1    *Escherichia coli* clonobiome: assessing the strains diversity in feces and urine by deep amplicon

2    sequencing.

3

4    Sofiya G. Shevchenko[a], Matthew Radey[a], Veronika Tchesnokova[a], Dagmara Kisiela[a], Evgeni V.

5    Sokurenko[a#]

6

7    [a]Department of Microbiology, University of Washington, Seattle, WA, USA

8

9    Running Head: Deep amplicon seq for characterizing *E. coli* diversity

10

11    #Address correspondence to Evgeni V. Sokurenko, evs@u.washington.edu

12    Present address: Evgeni V. Sokurenko, Dept. of Microbiology, University of Washington,

13    Seattle, WA 98195-7735

14

15

16

17    Abstract word count: 232

18    Importance word count: 148

19    Article word count: 5,112

20

21

22    **ABSTRACT**

23    While microbiome studies have focused on diversity on the species or higher level, bacterial

24    species in microbiomes are represented by different, often multiple strains. These strains could

25    be clonally and phenotypically very different, making assessment of strain content vital to a full

26    understanding of microbiome function. This is especially important with respect to antibiotic

27    resistant strains, the clonal spread of which may be dependent on competition between them and

28    susceptible strains from the same species. The pandemic, multi-drug resistant, and highly

29    pathogenic *E. coli* subclone ST131-*H*30 (H30) is of special interest, as it has already been found

30    persisting in the gut and bladder of healthy people. In order to rapidly assess *E. coli* clonal

31    diversity, we developed a novel method based on deep sequencing of two loci used for sequence

32    typing, along with an algorithm for analysis of resulting data. Using this method, we assessed

33    fecal and urinary samples from healthy women carrying *H*30, and were able to uncover

34    considerable diversity, including strains with frequencies at <1% of the *E. coli* population. We

35    also found that even in the absence of antibiotic use, *H*30 could complete dominate the gut and,

36    especially, urine of healthy carriers. Our study offers a novel tool for assessing a species' clonal

37    diversity (clonobiome) within the microbiome, that could be useful in studying population

38    structure and dynamics of multi-drug resistant and/or highly pathogenic strains in their natural

39    environments.

40    **IMPORTANCE**

41    Bacterial species in the microbiome are often represented by multiple genetically and

42    phenotypically different strains, making insight into subspecies diversity critical to a full

43    understanding of the microbiome, especially with respect to opportunistic pathogens. However,

44    methods allowing efficient high-throughput clonal typing are not currently available. This study

45    combines a conventional *E. coli* typing method with deep amplicon sequencing to allow analysis

46    of many samples concurrently. While our method was developed for *E. coli*, it may be adapted

47    for other species, allowing for microbiome researchers to assess clonal strain diversity in natural

48    samples. Since assessment of subspecies diversity is particularly important for understanding the

49    spread of antibiotic resistance, we applied our method to study of a pandemic multidrug-resistant

50    *E. coli* clone. The results we present suggest that this clone could be highly competitive in

51    healthy carriers, and that the mechanisms of colonization by such clones need to be studied.

52    **INTRODUCTION**

53    Microbiomes, both in terms of function and diversity, have recently been a topic of considerable

54    interest. The gut microbiome has gotten special attention due to its high complexity and

55    importance to health[1-9]. So far, studies have almost exclusively focused on species or higher-

56    level diversity. However, this paints an incomplete picture, since strains within the same species

57    can be of distinct clonal origin and have vastly different metabolic, pathogenic, and antibiotic

58    resistance profiles[10-19]. Importantly, multidrug-resistant bacterial strains have been found

59    competing with commensal strains in the gut, even without antibiotic pressure[18-23]. Thus, there is

60    a pressing need to identify strains in the human microbiome for species of critical health

61    importance.

62    *Escherichia coli* is one of the most common residents of the gut. While primarily a commensal

63    colonizer, extra-intestinal pathogenic *E. coli* clones are implicated in a variety of diseases,

3

64    including urinary tract infections (UTIs) - a leading cause of human antibiotic use[24-28]. The

65    spread of multi-drug resistant *E. coli* is now a major health concern, especially the pandemic

66    *fimH*30 subclone of sequence type ST131 (*H*30). Though recently-emerged, *H*30 is now globally

67    distributed and comprises up to half of all urinary and bloodstream isolates of *E. coli* that are

68    fluoroquinolone-resistant and produce extended-spectrum beta-lactamases (ESBL)[29-33].

