156   is less precise. As such, classification schemes using the Greengenes database can have good
157   recall values (range 50.0%-81.8%), but the precision values are very low (range 22.0%-36.0%),
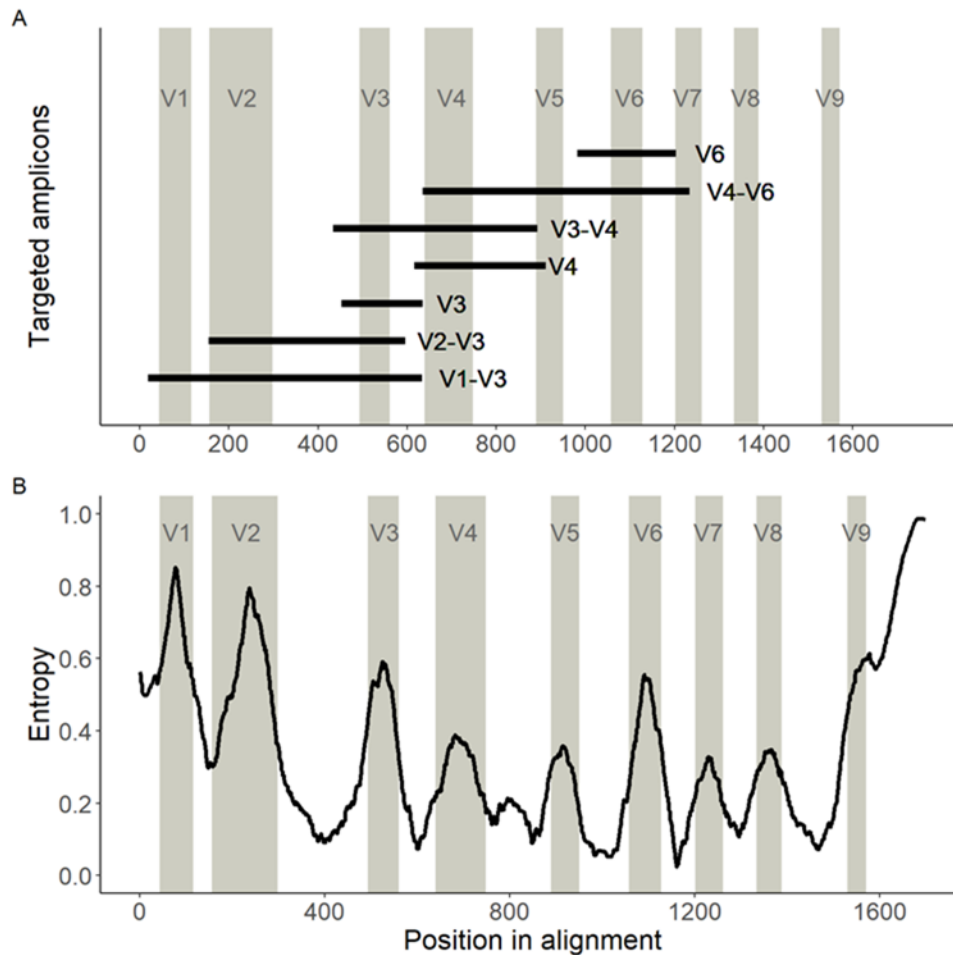158   indicating a large proportion of false matches to the number of true matches.

159



**Figure 2. Variable regions of the 16S rRNA gene used in this study**. A) Locations of the primers used in this study on the 16S rRNA gene. Locations of the predicted amplicons are shown as black bars in relation to the multi-sequence alignment (MSA) of the bacterial species described in Thomas-White et al. (2018). Gray columns are the locations of the known variable regions based on the sequence from *E. coli*. B) The information of variable regions, measured by entropy from a sliding window analysis of the MSA. Higher entropy indicates that the region has more variability across species, and therefore more information to identify a bacterial species. Lower entropy indicates that the region has little variability (i.e. is conserved) across species and therefore less information to identify a bacterial species.

160   When different variable regions are used as identifiers with Silva and NCBI 16S databases, there
161   are differences in classification scheme outcomes. Using the Silva database and Naive Bayes
162   classifier, the identifiers yielding the highest recall are the large V3-V4 (69.2%) and V4-V6
163   (69.2%) targeted amplicons. Using the Silva database and BLCA classifier, the V1-V3 and V4-
164   V6 amplicons have the highest recall (55.1% and 48.7%, respectively). In contrast, the identifiers
165   yielding the highest recall in classification schemes using the NCBI 16S database are the V1-V3
166   (90.3% on average) and V2-V3 (89.1% on average) targeted amplicons, regardless of the
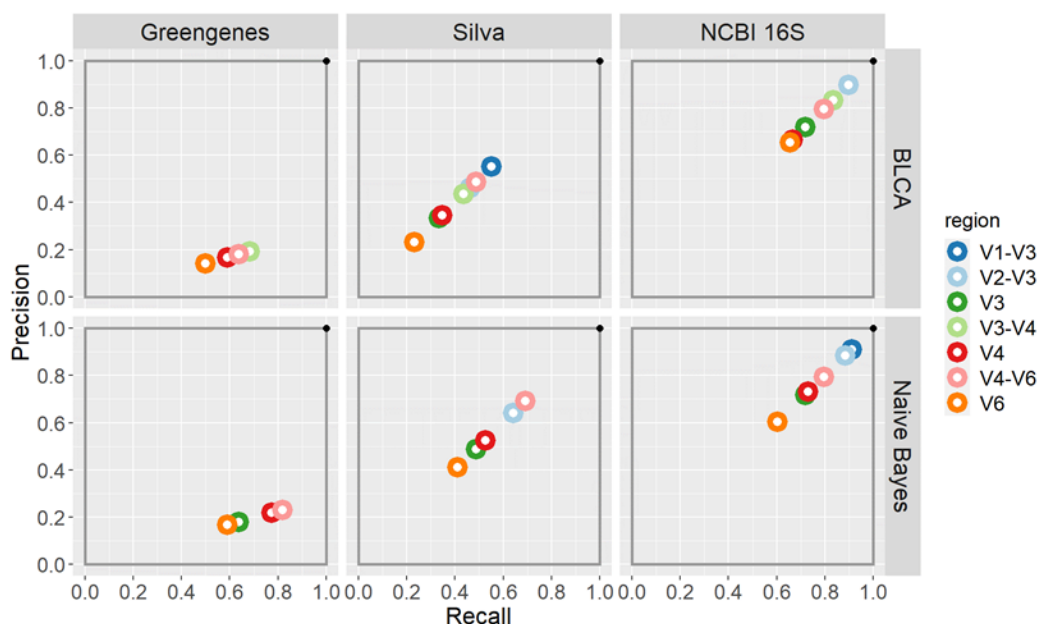167   classifier.

168

5

169



**Figure 3: Classification scheme evaluation when ignoring confidence scores.** The performance of each classification scheme is summarized by the precision (y axis) and recall (x axis) for each variable region (color). The best classification scheme would lie in the upper right-hand corner. Overall, classification schemes using the NCBI 16S Microbial database performed better than those using the Greengenes or Silva databases.

170 **Confidence scores affect classification**.

171 The BLCA and Naïve Bayes classifiers used in this study will classify an unknown sequence
172 even when the posterior probability for that taxon is very low. To account for this situation, a
173 confidence score is calculated that measures how much the classification changes through
174 random permutation (bootstrapping) and produces a value that reflects the "goodness of fit" of
175 that classification. When lacking any knowledge of how to choose the best confidence score that
176 minimizes the number of errors of a classification scheme (i.e. when a test set is not available),
177 using a predefined confidence score threshold is an option. Here, we evaluated the performance
178 of classification schemes when confidence score thresholds of 50% or 80% were used, such that
179 matches returned with confidence scores less than the threshold were considered non-matches.
180 **Figure 4** shows the effect of increasing the confidence score on the number of true matches
181 returned by each classification scheme.

182 Almost all classification schemes had a decrease in recall when using a default confidence score
183 of 80% (**Supplemental Figure 1**). This effect is especially marked for the classification schemes
184 that use the Silva database, which shows a 79.3% reduction in recall, on average. Classification
185 schemes that use the NCBI 16S database are unequally affected, for example the V1-V3
186 identifier shows a slight reduction in recall (7.1% on average), while the V6 identifier shows the
187 largest (43.3% on average). Classification schemes that use the Greengenes database are slightly
188 affected. These reductions in recall are mirrored in all classification schemes when a confidence
189 score of 50% is used as a threshold, but at a smaller magnitude.