69    Additionally, it is strongly associated with drug-bug mismatches and adverse outcomes in elderly

70    and immunocompromised individuals[31-34]. Somewhat paradoxically, *H*30 is also a persistent gut

71    colonizer of healthy people and frequently causes asymptomatic bacteriuria (ABU) in such

72    carriers[35]. Yet, the relative clonal predominance of *H*30 strains among *E. coli* colonizing the gut

73    or bladder in healthy carriers remains unknown. Answering these questions could have a

74    significant impact on understanding the spread of antibiotic resistance and its reservoirs.

75    Currently, microbiome diversity is studied by sequencing the 16S rRNA gene, but this cannot

76    capture clonal diversity[36, 37]. Conventional methods for assessing clonal diversity, such as

77    metagenomic sequencing and single colony typing, are costly and labor intensive. For reliable

78    clonal diversity analysis, metagenomic sequencing requires very high coverage per sample,

79    while single colony typing requires handpicking large numbers of colonies for multi-locus

80    sequence typing (MLST)[38-42.] In *E. coli*, MLST requires assessment of 7 genes per isolate which

81    is analytically complex, costly, labor intensive, and therefore difficult to implement. Previously,

82    we reported an alternative clonotyping method that requires sequencing regions of only 2 genes –

83    *fumC* which is part of the MLST scheme and *fimH* that encodes a rapidly-evolving fimbrial

84    adhesin[43]. The *fumC/fimH*-based (CH) typing of *E. coli* is widely accepted due to its simplicity

85    and ability to not only identify specific STs but subdivide them into smaller subclones[43].

4

86    Specifically, *H*30 is identified using the allele combination *fumC*40*/fimH*30, while other less

87    resistant ST131 strains have the same *fumC* but different *fimH* alleles.

88    Here, we report a high-throughput method for clonal typing of *E. coli* strains by combining CH

89    typing and deep amplicon sequencing. We developed a new algorithm - Population-Level Allele

90    Profiler (PLAP) - for detecting alleles and predicting the relative prevalence of each allele in a

91    sample. We were able to assess the prevalence of clonal groups (including *H*30) in multiple fecal

92    and urine samples concurrently, with a limit of relative abundance detection at <1% of the total

93    population.

94    **RESULTS**

95    **Deep amplicon sequencing of defined samples**

96    To validate our approach and establish a limit of detection for strain presence, we first tested our

97    deep amplicon sequencing procedure on a set of defined samples. To create the defined samples,

98    we first selected a fecal sample from our lab collection known to contain *H*30 and ST101. Next,

99    we isolated a single colony from each and confirmed them to be strains of *H*30 (*fumC*40*/fimH*30)

100   and ST101 (*fumC*41*/fimH*86) using CH typing. From these single colonies, we first created *H*30-

101   only and ST101-only mixtures of *fumC* and *fimH* amplicons. We also created four ST101/H30

102   mixed samples by combining the *fumC* and *fimH* amplicons from ST101 and H30 in ST101:*H*30

103   ratios of 1:1, 1:4, 1:100, and 1:1000.

104   Analysis of raw sequencing data from *H*30-only and ST101-only samples showed the average

105   coverage of erroneous bases was 0.08% ± 0.09% for both strains. Erroneous bases were observed

106   in both genes across most nucleotide positions. The highest coverage for an erroneous base was

107   0.66% of aligned reads in *fumC* and 0.45% in *fimH* for *H*30, and 0.68% in *fumC* and 0.46% of

5

108    reads in *fimH* for ST101. The frequency distribution for erroneous base coverage is presented in

109    Supplemental Figure 1.

110    Analysis of raw sequencing data from ST101/*H*30 mixes showed that both *H*30 and ST101

111    alleles were detectable in the 1:1, 1:4, and 1:100 mixes. In the 1:1000 mix, only alleles of the

112    dominant *H*30 strain were observed. In the 1:1, 1:4, and 1:100 mixes, input and observed allele

113    prevalence was highly correlated for both *fumC* and *fimH* ($R^2$=0.996 and 0.997 respectively,

114    Suppl. Fig. 2). Erroneous bases were observed at 0.09% ± 0.1% and 0.08% ± 0.09% of aligned

115    reads in *fumC* and *fimH*, respectively (Suppl. Fig. 1). The highest coverage for erroneous bases

116    among all mixes was 0.79% of aligned reads for *fumC* and 0.57% of aligned reads for *fimH*.