6

190 Changes in the precision of the classification schemes are affected the most by the database used
191 (**Supplemental Figure 1**). For the classification schemes that use the NCBI 16S database,
192 precision is generally improved regardless of confidence score, but at unequal amounts. For
193 example, using the 80% threshold, the V1-V3 identifier shows a slight increase of 3.5% on
194 average, while the V6 identifier shows a large 39.7% increase on average. Classification schemes
195 that use the Silva database are unequally affected, with both reduction and gains in precision. A
196 dramatic increase in precision is shown by the classification scheme composed of the Silva
197 database, V4 identifier, and the Naive Bayes classifier. When ignoring a confidence score, this
198 classification scheme has a precision of 52.6%, but shows a 63.1% gain when using a confidence
199 score of 80% as a threshold. In general, precision is reduced when using the BLCA classifier and
200 the Silva database. As with recall, classification schemes that use the Greengenes database show
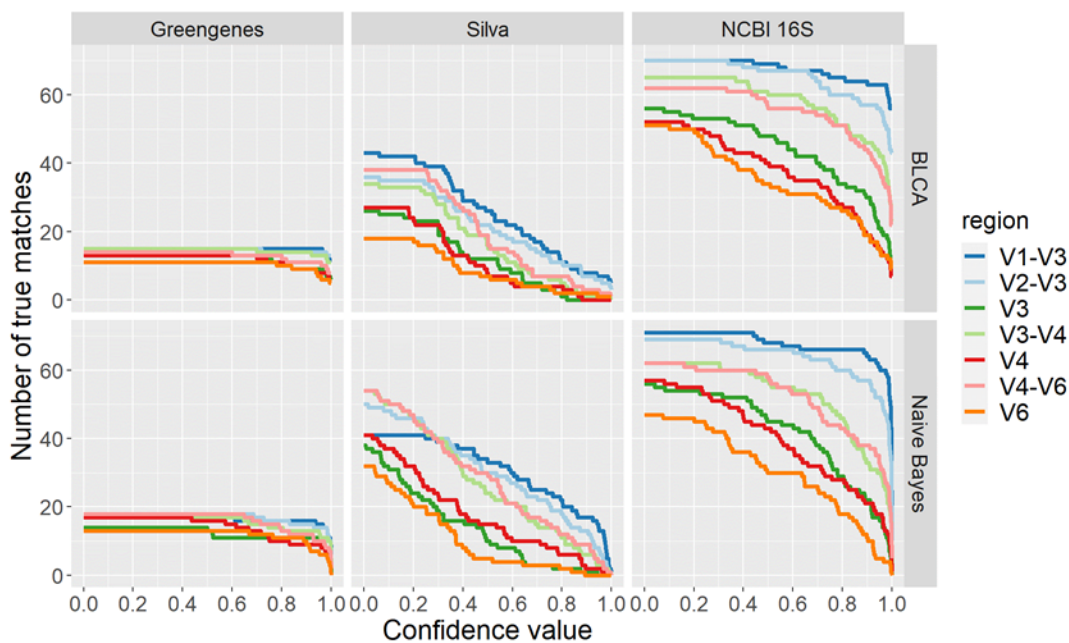201 slight changes in precision.



**Figure 4: The number of true matches returned for each classification scheme across all confidence score values.** As the confidence score value is increased, the number of true matches dramatically decreases, especially for schemes using the Silva database.

202 The overall changes in how these classification schemes perform when using a 50% or 80%
203 confidence score can be summarized by comparing the F-measure values shown in
204 **Supplemental Figure 2**. In almost every classification scheme, the F-measure value decreases
205 when a threshold is used, indicating a larger proportion of false matches and false non-matches
206 to the number of true matches. The classification schemes that use the Silva database clearly
207 demonstrate this effect, which show a 66.9% reduction in F-measure values on average. The
208 classification schemes that use the NCBI 16S database show slight decreases in the F-measure
209 values, with the exceptions of those that use the V3, V4 and V6 regions as identifiers. Those
210 classification schemes show a large 27.2% reduction in F-measure values on average. Finally,
211 the classification schemes that use the Greengenes database have slight changes in their F-
212 measure values, regardless of using a threshold or not.

213 **Amplicons spanning more than one variable region identify a higher number of bladder**
214 **bacterial species.** Amplicons spanning more than one variable region identified more unique
215 bladder bacteria at the species level than amplicons spanning a single variable region. For
216 example, with the commonly used V4 variable region and Naïve Bayes classifier, 21.8% of
217 bladder bacteria are correctly identified with the Greengenes database, whereas 52.6% are
218 identified with the Silva database and 73.1% with the NCBI 16S database (**Figure 5**). However,
219 with the NCBI database, when using amplicons spanning more than one variable region, such as
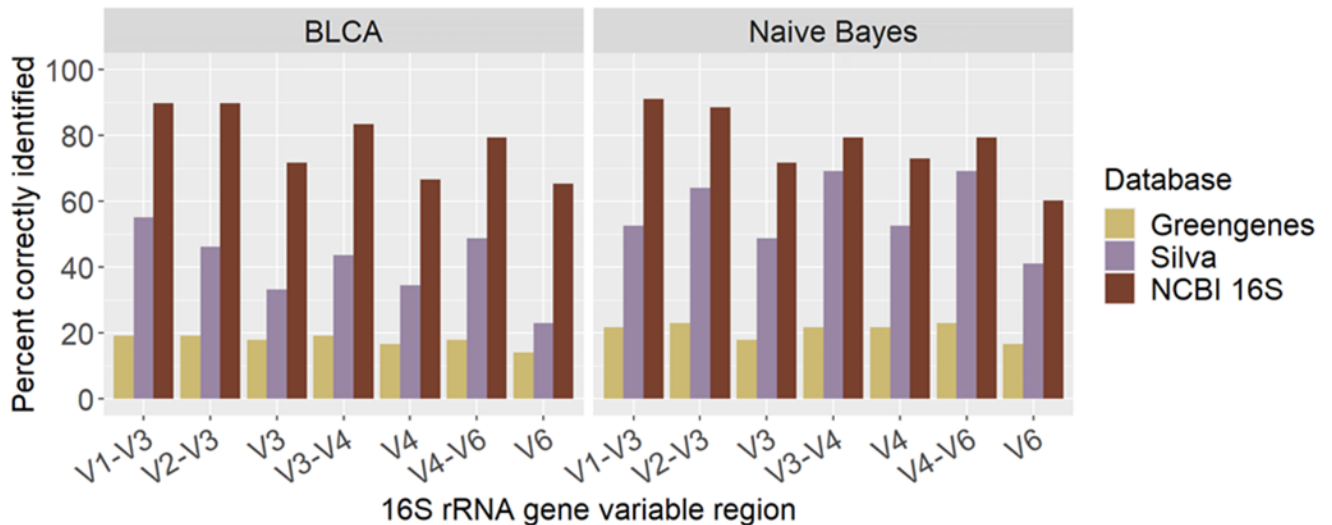220 the V1-V3 region, 91.0% of bacteria are correctly identified at the species level.



**Figure 5: Percent of bladder bacteria correctly identified for each classification scheme.** With the commonly used V4 variable region and BLCA classifier, 17% of bladder bacteria are correctly identified using the Greengenes database, compared with 35% correctly identified using the Silva database and 67% using the NCBI 16S database. A similar trend is seen with the Naïve Bayes classifier. Using other variable regions can lead to improved species-level resolution to a maximum number of 91% correctly identified.

221 **Species identified depends on choice of database and variable region.** While the results thus
222 far have focused on summarizing overall performance of classification schemes for identifying
223 bladder bacteria at the species level, we also sought to determine which classification schemes
224 could be used to identify specific bacteria (**Table 1, Supplemental Figure 3**). Although the
225 NCBI database contains the largest representation of bladder species, some species were not
226 identified with certain variable regions, if at all. For example, *Lactobacillus* species were overall
227 best represented within the NCBI database, with 8 out of 9 species being identified with the V1-
228 V3 and V2-V3 variable regions (**Figure 6**). However, the other variable regions only identified
229 between 4 and 6 *Lactobacillus* species when using the NCBI database. Interestingly,
230 *Lactobacillus crispatus* was identifiable with the Silva and NCBI databases when using the V4-
231 V6 regions, but only with the NCBI database using the V1-V3 and V2-V3 regions, and only the
232 Silva database when using the V4 and V6 regions independently. *Lactobacillus iners* was not
233 correctly identified from our dataset with any classification scheme.

234 Additionally, we found that there were important discrepancies for bacteria that are thought to
235 play a role in bladder health and disease (**Supplemental Figure 3**). Several bladder species, such
236 as *Gardnerella vaginalis,* were only detected with the NCBI and Silva databases. *Staphylococcus*

8

237    species were poorly identified with the V4 region but were distinguishable with all other regions.
238    *Streptococcus* and *Corynebacterium* species were best identified with NCBI. *Escherichia coli* is
239    not well represented in any of the databases, and was only detected with the V4 region and the
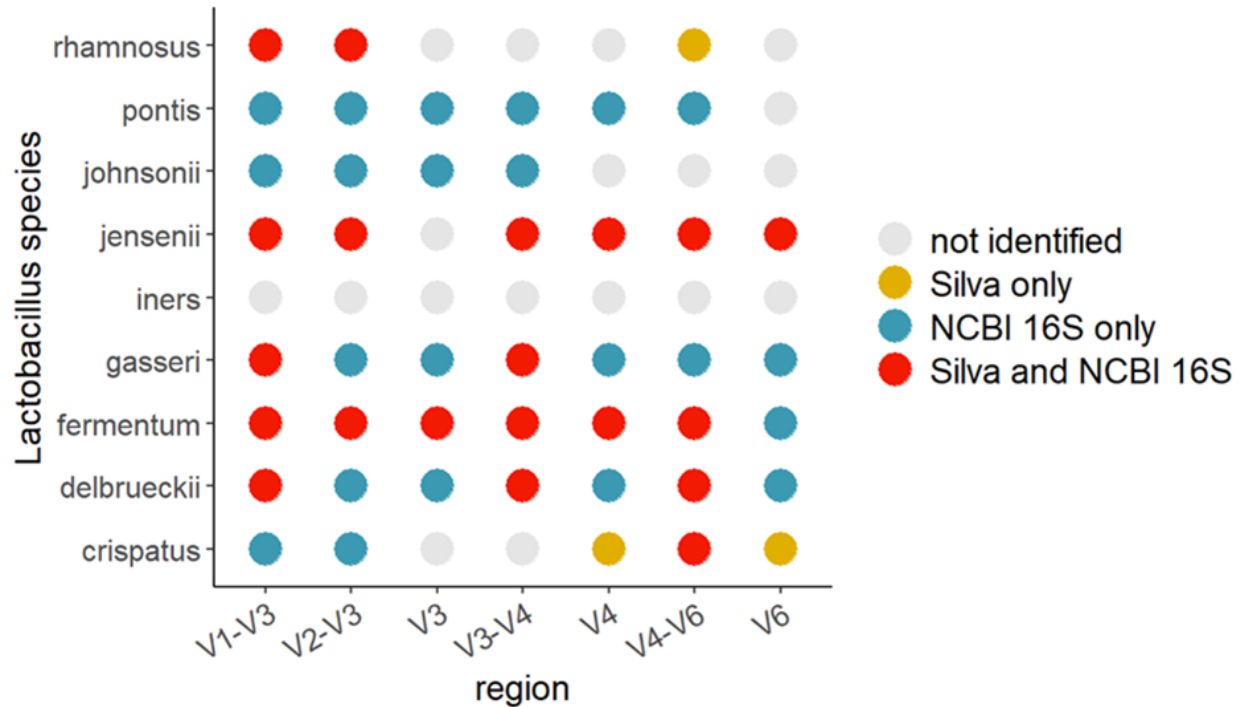240    NCBI database.