117    Since 0.79% of aligned reads was the highest coverage for an erroneous base, we established

118    0.8% as a cutoff for correct base calling in both genes. This cutoff was used for all further PLAP

119    analysis.

120    **Deep sequencing of study samples and allele prediction**

121    Next, we applied PLAP to 67 participant samples (43 fecal and 24 urine) collected from a

122    previous study[35]. A total of 128 *fumC* and 129 *fimH* alleles were predicted across all samples, of

123    which 123 (96.1%) and 125 (96.9%) were previously known *fumC* and *fimH* alleles,

124    respectively. 5 novel *fumC* and 4 novel *fimH* alleles were potentially detected. All novel *fumC*

125    and *fimH* alleles were phylogenetically distant from other alleles predicted in the sample,

126    indicating that these alleles are not artifacts of sequencing (Suppl. Fig. 3, 4). These novel alleles

127    nonetheless clustered with other *E. coli fumC* and *fimH* alleles, indicating that these are novel *E.*

128    *coli* alleles rather than alleles belonging to other species.

6

129    The average number of alleles predicted per sample was 1.91 ± 0.96 for *fumC* and 1.93 ± 1.01

130    for *fimH*. 43 samples had same numbers of predicted *fumC* and *fimH* alleles; 24 samples had

131    different numbers of predicted *fumC* and *fimH* alleles (Fig. 1). Overall, the number of predicted

132    *fumC* alleles correlated to the number of predicted *fimH* alleles with an $R^2$ of 0.88 (Fig. 1).

133    To assess the performance of PLAP for predicting alleles, we used samples containing criterion

134    clones - strains previously identified by single colony typing. PLAP detected criterion *fimH* and

135    *fumC* alleles in 52 of these samples (90%). In the 6 samples where criterion allele(s) were not

136    found, the criterion clones were ciprofloxacin-resistant, but their isolation from the sample

137    required ≥2 plating attempts. This leads us to believe that these alleles were not detected because

138    they were absent in the MacConkey-plated population prior to deep sequencing.

139    A total of 72 non-criterion (previously unidentified) *fumC* and 71 non-criterion *fimH* alleles were

140    predicted by PLAP across all 67 samples. To assess the performance of PLAP on non-criterion

141    alleles, we analyzed 14 samples (10 fecal, 4 urine) predicted to contain 22 non-criterion *fumC*

142    and 22 non-criterion *fimH* alleles. 12 of these samples had at least one non-criterion allele

143    alongside criterion alleles; the remaining 2 had multiple non-criterion alleles in each gene only.

144    For each sample ≥40 single colonies were isolated and CH type determined using 7-SNP qPCR,

145    with each CH type verified by sequencing. With these data, we confirmed 19 (86%) predicted

146    non-criterion alleles for each gene. This included one predicted novel *fumC* allele. Of the

147    unconfirmed alleles, one was not distinguishable by 7-SNP qPCR and had a predicted prevalence

148    of 1%; therefore, we did not attempt to locate it. The remaining unconfirmed alleles had

149    predicted prevalences of <3% and therefore may have been missed due to insufficient sampling.

150    Additionally, all criterion alleles in these samples, 12 per gene, were predicted by PLAP.

151    **Prediction of allele prevalence in multi-allele samples**

7

152    We have also designed PLAP to predict the within-sample prevalence of each allele. The average

153    allele prevalence in fecal samples was 47.3% ± 4.3% SEM (range 0.88 – 100%) for *fumC* and

154    48.4% ± 4.22% SEM (range 1 – 100%) in *fimH*. The average allele prevalence in urine samples

155    was 64.8% ± 6.91% SEM (range 1.4 – 100%) for *fumC* and 58.3% ± 7.18% SEM (range 1 –

156    100%) in *fimH*.

157    In order to verify that the prevalences predicted by PLAP were accurate, we compared

158    predictions to actual in-sample prevalence using two different methods.

159    In the first method, we used *H*30 since ascertaining its prevalence is relatively simple. By plating

160    the sample on MacConkey agar then patching onto LB-ciprofloxacin, it is possible to compare

161    the number of cipro-resistant (*H*30) colonies to the total number of *E. coli* colonies. The ratio of

162    these two numbers provides the *H*30 load in a sample. We compared the predicted prevalences of

163    *fumC*40 and *fimH*30 to the *H*30 load in 17 fecal samples containing cipro-resistant *H*30.

164    Correlations between the *H*30 load and the predicted prevalence of *fumC*40 and *fimH*30 were

165    0.86 and 0.84 respectively (Fig. 2), indicating that prevalences given by PLAP were

166    representative of actual allele prevalences. To determine whether outliers were present, we

167    calculated the 99% CI range for every sample (see Methods). Three outlier samples were

168    identified (open circles, Fig. 2). Since it is possible that these outliers contain ciprofloxacin-

169    sensitive non-*H*30 *fimH*30-containing clones, *fumC*-null or *fimH*-null clones, and/or

170    ciprofloxacin-sensitive *H*30, we decided to employ screening of a large number of single

171    colonies.

172    In this second method, we used single colony typing for the in-depth characterization of 14

173    multi-allele samples described above, alongside 4 additional single-allele samples (2 fecal, 2

174    urine) for which only one allele per gene was predicted. This set of 18 samples included 11 of

8

175     the 17 fecal samples used for the *H*30-based analysis above, including one of the outlier samples.

176     For all 18 samples, we used CH typing of ≥40 single colonies per sample to determine the

177     prevalence of each *fumC* and *fimH* allele. Correlation between the PLAP-predicted prevalence

178     and the experimental allele prevalence was 0.98 for both *fumC* and *fimH* alleles (Fig. 3). As in

179     the *H*30 analysis above, we determined whether outliers were present using the 99% CI range for

180     every sample. Only one outlier was detected, corresponding to the only sample that contained

181     colonies from which *fimH* could not be amplified (*fimH*-null colonies). Furthermore, the sample

182     that was an outlier in the *H*30-based analysis was found to contain a relatively rare ciprofloxacin-

183     sensitive *H*30.

**Matching *fumC* and *fimH* alleles to predict sample strain content**

185     In CH typing, unique combinations of *fumC* and *fimH* alleles are used to determine the identities

186     of strains in a sample. Since a strain contains one copy of *fumC* and *fimH*, the prevalences of

187     alleles of these two genes in the sequencing data should be identical. For example, in a sample

188     containing 30% *H*30 (*fumC*40/*fimH*30) and 70% ST101 (*fumC*41/*fimH*86), we expect to see 30%

189     of *fumC* reads to be *fumC*40 and 30% of *fimH* reads to be *fimH*30. In reality, however, the

190     prevalences will be slightly different due to PCR and sequencing errors. To establish an

191     acceptable difference between the prevalences of same-strain *fumC* and *fimH* alleles, we looked

192     at 11 samples containing unique CH types (i.e. without allele sharing). In these 11 samples, the

193     predicted prevalences of *fumC* and *fimH* were highly correlated (0.99, Fig. 3). First, we

194     calculated the absolute difference between the predicted *fumC* and *fimH* prevalence for each

195     matched pair of alleles. Next, each absolute difference was divided by the predicted *fumC* or

196     *fimH* prevalence to obtain a relative deviation (Fig. 4). Finally, we used the relative deviations to

9

197    derive an equation for the maximum acceptable difference between matching *fumC* and *fimH*

198    alleles (Fig. 4).

199    While some samples, like those discussed above, contain only unique CH types, others contain

200    CH types with shared alleles. For example, in a sample containing 30% *H*30 and 70% ST131,

201    which share *fumC*40, the prevalence of *fumC*40 is not representative of either *H*30 or ST131

202    prevalence. For such samples, the minority rule was applied to resolve the strain content. Thus,

203    under the minority rule, the percentage of *H*30 in the example above would be determined by

204    *fimH*30, rather than *fumC*40, since the *fimH*30 prevalence is smaller. We tested this approach on

205    both the *H*30 and the 18-sample analysis described above to see if this resolved outliers. In both

206    cases, using the minority rule removed outliers and improved the correlation between predicted

207    and experimental prevalence (Suppl. Fig. 5). Thus, we were able to assign strain content and

208    strain prevalence in all samples, including samples with allele sharing.