**Figure 6: The ability of classification schemes to distinguish between different Lactobacillus species**. Results shown for classification schemes using the BLCA classifier. Classification schemes using the NCBI 16S database have the most coverage, regardless of variable region chosen. The Greengenes database is not shown since it only classified two species (*L. pontis* and *L. delbrueckii*).
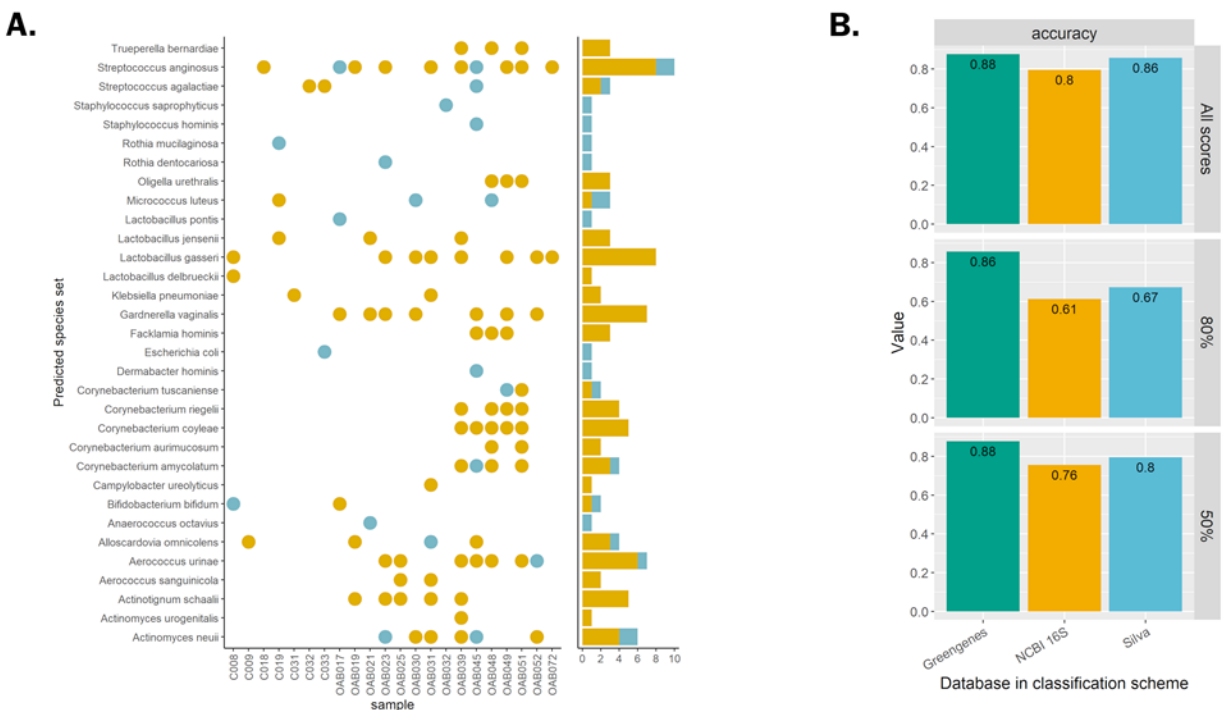
241

242 **Table 1.** Number of species identified with each database by variable region (BLCA).

| Genus | V1-V3 | V2-V3 | V3 | V3-V4 | V4 | V4-V6 | V6 |
|---|---|---|---|---|---|---|---|
| **Actinomyces** | | | | | | | |
| NCBI 16S | 3 | 3 | 3 | 3 | 2 | 3 | 3 |
| Silva | 3 | 4 | 1 | 1 | 1 | 1 | 1 |
| Greengenes | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **Aerococcus** | | | | | | | |
| NCBI 16S | 3 | 3 | 3 | 3 | 3 | 3 | 3 |
| Silva | 3 | 3 | 3 | 3 | 3 | 3 | 0 |
| Greengenes | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **Bifidobacterium** | | | | | | | |
| NCBI 16S | 3 | 3 | 3 | 3 | 2 | 2 | 2 |
| Silva | 3 | 2 | 2 | 2 | 2 | 1 | 1 |
| Greengenes | 3 | 3 | 3 | 3 | 2 | 2 | 2 |
| **Corynebacterium** | | | | | | | |
| NCBI 16S | 7 | 7 | 6 | 7 | 6 | 7 | 5 |
| Silva | 1 | 2 | 1 | 2 | 0 | 1 | 1 |
| Greengenes | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **Lactobacillus** | | | | | | | |
| NCBI 16S | 8 | 8 | 5 | 6 | 5 | 6 | 4 |
| Silva | 5 | 3 | 1 | 4 | 3 | 5 | 2 |
| Greengenes | 2 | 2 | 2 | 2 | 2 | 2 | 1 |
| **Staphylococcus** | | | | | | | |
| NCBI 16S | 4 | 5 | 4 | 4 | 2 | 4 | 3 |
| Silva | 3 | 2 | 1 | 2 | 1 | 2 | 2 |
| Greengenes | 1 | 1 | 1 | 1 | 0 | 1 | 1 |
| **Streptococcus** | | | | | | | |
| NCBI 16S | 9 | 9 | 6 | 7 | 6 | 7 | 6 |
| Silva | 4 | 5 | 2 | 2 | 2 | 2 | 2 |
| Greengenes | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| **Other** | | | | | | | |
| NCBI 16S | 33 | 32 | 26 | 32 | 26 | 30 | 25 |
| Silva | 21 | 15 | 15 | 18 | 15 | 23 | 9 |
| Greengenes | 8 | 8 | 8 | 8 | 8 | 8 | 6 |
| **Total** | | | | | | | |
| NCBI 16S | 70 | 70 | 56 | 65 | 52 | 62 | 51 |
| Silva | 43 | 36 | 26 | 34 | 27 | 38 | 18 |
| Greengenes | 15 | 15 | 14 | 15 | 13 | 14 | 11 |

243

244    **Validation of computational findings on V4 amplicon data.** To evaluate the performance of
245    our computational findings on actual data, we acquired targeted amplicon sequencing data from
246    24 urine samples. These urine samples were a subset of those that were used to derive cultures in
247    the Thomas-White dataset and thus should contain the same bacteria.  Sequencing data were
248    generated as part of two other studies using Illumina sequencing of the V4 region of the 16S
249    rRNA gene(4,11). We reprocessed the raw sequencing data (see **Methods**) and performed
250    taxonomic classification to assess the performance of our computational findings. Since 16S
251    rRNA gene sequencing will detect many more bacteria than those identified even with enhanced
252    culture, we restricted the evaluation to only the bacteria that grew in culture from a given sample.
253    We used accuracy to assess the number of predicted matches that were correctly identified in the
254    V4 dataset, using classification schemes composed of the V4 identifier, each of three databases,
255    and the BLCA classifier (**Figure 7**). All databases had good accuracy with high proportions of
256    accurate identifications at the species level (80% for NCBI 16S, 86% for Silva, and 88% for
257    Greengenes). Accuracy was reduced when the default confidence score of 80% was applied
258    (61% for NCBI 16S, 67% for Silva, and 86% for Greengenes). The default confidence score of
259    50% reduced the accuracy of two of the classification schemes (76% for NCBI 16S and 80% for
260    Silva). We also evaluated classification schemes with the Naive Bayes classifier and found
261    similar results (**Supplemental Figures 4 and 5**)



262

**Figure 7.** Taxonomic classification of the V4 validation dataset. A) Results when using a classification scheme including the V4 identifier, NCBI 16S database, and BLCA classifier. Blue dots represent species identified in cultured isolates, but not identified in targeted amplicon sequencing using this classification scheme. Yellow dots represent the species that were present in cultured isolates and successfully identified by the classification scheme. B) Summary of accuracy for classification schemes that use the V4 rRNA identifier, BLCA classifier, and the three databases (Greengenes, Silva and NCBI 16S). Rows show accuracy results when ignoring confidence scores, and when using confidence scores of 50% or 80% as thresholds.

11

## Discussion

264  Our study demonstrates that it is possible to gain higher resolution results at the species level
265  with existing resources when performing targeted amplicon sequencing of urinary specimens.
266  Though higher resolution is possible, it requires a carefully chosen classification scheme. Within
267  the classification scheme, the reference database strongly influences the identification of bacteria
268  at the species level.  Overall, we found the NCBI 16S database performs the best, whereas the
269  Greengenes database performs the worst, primarily because it does not currently contain
270  representatives of bladder bacteria. The identifier, or 16S rRNA variable region that is chosen,
271  can also influence the types of bacterial species that are identified. The choice of classifier did
272  not drastically affect the identification of species and thus is less critical within the classification
273  scheme.

274  The largest limitation of any reference database is that the number of records of accurately
275  classified bacteria is dwarfed by the number and diversity of unidentified sequences obtained
276  through metagenomic sequencing of environmental samples. Because of the considerable
277  amount of work required to construct and maintain databases, they will undoubtedly
278  incompletely represent existing bacteria.

279  For species level taxonomy assignments, the reference database must contain species-level
280  information. In other words, if species of bacteria are expected in a sample, it must be verified
281  that the database contains those species. For example, we found that the Greengenes database
282  does not currently contain many bacterial species that are found in the human bladder. In
283  contrast, the NCBI 16S Microbial and Silva databases had representation of all species that were
284  identified from prior studies of bladder bacteria. Thus, the latter two databases are better choices
285  for evaluating bacterial species from the bladder.

286  While the databases reviewed in this study do have species-level information associated with the
287  records, additional work was needed before species-level identification could be achieved with
288  the Naïve Bayes classifier. This classifier requires a database that has undergone the "training"
289  steps that convert the DNA sequences to the calculated frequencies that each $k$-mer occurs in a
290  taxon. For available classification algorithms like the RDP classifier[25] and QIIME2[17], the
291  training is only currently done to reliably identify bacteria to the genus level. For this study, it
292  was necessary to train the Silva and NCBI 16S databases to the species level for use with the
293  Naive Bayes classifier. While training the reference databases did take significant computational
294  effort, once completed it was used repeatedly.

295  The classifiers used in this study are examples of two different strategies designed to overcome
296  the common challenges of searching an extremely large dataset in order to find matching pairs of
297  query sequences and reference records. While these two approaches are different in concept, we
298  did not find significant differences in their performance for species-level classification of bladder
299  bacteria.

300  BLCA is an example of sequence comparison by pairwise alignment. The strength of this
301  method is due to the fact that the similarities between two DNA samples are directly compared.
302  This is the most effective way to compare the characteristics of a sample to those that define a
303  taxon; however, until recent advances in computer technology, it remained impractical because

12

304     of the computational burden. The Naive Bayes classifier is an example of a *k*-mer-based
305     classification approach, and was designed to circumvent the computational challenges that are
306     faced with use of a pairwise alignment classifier. However, there are limitations when using
307     Naive Bayes for species-level identification. The first limitation arises from the database training
308     process. If one taxon has more training examples than another, Naive Bayes generates
309     unfavorable weights for the decision boundary(26). The second limitation is that all features (i.e.
310     the *k*-mers generated from the DNA sequences) are assumed to be independent, and weights for
311     taxa with strong dependencies among the associated *k*-mers are larger than those taxa with
312     weakly dependent *k*-mers(26).

313     Finally, as shown by both the computational and V4 validation results, the use of the 50% or
314     80% confidence score thresholds significantly reduced the recall and accuracy of the
315     classification schemes. Precision increased in several cases, for example with classification
316     schemes that use the NCBI 16S database or those that use the Silva database and Naïve Bayes
317     classifier, but at the cost of severely decreasing the number of species identified. These results
318     show that the default settings of 50% or 80% are restrictive, and limit the ability to detect bladder
319     species, especially when using the Silva and Greengenes databases. This could be resolved
320     through the use of a comparative data set to find the confidence score values yielding optimal
321     performance of these classification schemes.

322     Affordable sequencing of large-scale data is presently done with short read sequencing
323     technology, such as Illumina MiSeq. This is currently limited to sequencing reads up to 300
324     nucleotides in length. Until full-length 16S rRNA gene sequencing can be achieved affordably
325     on a large scale (such as with Oxford Nanopore and PacBio technologies), choosing the optimal
326     region of the 16S rRNA gene for identification purposes remains a significant part of the
327     experimental design. Thus, the variable regions that are used as identifiers require some
328     consideration.

329     Our findings show that use of the V2-V3, and V1-V3 regions of the 16S rRNA gene allowed for
330     the correct identification of the most bladder bacterial species when combined with the NCBI
331     16S database and either classifier. In general, amplicons that span more than one variable region
332     perform better than those that contain single variable regions. This is likely due to the increased
333     information available with longer reads. It is important to note that longer reads can also have
334     limitations, which are discussed in more detail below.  While shorter variable regions, such as
335     the V4 region, did not perform as well as longer amplicons, they were able to identify many
336     bladder bacteria at the species level (52 out of 78). These shorter amplicons are widely used with
337     Illumina sequencing and may be valid, depending on the study design and level of precision
338     desired. However, other variable regions may be explored for practical application, or when
339     more detailed information is desired.

340     By taking a computational approach to evaluating classification schemes that are capable of
341     identifying bladder bacterial species, we were able to thoroughly assess the ability of
342     classification schemes to identify known bacterial species. However, there are several practical
343     limitations of amplicon sequencing that were not captured in this approach.

344     Targeted amplicons are generated by priming the polymerase chain reaction with specially
345     designed oligonucleotides (PCR primers). The challenge of PCR primer design is to identify a

13

346    sequence of nucleotides that will anneal to only one location on the template DNA. Finding
347    suitable annealing sites that flank the variable region of interest becomes very difficult when
348    considering the 16S rRNA gene sequence of many species. We used published primers to create
349    computational amplicons, but this may not reflect the actual experimental efficiency. Finally,
350    quality control with DNA sequence processing must be conducted before classification is
351    performed. Targeted amplicon sequencing generates a large number of overlapping reads and
352    provides the data for methods to correct for errors introduced by the polymerase enzyme. After
353    error correction, similar reads are aggregated into operational taxonomic units or amplicon
354    sequence variants. The last step is to attempt to merge reads that are complementary before
355    attempting to classify them. If the sequence reads do not overlap, loss of phylogenetic
356    information occurs in the gaps and impacts the accuracy of identification, which may occur for
357    longer amplicons, such as the V1-V3 and V4-V6 regions.

358    In our study, we identified the V1-V3 region of the 16S rRNA gene as having the greatest
359    taxonomic resolution for the bacteria that are found in the bladder. This may be attributed to the
360    high occurrence of insertions and deletions (indels) in the conserved regions between the first
361    three variable regions across the bacteria in the Thomas-White dataset. Designing one degenerate
362    primer set that would amplify the entire dataset may not be possible for this region. A future
363    research direction could be to stratify the Thomas-White dataset into smaller, more closely
364    related phylogenetic groups for more specific primer design.

365    **Conclusion**

366    Species level taxonomy assignment will greatly benefit studies focused on the urobiome and its
367    relationship to bladder health and disease. Our results show that it is possible to reliably classify
368    bladder bacterial species using targeted amplicon sequencing of the 16S rRNA gene variable
369    regions with existing classification algorithms and databases. We determined that the most
370    important component of the classification scheme is the database used, and that the NCBI
371    database allows for best identification of bladder species. Our validation with V4 amplicon data
372    demonstrates that the predicted computational outcomes are a good approximation for how a
373    classification scheme will perform on real data. The knowledge that a majority of the predicted
374    matches reflect reality is encouraging. It can be expected that the alternate variable regions
375    covered in this study, such as the V2-V3 region of the 16S rRNA gene, would have similar
376    outcomes.

377    Importantly, we found that no single variable region gives 100% coverage of all bladder bacteria
378    species. Thus, the choice of variable region may significantly affect the results of a given study.
379    One approach to resolve this could be to use multiple amplicon sequencing or long read
380    sequencing technology. These technologies are emerging and may prove to be beneficial for the
381    urobiome community. Furthermore, no database has 100% coverage across a variable region.
382    This could be resolved by using more than one database for classification, though this approach
383    is complicated by differences in databases in terms of formatting, as well as conflicting
384    classifications. Both of these components are important for planning experimental and
385    computational aspects of urobiome studies, and should be considered when comparing results
386    across studies.

387

14

388 **Material and Methods**

389 **Code resources.** All scripts that were written for this project can be found in the GitHub
390 repository (https://github.com/lakarstens/BladderBacteriaSpecies). All scripts sourced from this
391 repository are referred to as "custom."

392 **The Thomas-White dataset.** The 78 species of bladder bacteria used in this study were
393 identified by culturing 149 urine samples and performing whole-genome sequencing, as
394 described in Thomas-White et al.(20). This set of identified species served as the basis for our
395 computational analysis and is referred to as the Thomas-White dataset. For each species
396 identified, the 16S rRNA gene sequence of the corresponding type strain was downloaded from
397 the Silva v132 release (https://www.arb-silva.de/) on 4/27/2019. A *type strain* is the sequence of
398 the cultured isolate that was subject to the metabolic, genotypic and phenotypic evaluations taken
399 to define the bacterial species(27), and is the agreed bacterial organism to which the taxonomic
400 name refers. Sequences were searched using the "[T]" filter setting, and sequences longer than
401 1450 nt with alignment and pintail quality scores greater than 95% were selected. For the species
402 that had no hits, the taxonomic synonym (see below) was used as the search query, if available.
403 One unidentified *Corynebacterium* species had no type strain available, and was excluded from
404 the analysis.

405 **The V4 validation dataset.** Targeted amplicon sequences from 24 urine samples, using the
406 V4 region of the 16S rRNA gene sequence, is referred to as the V4 validation dataset. These 24
407 urine samples originated from a subsample of the women whose samples comprised the Thomas-
408 White dataset. Sequencing data were generated as part of two other published studies using
409 Illumina sequencing of the V4 region of the 16S rRNA gene(4,11). The raw sequence reads were
410 processed with DADA2 version 1.14.1(18) to generate amplicon sequence variants (ASVs). The
411 ASVs were subjected to taxonomic classification with the BLCA algorithm.