209    **Predicted strain diversity of fecal and urine samples**

210    Using the equation described above, we were able to classify all samples in our study into 4

211    categories (see Fig. 5): samples with only one CH type (uniclonal); samples with multiple unique

212    CH types (unambiguous); samples with one dominant unique CH type and multiple minor non-

213    unique CH types (ambiguous-simple), and samples where the dominant CH type was not unique

214    (ambiguous-complex). Fecal samples were 33% uniclonal, 23% unambiguous, 21% ambiguous-

215    simple, and 23% ambiguous-complex. Urine samples were 54% uniclonal, 8% unambiguous,

216    25% ambiguous-simple, and 12.5% ambiguous-complex.

217    Overall, 107 fecal and 48 urine strains were predicted, corresponding to 68 clones in fecal

218    samples and 33 clones in urine samples. Of these clones, 50 (73.5%) and 24 (73%) were found in

219    Enterobase, respectively.

220    Out of the 155 total strains predicted, 6 were *fumC*-null (3.9%) and 2 were *fimH*-null (1.3%).

221    This is congruent with the occurrence of null alleles in our 18-sample subset, where 1 (3%) out

222    of 35 total strains predicted was a null-allele strain.

223    The average number of strains per sample was $2.47 \pm 1.32$ for fecal samples and $1.96 \pm 1.40$ for

224    urine samples. Based on Enterobase's ST-phylogroup data, we determined that B2 was the most

225    common (14 out of 47, 30%) among non-criterion fecal strains. Other phylogroups included A

226    (26%), B1 (19%), C (8.5%), D (11%), E (2%), and F (4%). Non-criterion strains in urine

227    samples included strains from phylogroups B2 (8 out of 16, 50%), B1 (19%), D (19%), A and F

228    (6% each).

229    **Novel clones**

230    17 fecal samples (40%) and 8 urine samples (33%) in our study were found to contain at least

231    one novel CH type. This included 19 fecal and 9 urine CH types not found in Enterobase. Of

232    these, 5 fecal and 3 urine CH types included at least one novel allele, and 14 fecal and 6 urine

233    CH types were combinations of *fumC* and *fimH* that were not previously observed (novel CH

234    combinations). Both CH types involving novel alleles and novel CH combinations were

235    observed to be primarily low-frequency clones. The average predicted prevalence for novel CH

236    combinations was $8.7\% \pm 3.5\%$ SEM (range 1-64.2%), and 13 out of 20 novel CH combinations

237    had predicted prevalences of <5%. One such combination was confirmed in our 14 characterized

11

238    sample set, consisting of *fumC*24 and *fimH*9, with a predicted prevalence of 1.6% and

239    experimental prevalence of 1.2%.

240    Similarly, 7 out of 8 novel allele-containing CH types had predicted prevalences of <2%. The

241    remaining CH type had a predicted prevalence of 70.7% and was detected using single colony

242    typing. The novel *fumC* allele was paired with *fimH*47 and was verified to be 8 SNPs away from

243    the closest known allele. The remaining MLST gene alleles for this strain were *adk*46, *icd*260,

244    *mdh*160, *gyrB*266, *purA*1, and *recA*221.

**Clones below error threshold**

246    To ascertain if we could identify alleles at prevalences below our defined error threshold of

247    0.8%, we ran PLAP on the set of 14 multi-allele samples using an error threshold of 0.5%. In 8

248    and 6 samples, respectively, prevalence of *fumC* and *fimH* alleles was <0.8%. None of the alleles

249    corresponded to known *fumC* or *fimH* alleles. These apparent novel alleles clustered alongside

250    known alleles identified in the sample (Suppl. Fig. 6, 7), leading us to conclude that these arose

251    due to sequencing or amplification error rather than belonging to clonally different strains.

**Predicted strain diversity in urine and fecal samples**

253    Strain diversity in first fecal samples was comparable with diversity in second fecal samples

254    (paired t-test, p>0.1). Distinguishing between *H*30-containing and non-*H*30 samples showed that

255    there was no statistical difference in strain diversity between *H*30-containing and non-*H*30 fecal

256    samples of either kind (unpaired t-test, p>0.1), and that there was no difference in diversity

257    between first and second fecal samples in either non-*H*30 or *H*30-containing samples (Fig. 6,

258    paired t-test, p>0.1). Both *H*30 and non-*H*30 urine samples were less diverse than corresponding

12

259     fecal samples (paired t-test, p<0.01 and 0.02, respectively). However, *H*30 urine samples were

260     less diverse than non-*H*30 urine samples (t-test, p=0.04).

261     It is also noteworthy that in 6 out of 23 *H*30-containing fecal samples, *H*30 was the only strain

262     predicted, indicating that it may be fully dominant in the gut niche in these participants.

263     **Strain turnover in fecal samples**

264     There was no correlation between number of strains in the first and second fecal sample, as well

265     as no correlation between number of strains in the urine sample and either fecal sample (Fig. 7).

266     When comparing the strain content of first and second fecal samples, we found that 92% of non-

267     criterion strains appeared to be transient i.e. were detected in one of the fecal samples only.

268     Transient non-criterion strains were also skewed towards lower-frequency strains (t-test,

269     p<0.001, Fig. 8B). It is possible that these strains are present in both fecal samples but are below

270     our limit of detection in one. However, we find that in one participant (P2, Suppl. Data) the first

271     fecal sample contains 3 ciprofloxacin-sensitive non-criterion strains while the second fecal

272     sample contains only ciprofloxacin-resistant *H*30 as verified by single colony testing. This leads

273     us to believe that there may be significant strain turnover in our fecal samples overall.

274     **DISCUSSION**

275     We combined conventional *fumC/fimH* typing with deep amplicon sequencing to assess *E. coli*

276     clonal diversity in a high-throughput manner. Our method has several advantages over existing

277     protocols. Firstly, our method has high sequencing resolution for target species. Since we only

278     sequence *E. coli fumC* and *fimH*, we can generate $\geq$0.5 million reads per sample, yielding $\geq$5,000

279     reads per base. In contrast, metagenomic sequencing, which is nonspecific to target species,

280     yields only 20 reads per base per genome (assuming a 5Mb genome). Secondly, our method

13

281    assessed up to 46 samples per sequencing run. In contrast, MLST requires typing ≥100 single

282    colonies per sample to capture the low-prevalence strains that PLAP detects. Finally, while we

283    developed PLAP for *E. coli*'s CH typing, PLAP is not limited to *E. coli* clonotyping and may be

284    generalized to other MLST schemes.

285    Despite studies showing that the healthy gut *E. coli* population typically includes multiple

286    clones, we show that the pandemic multidrug-resistant subclone *H*30 can dominate the gut in

287    healthy women, sometimes as the only detectable clone[42, 44-48]. This builds upon previous

288    research which has found multidrug-resistant bacteria in healthy people, and healthy people who

289    appear to harbor only one gut clone[44-48]. Total dominance is especially concerning since

290    antibiotic pressure was absent, indicating that *H*30 is potentially outcompeting other clones by

291    alternative means. Whether these mechanisms are metabolic, or whether certain virulence factors

292    give *H*30 an advantage is unclear, though previous studies have speculated that some virulence

293    factors may be beneficial for *E. coli* gut survival[49]. Additionally, our study involved a small

294    number of participants in which *H*30 was present in the gut and bladder. Therefore, it is possible

295    that host differences play a significant role. Another novel observation was that *H*30 was the sole

296    detected urinary strain more frequently than other clones, regardless of *H*30 gut dominance/non-

297    dominance. This may indicate that *H*30 might be an especially well-adapted uropathogen,

298    potentially explaining its association with UTI. Since it is unknown how ABU converts to UTI,

299    further study into *H*30 dominance in both ABU and UTI are needed.

300    We also uncovered substantial diversity in our samples. This includes significant *E. coli* diversity

301    in non-*H*30 urine samples from healthy women. Reports of multi-strain bacteriuria are rare,

302    likely due to the convention of selecting one isolate per urine sample[46, 47]. Therefore, it is

303    unknown how common multi-strain bacteriuria may truly be. Remarkably, we also detected low-

14

304    prevalence strains in the gut, some of which were novel clones, with up to 6 clones in a single

305    sample. Gut *E. coli* diversity of this magnitude is supported by studies typing >200 single

306    colonies per sample[42]. Studies using smaller counts usually report fewer clones, indicating that

307    there may be undescribed *E. coli* diversity when manageable numbers of colonies are used[44, 45].

308    Therefore, we believe that microbiome-like approaches to *E. coli* diversity are necessary to fully

309    understand intra-species dynamics in both the gut and bladder.