412 **Synonyms of species.** Species names have changed in response to advances in bacterial
413 systematics. All currently known species synonyms were downloaded from the Prokaryotic
414 Nomenclature Up-to-Date(28) (PNU) website on 1/5/2020. PNU includes information down to
415 the strain level, but these entries were consolidated to the species level. For example, entries like
416 *Enterobacter cloacae* and *Enterobacter cloacae dissolvens* are treated as synonyms of
417 *Enterobacter cloacae*. Classification results were then checked for synonyms using the custom
418 "validate_match_batch.py" script.

419 **Databases.** The Greengenes database version 13_5 was downloaded on 9/23/19 from
420 (http://greengenes.secondgenome.com/?prefix=downloads/greengenes_database/gg_13_5/). For
421 use with BLCA, the database was processed using the provided "1.subset_db_gg.py" script
422 (https://github.com/qunfengdong/BLCA/). For use with the Qiime2 package, the FASTA file
423 was reformatted to work with Qiime2 using the custom "write_qiime_train_db.py" script, and
424 trained to work with the Naive Bayes classifier with the provided "fit-classifier-naive-bayes"
425 script.

426 The Silva database version 132 was downloaded on 9/14/19 from (https://www.arb-
427 silva.de/no_cache/download/archive/release_132/Exports/) as a FASTA formatted file. The
428 FASTA file was compiled into a database that could be used with BLCA by using the

15

429    "makeblastdb" utility provided in the Blast+ suite. The taxonomy file that was required by
430    BLCA was generated with the custom "write_taxonomy.py" script. For use with the Qiime2
431    package, the FASTA file was reformatted to work with Qiime2 using the custom
432    "write_qiime_train_db.py" script, and trained to work with the Naive Bayes classifier with the
433    provided "fit-classifier-naive-bayes" Qiime2 script.

434    The 16SMicrobial database is bundled with the BLCA package, but is available from
435    (ftp://ftp.ncbi.nlm.nih.gov/blast/db/). For use with BLCA, the database was processed using the
436    provided "1.subset_db_acc.py" script included with BLCA. For use with the Qiime2 package, a
437    FASTA file was extracted from the bundled BLCA database using "blastdbcmd" utility provided
438    in the Blast+ suite, and reformatted to work with Qiime2 using the custom
439    "write_qiime_train_db.py" script. The database was trained to work with the Naive Bayes
440    classifier with the provided "fit-classifier-naive-bayes" script included in Qiime2.

441    **Presence of Thomas-White species in databases.** To verify that all species from the
442    Thomas-White dataset were present in the databases used in this study, each database was first
443    converted to a FASTA file (if needed) using the "blastdbcmd" utility included in the Blast+ suite.
444    The FASTA file was then searched using regular expressions and the Linux command-line
445    program *grep* for a match of each species in the dataset. The commands were implemented using
446    the custom "species_in_db.bash" script. The presence or absence of each species was recorded.

447    **Multisequence alignment**. The 16S gene sequences from the Thomas-White dataset were
448    formed into a multi-sequence alignment using the T-coffee program(29). T-coffee version
449    12.00.7fb08c2 was downloaded from (http://tcoffee.org/Packages/Stable/Latest/) on 4/5/2019.
450    Alignments were performed using the default settings.

451    **Sliding window analysis**. Comparing the 16S rRNA gene sequences of the species in the
452    Thomas-White dataset reveals regions of conserved sequence and regions of variability. The
453    degree that variable regions of species differ from each other can aid the identification of each
454    species; therefore, quantifying the amount of variability of a region across a set of species is
455    important.

456    Sliding window analysis (SWA) is the method by which a list of subsequences are generated by
457    taking successive groups of equal size, in the manner of a window of fixed length sliding across
458    the full sequence. Quantifying the amount of variability along a MSA is achieved by combining
459    SWA with calculating the Shannon Entropy contained in each column framed by the window.

460    The minimum Shannon entropy occurs when all nucleotides in a position (column) of the MSA
461    are the same. The maximum occurs when all possible nucleotides in the MSA are present at that
462    position. However, the Shannon Entropy treats gaps in a sequence as relevant, where in practice
463    gaps reflect an absence of useful information. Multisequence alignments can generate many
464    columns of gap characters due to insertions or deletions (indels) in the respective sequences that
465    make up the MSA. A consequence of treating gaps as relevant is the Shannon Entropy will
466    interpret these indel regions as conserved sequence. This limitation was overcome by weighting
467    the entropy scores against gaps(30). The locations of known variable regions of the 16S gene
468    sequence were validated, and the relative amount of variability was quantified, using the custom
469    "weighted_ent.py" script.

16

470    **Primers**. Amplicons were computationally generated from the Thomas-White genome
471    sequencing dataset for the V1-V3, V2-V3, V3-V4, V4-V6, and V4 variable regions using
472    published primers and the V3 and V6 regions using designed primers. The primer sequences
473    used, listed in order of amplicon spanning variable region(s), forward primer name and sequence,
474    reverse primer name and sequence are: **V1-V3**: A17F 5'-GTT TGA TCC TGG CTC AG-3', 515R
475    5'-TTA CCG CGG CMG CSG GCA-3'(31,32). **V2-V3**: 16S_BV2f 5'-AGT GGC GGA CGG
476    GTG AGT AA-3', HDA-2 5'-GTA TTA CCG CGG CTG CTG GCA C-3'(33,34). **V3:** v3_579F
477    5'-THT TSS RCA ATG GRS GVA-3', v3_779R 5'-GKN SCR AGC STT RHY CGG-3'. **V3-V4:**
478    V3F 5'-CCT ACG GGA GGC AGC AG-3', V4R 5'-GGA CTA CHV GGG TWT CTA AT-
479    3'(35). **V4**: F515 5'-GTG CCA GCM GCC GCG GTA A-3', R806 5'-CCT GAT GHV CCC
480    AWA GAT TA-3'(36). **V4-V6:** 519F 5'-GTG CCA GCT GCC GCG GTA ATA-3', 1114R 5'-
481    GGG GTT GCG CTC GTT GC-3'(32). **V6:** v6_1183F 5'-CCG CCT GGG GAS TAC GVH-3',
482    v6_1410R 5'-AGT CCC RYA ACG AGC GCA-3'. Degenerate primer design was employed to
483    generate primer sets for the V3 and V6 regions of the 16S rRNA gene that would anneal to as
484    many species in the Thomas-White dataset as possible with DegePrime(37)
485    (https://github.com/EnvGen/DEGEPRIME.git). DegePrime has the option to ignore columns of a
486    MSA if the number of "-" characters exceed a user-defined threshold. The MSAs were
487    preprocessed with this threshold set to .01. The main script of DegePrime was run using a
488    degeneracy setting of 4096 and a window length of 18.

489    **Extracting computational amplicons.** For each primer set, the DNA sequence
490    bracketed by the forward and reverse primers was extracted from the multisequence alignment.
491    Coordinates of the MSA were identified by searching the *E. coli* sequence (accession number
492    EU014689.1.1541) included in the MSA for a match to the forward and reverse primer
493    sequences, and then mapping those position to the MSA of the Thomas-White dataset. This
494    procedure was done using the custom "extract_16s_vr.py" script and output as a multi-record
495    FASTA formatted file.

496    **Taxonomic classifiers.** Taxonomic classification was performed with Bayesian lowest
497    common ancestor (BLCA) and Naïve Bayes classifiers. BLCA(19) was cloned from the
498    GitHub repository https://github.com/qunfengdong/BLCA.git. For the 16S variable regions, the
499    BLCA was run using default settings but pointing to the selected reference database, either
500    Greengenes, Silva, or NCBI 16S. The Naïve Bayes classifier as implemented by Qiime2(17) was
501    used with the Greengenes, Silva, and NCBI 16S databases and a confidence setting of 0, 50, and
502    80, but otherwise default settings.

503    **Evaluating computational results.** To evaluate the taxonomic classification results of
504    each classification scheme on the computational amplicon dataset, the custom
505    "new_taxonomy_results_2020-3-14.Rmd" file was used. These scripts compare the results of
506    each record pair (each comparison between the query sequence and sequence held in the
507    reference database) from the classification scheme to the known identify of the query sequence
508    from the Thomas-White dataset. All record pairs that were assigned a match by the classification
509    schemes were evaluated according to the following definitions (**Figure 8**):

510        ***True match*** - All record pairs assigned as a match that have identical genus and species
511        labels.

17

512    *False match* - All record pairs assigned as a match that did not have identical genus and
513    species labels.
514    *False non-match* - If a record representing a species in the Thomas-White dataset was
515    present in the database, but was not assigned as a match, the record was evaluated as a
516    false non-match.
517    *True non-match* - All records in the reference database that were not in the Thomas-
518    White dataset. While records assigned to this category were not used in evaluating the
519    classification schemes in this manuscript, the definition is still included for completeness.
520
521
522



**Figure 8. Example of classification evaluation used in this study**. Suppose there is a classification scheme comprising a set of query sequences (the rows E,F,G) and the set of reference sequences (the columns E,F,L,M) held in a reference database. In this example, the number of reference records is greater than the query records, and the reference is missing a corresponding G record from the query set. **A**) If the query and reference record letters are the same, then they are designated as a **match**. If they are different they are designated as a **non-match**. **B**) Next, the classifier is allowed to assign record pairs as matches or non-matches for all query sequences, represented as green plus signs for matches and blank cells as non-matches. Some results are correct, and some are not. Note that despite the lack of a matching record in the reference database, the classifier still designated the (G:M) pair as a match. **C**) Using the definitions for assigning the classifications to the confusion matrix, there is one **true match** (green square), two **false matches** (red squares), one **false non-match** (yellow square), and 8 **true non-matches** (white squares). D) The cell values of the confusion matrix are then filled out, and performance measurements can be calculated. For this classification scheme, the precision is $1/(1+2)=.33$, recall is $1/(1+1)=.5$, and the F-measure is $(2*.33*.5)/(.33+.5)=.40$.

523    **Performance measures.** Recall, precision and the F-measure were used to evaluate the
524    performance of each classification scheme implemented. Recall refers to the proportion of
525    matches that the classification scheme correctly identified (true matches) out of all possible
526    matches (true matches plus false non-matches). Precision refers to the proportion of matches that
527    the classification scheme called correctly (true matches) out of all classified matches (true
528    matches plus false matches). The F-measure is the equally weighted harmonic mean of recall and
529    precision. For this study, we chose to maximize recall over precision, because the number of true
530    matches impacts the subsequent work on diversity measures, such as species richness and
531    evenness(38).

532    **Evaluating V4 validation results.** The species of bacteria in the V4 seqeuncing data
533    were identified using classification schemes composed of the V4 sequencing results as the
534    identifier, BLCA classifier, and the Greengenes, Silva, and NCBI 16S microbial databases. To
535    determine the expected bacterial species in each sample, the results of the whole genome
536    sequencing on the isolates cultured from the corresponding subject was used.  For each

537  classification scheme, accuracy was calculated by enumerating the number of species identified
538  by WGS that were also identified by the V4 16S targeted amplicon sequencing using the custom
539  "real_world_data_2020-4-17.Rmd" file.

540  The results of the V4 validation set were evaluated according to the following definitions
541  (**Figure 9**):

542      ***True match*** - All matches from the computational classification scheme that were
543          correctly identified by V4 16S targeted amplicon sequencing
544      ***False match*** - All species identified by V4 16S targeted amplicon sequencing that were
545          not identified by the computational classification scheme
546      ***False non-match*** - All matches from the computational classification scheme that were
547          not identified by V4 16S targeted amplicon sequencing
548      ***True non-match*** - All species that were not identified by either the computational
549          classification schemes or the V4 16S targeted amplicon sequencing



**Figure 9: Definitions of how the classification scheme outcomes are assigned to the cells of the confusion matrix for the V4 validation results.** This example shows the classification scheme composed of the Greengenes database, BLCA classifier, and the V4 region of the 16S rRNA gene as the identifier. When the Thomas-White dataset is subsetted by the 24 samples that underwent targeted amplicon sequencing, a smaller set of 49 species remains. The light yellow rows indicate the species correctly identified by the computational classification scheme. Blue dots represent species identified in the collected samples by whole genome sequencing after expanded urine culturing and isolation. Yellow dots indicate the species were identified in those samples by V4 targeted amplicon sequencing. Yellow dots in light yellow rows are true matches, when found elsewhere they are false matches. Blue dots in the light yellow rows are false non-matches, when found elsewhere they are true non-matches.

19

**Author contributions.**
CH and LK conceived and designed experiments, and wrote the manuscript. CH performed analyses, LK reviewed analyses. AJW supplied the data. LK, CH, MM, HS, TG, IF, NS, and AJW interpreted results and revised the manuscript.

# References

1. Lederberg J, McCray AT. 'Ome Sweet 'Omics-- A Genealogical Treasury of Words. The Scientist. 2001 Apr 2;15(7):2.

2. Wolfe AJ, Toh E, Shibata N, Rong R, Kenton K, FitzGerald M, et al. Evidence of Uncultivated Bacteria in the Adult Female Bladder. J Clin Microbiol. 2012 Apr 1;50(4):1376–83.

3. Khasriya R, Sathiananthamoorthy S, Ismail S, Kelsey M, Wilson M, Rohn JL, et al. Spectrum of Bacterial Colonization Associated with Urothelial Cells from Patients with Chronic Lower Urinary Tract Symptoms. J Clin Microbiol. 2013 Jul 1;51(7):2054–62.

4. Pearce MM, Hilt EE, Rosenfeld AB, Zilliox MJ, Thomas-White K, Fok C, et al. The Female Urinary Microbiome: a Comparison of Women with and without Urgency Urinary Incontinence. Blaser MJ, editor. mBio. 2014 Jul 8;5(4):e01283-14.

5. Fouts DE. Next Generation Sequencing to Define Prokaryotic and Fungal Diversity in the Bovine Rumen. PLOS ONE. 2012;7(11).

6. Price TK, Dune T, Hilt EE, Thomas-White KJ, Kliethermes S, Brincat C, et al. The Clinical Urine Culture: Enhanced Techniques Improve Detection of Clinically Relevant Microorganisms. Forbes BA, editor. J Clin Microbiol. 2016 May;54(5):1216–22.

7. Ackerman AL, Underhill DM. The mycobiome of the human urinary tract: potential roles for fungi in urology. Ann Transl Med. 2017 Jan;5:31–31.

8. Karstens L, Asquith M, Davin S, Stauffer P, Fair D, Gregory WT, et al. Does the Urinary Microbiome Play a Role in Urgency Urinary Incontinence and Its Severity? Front Cell Infect Microbiol. 2016 Jul 27;6.

9.  Pearce MM, Zilliox MJ, Rosenfeld AB, Thomas-White KJ, Richter HE, Nager CW, et al. The female urinary microbiome in urgency urinary incontinence. Am J Obstet Gynecol. 2015 Sep;213(3):347.e1-347.e11.

10. Thomas-White KJ, Gao X, Lin H, Fok CS, Ghanayem K, Mueller ER, et al. Urinary microbes and postoperative urinary tract infection risk in urogynecologic surgical patients. Int Urogynecology J. 2018 Dec;29(12):1797–805.

11. Thomas-White KJ, Hilt EE, Fok C, Pearce MM, Mueller ER, Kliethermes S, et al. Incontinence medication response relates to the female urinary microbiota. Int Urogynecology J. 2016 May;27(5):723–33.

12. Kliem M, Sauer S. The essence on mass spectrometry based microbial diagnostics. Curr Opin Microbiol. 2012 Jun;15(3):397–402.

13. Pace NR. A Molecular View of Microbial Diversity and the Biosphere. Science. 1997 May 2;276(5313):734–40.

14. Yang B, Wang Y, Qian P-Y. Sensitivity and correlation of hypervariable regions in 16S rRNA genes in phylogenetic analysis. BMC Bioinformatics. 2016 Dec;17(1):135.

15. Hugenholtz P, Skarshewski A, Parks DH. Genome-Based Microbial Taxonomy Coming of Age. Cold Spring Harb Perspect Biol. 2016 Jun;8(6).

16. Pruesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, et al. SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. Nucleic Acids Res. 2007 Nov 14;35(21):7188–96.

17. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, et al. Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. Nat Biotechnol. 2019 Aug;37(8):852–7.

18. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-resolution sample inference from Illumina amplicon data. Nat Methods. 2016 Jul;13(7):581–3.

19. Gao X, Lin H, Revanna K, Dong Q. A Bayesian taxonomic classification method for 16S rRNA gene sequences with improved species-level accuracy. BMC Bioinformatics. 2017 Dec;18(1):247.

20. Thomas-White K, Forster SC, Kumar N, Kuiken MV, Putonti C, Stares MD, et al. Culturing of female bladder bacteria reveals an interconnected urogenital microbiota. Nat Commun. 2018 Apr 19;9(1):1557.

21. Breitwieser FP, Lu J, Salzberg SL. A review of methods and databases for metagenomic classification and assembly. Brief Bioinform. 2019 Jul 19;20(4):1125–36.

21

625   22. Hugerth LW, Andersson AF. Analysing Microbial Community Composition through
626        Amplicon Sequencing: From Sampling to Hypothesis Testing. Front Microbiol.
627        2017;8.

628   23. DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, et al.
629        Greengenes, a Chimera-Checked 16S rRNA Gene Database and Workbench
630        Compatible with ARB. Appl Env Microbiol. 2006 Jul 1;72(7):5069–72.

631   24. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA
632        ribosomal RNA gene database project: improved data processing and web-based
633        tools. Nucleic Acids Res. 2012 Nov 27;41(D1):D590–6.

634   25. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian Classifier for Rapid
635        Assignment of rRNA Sequences into the New Bacterial Taxonomy. Appl Environ
636        Microbiol. 2007 Aug 15;73(16):5261–7.

637   26. Rennie JDM, Shih L, Teevan J, Karger DR. Tackling the Poor Assumptions of Naive
638        Bayes Text Classifiers. 2003;8.

639   27. Rainey FA. How to Describe New Species of Prokaryotes. In: Methods in
640        Microbiology. Elsevier; 2011. p. 7–14.

641   28. Parte AC. LPSN – List of Prokaryotic names with Standing in Nomenclature
642        (bacterio.net), 20 years on. Int J Syst Evol Microbiol. 2018 Jun 1;68(6):1825–9.

643   29. Notredame C, Higgins DG, Heringa J. T-coffee: a novel method for fast and
644        accurate multiple sequence alignment 1 1Edited by J. Thornton. J Mol Biol. 2000
645        Sep;302(1):205–17.

646   30. Valdar WSJ. Scoring residue conservation. Proteins Struct Funct Bioinforma.
647        2002;48(2):227–41.

648   31. Kumar PS, Griffen AL, Moeschberger ML, Leys EJ. Identification of Candidate
649        Periodontal Pathogens and Beneficial Species by Quantitative 16S Clonal Analysis.
650        J Clin Microbiol. 2005 Aug 1;43(8):3944–55.

651   32. Kumar P, Brooker M, Dowd S, Camerlengo T. Target Region Selection Is a Critical
652        Determinant of Community Fingerprints Generated by 16S Pyrosequencing. PLOS
653        ONE. 2011 Jun;6(6).

654   33. Bukin YuS, Galachyants YuP, Morozov IV, Bukin SV, Zakharenko AS, Zemskaya TI.
655        The effect of 16S rRNA region choice on bacterial community metabarcoding
656        results. Sci Data. 2019 Mar;6(1):190007.

657   34. Walter J, Tannock GW, Tilsala-Timisjarvi A, Rodtong S, Loach DM, Munro K, et al.
658        Detection and Identification of Gastrointestinal Lactobacillus Species by Using
659        Denaturing Gradient Gel Electrophoresis and Species-Specific PCR Primers. Appl
660        Environ Microbiol. 2000 Jan 1;66(1):297–303.

661  35. Graspeuntner S, Loeper N, Künzel S, Baines JF, Rupp J. Selection of validated
662      hypervariable regions is crucial in 16S-based microbiota studies of the female
663      genital tract. Sci Rep. 2018 Dec;8(1):9678.

664  36. Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA, Turnbaugh PJ,
665      et al. Global patterns of 16S rRNA diversity at a depth of millions of sequences per
666      sample. Proc Natl Acad Sci. 2011 Mar 15;108(Supplement_1):4516–22.

667  37. Hugerth LW, Wefer HA, Lundin S, Jakobsson HE, Lindberg M, Rodin S, et al.
668      DegePrime, a Program for Degenerate Primer Design for Broad-Taxonomic-Range
669      PCR in Microbial Ecology Studies. Löffler FE, editor. Appl Environ Microbiol. 2014
670      Aug 15;80(16):5116–23.

671  38. Morris EK, Caruso T, Buscot F, Fischer M, Hancock C, Maier TS, et al. Choosing
672      and using diversity indices: insights for ecological applications from the German
673      Biodiversity Exploratories. Ecol Evol. 2014 Sep;4(18):3514–24.

674

675

676

677

678

679

680

681

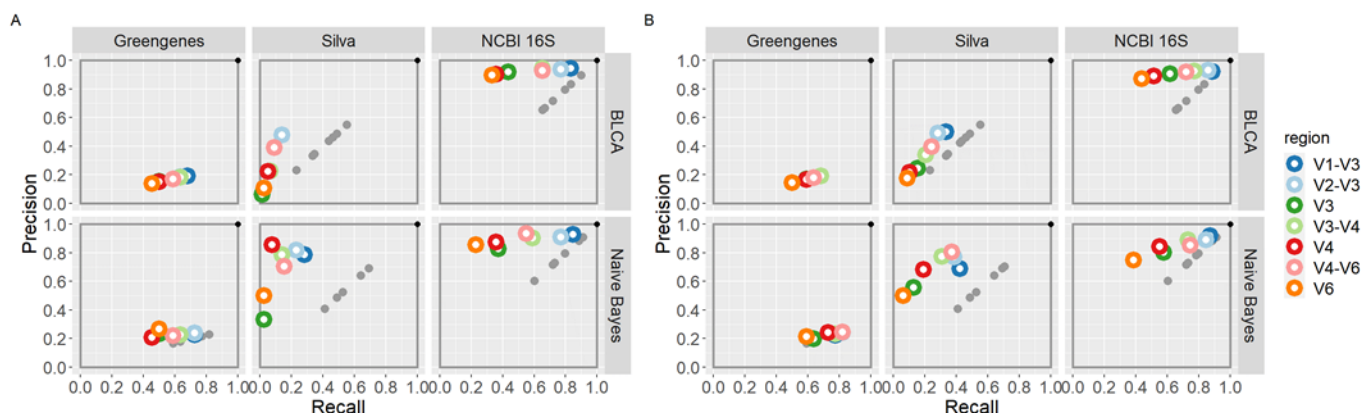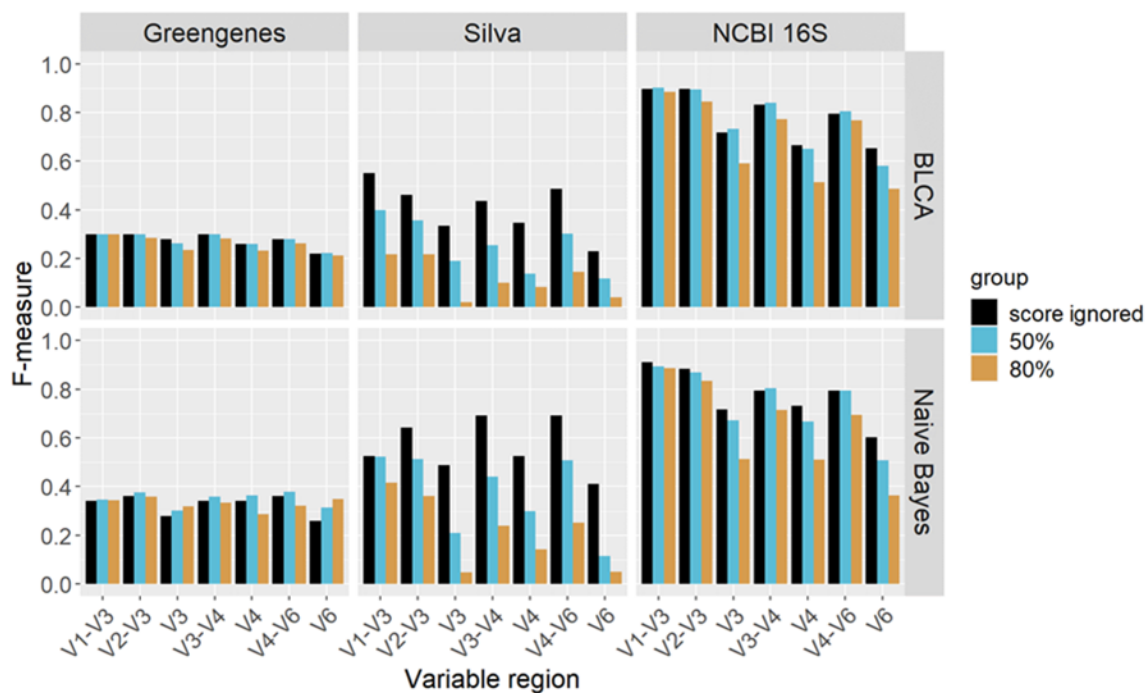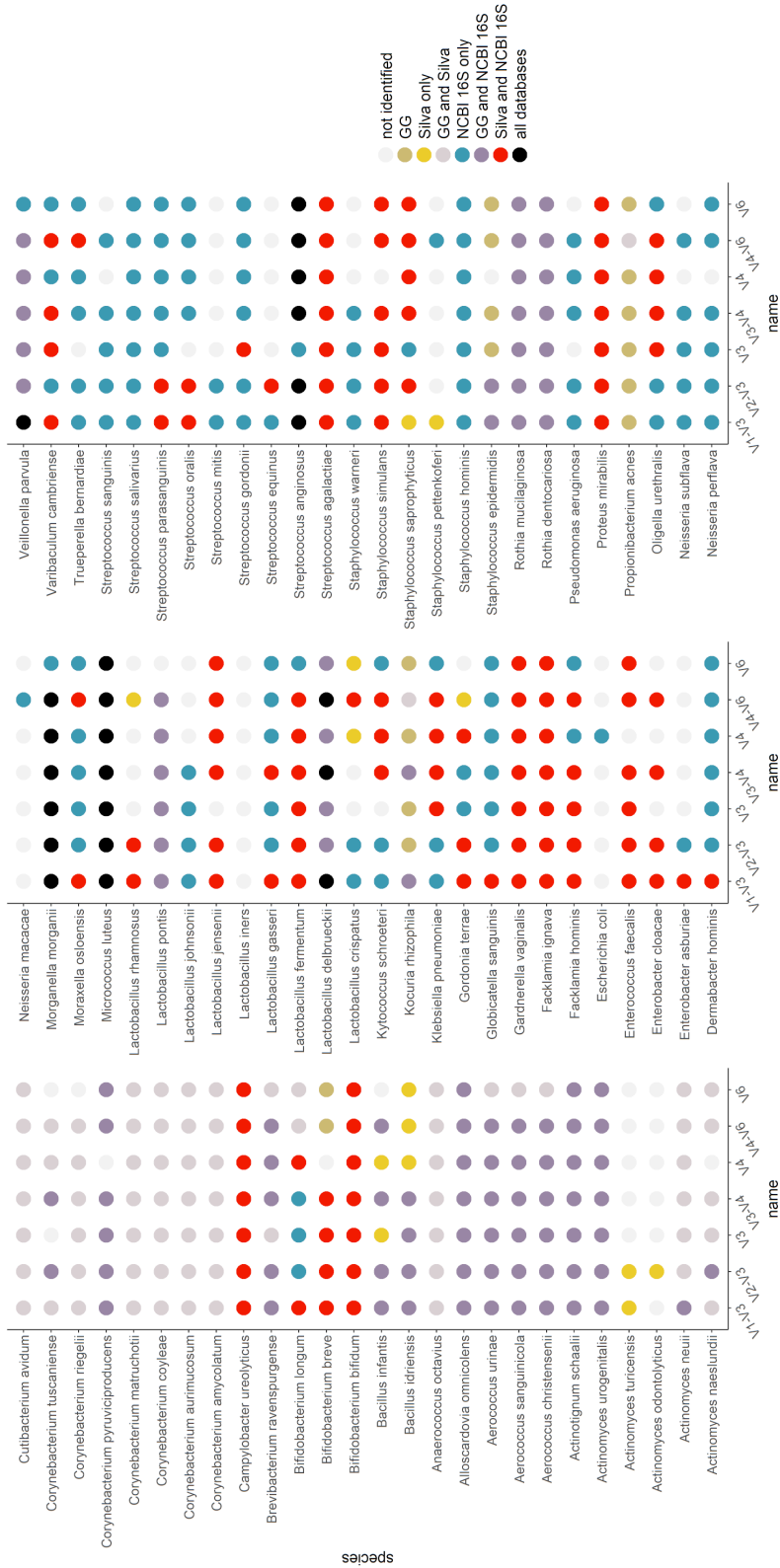682

683

684

685

686

687

688

## Supplemental Figures



**Supplemental Figure 1.** Classification scheme precision and recall when using a confidence scores of (A) 80% and (B) 50% as a threshold. The schemes that use the Silva database have very low recall compared to when the confidence score is ignored (gray dots), whereas schemes that use Greengenes and NCBI 16S are not as affected.
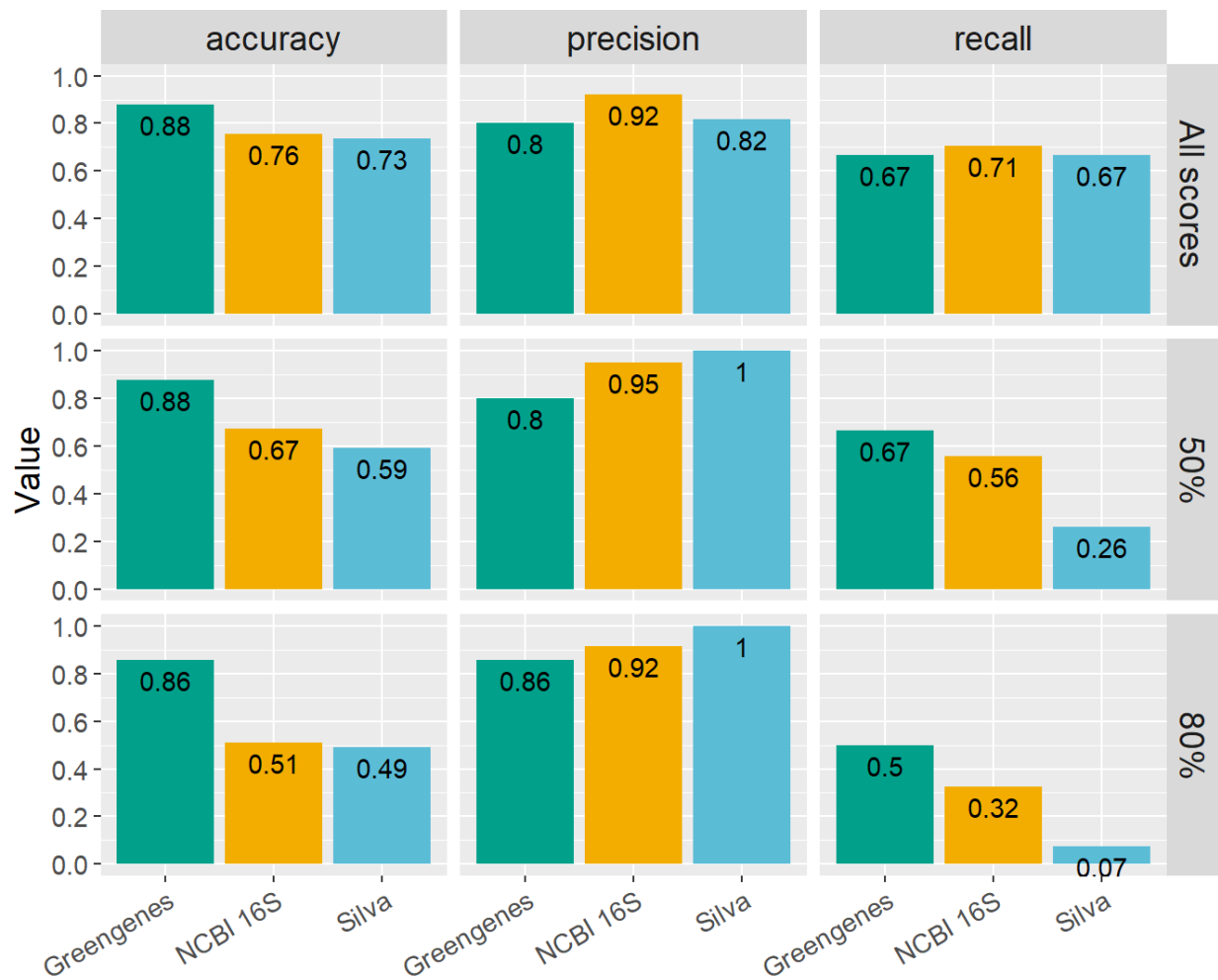


**Supplemental Figure 2.** F-measure values for all classification schemes. Values shown are for confidence scores of 50%, 80%, and when confidence scores are ignored.
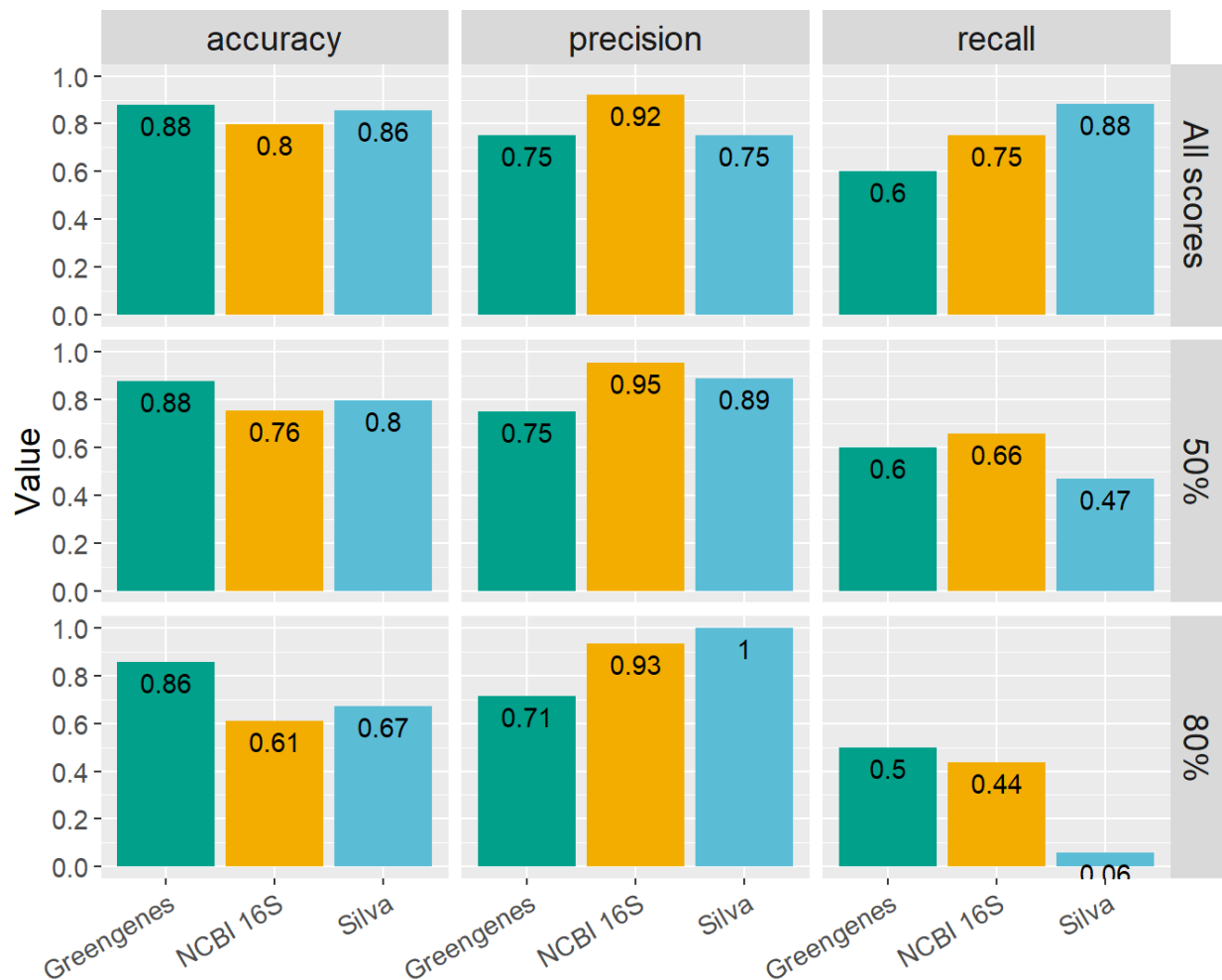
**Supplemental Figure 3.** Bladder bacterial species identified by database, variable region, and BLCA classifier.

25

703
704



705
706 **Supplemental Figure 4.** Values for accuracy, precision and recall (columns) when assigning
707 taxonomy with the V4 identifier, Naive Bayes classifier and all databases. Rows are values when
708 ignoring confidence scores, and when using confidence scores of 50% or 80% as thresholds.
709
710

26

711
712
713 **Supplemental Figure 5.** Values for accuracy, precision and recall (columns) when assigning
714 taxonomy with the V4 identifier, BLCA classifier and all databases. Rows are values when
715 ignoring confidence scores, and when using confidence scores of 50% or 80% as thresholds.