310    Our approach does have limitations. Firstly, our lowest detectable strain prevalence is 0.8% of

311    the *E. coli* population. This limit may be addressed in several ways including use of a high-

312    fidelity polymerase and preferential selection of *E. coli* colonies. However, we also recognize

313    that detection of rare strains may still prove difficult and that methods like ours may not fully

314    replace current techniques. Secondly, our method relies on sub-culturing *E. coli*. We are aware

315    that, theoretically, some strains could be suppressed during growth on selective media, forming

316    no/smaller colonies and skewing prevalence results. However, we did not encounter this during

317    our study. While amplification of *fumC* and *fimH* may be applied to urine samples without

318    culturing, attempts at doing this directly from fecal samples were unsuccessful, possibly due to

319    *E. coli* comprising <1% of the gut microbiome, making *E. coli* DNA too rare to effectively

320    amplify. Therefore, we used culturing for all samples. These issues lower the reliability of our

321    approach, but we believe that it remains an important step towards development of

322    comprehensive clonal diversity (clonobiome) assessment tools for any species of interest.

323    **MATERIALS AND METHODS**

324    **Study design and sample processing**

15

325   We selected a subset of participants from a previous study carried out by Kaiser Permanente

326   Washington and University of Washington (Seattle, WA)[35]. That study identified healthy gut

327   carriers of ciprofloxacin-resistant *E. coli*, including *E. coli H*30. These *E. coli* were found in

328   initial fecal samples by plating on LB-ciprofloxacin and CH typing of 1 to 8 single colonies.

329   After the initial fecal sample was analyzed, *H*30 carriers as well as carriers of some other strains

330   were asked to provide urine samples. These were received on average $152 \pm 55.9$ days after the

331   initial sample (85% responded). The respondents were then asked to provide follow-up fecal

332   samples, which were received on average $82 \pm 41.1$ days after the urine sample (84%

333   responded). All fecal and urine samples were tested for ciprofloxacin-resistant *E. coli* as with

334   initial samples. For this study, we chose 28 individuals who supplied all three samples. In 11

335   participants, *H*30 was identified in all three samples; in 4 additional participants *H*30 was

336   isolated in two samples. In 8 participants ciprofloxacin-resistant ST1193 was found in at least

337   two samples. In 5 participants the same ciprofloxacin-susceptible clone was found in at least two

338   samples. The sample types, strains clonal identity, and sampling times for all participants are

339   shown in Supplemental Figure 8. Average age of participants was $66.7 \pm 15.7$ years.

340   **Preparation of predefined control samples**

341   For control experiments, two predefined strains were chosen - *H*30 (*E. coli* FESS614.ds6) and

342   clonal group ST101 (*E. coli* FESS614.ds4). DNA from these strains was extracted and *fumC* and

343   *fimH* was amplified by PCR using the following conditions: 3min denaturation (95°C), 35 cycles

344   of annealing (95°C for 45sec, 57°C for 45sec, 72°C for 45sec), 5min extension (72°C), 4°C hold.

345   The primers (10 uM) used were as follows: 5'-TCACAGGTCGCCAGCGCTTC-3' (*fumC*

346   forward), 5'-GTACGCAGCGAAAAGATTC3' (*fumC* reverse), 5'-

347   TCAGGGAACCATTCAGGCA-3' (*fimH* forward), 5-ACAAAGGGCTAACGTGCAG-3' (*fimH*

16

348    reverse). Amount of PCR product was measured by Qbit. To create *H*30-only and ST101-only

349    samples, the corresponding *fumC* and *fimH* PCR products were pooled together at a 1:1 ratio. To

350    create mixes, *H*30 and ST101 amplicons of *fumC* were mixed together in ST101:*H*30 ratios of

351    1:1, 1:4, 1:10, 1:100, and 1:1000. The same was performed with *fimH* amplicons. The *fumC* and

352    *fimH* mixes were then pooled together by ratio type to create mixes that had equal concentrations

353    of total *fumC* and *fimH*. The DNA mixes were prepared for sequencing using Nextera XT DNA

354    library prep kit using standard protocol. The resulting library was sequenced on the Illumina

355    MiSeq (v3 kit). All mixes, except 1:10, reached coverage of ≥9,000X and were analyzed.

**Deep sequencing and allele analysis of the fecal and urine samples**

357    Each fecal and urine sample was plated on MacConkey agar to reach ~1,000 *E. coli* single

358    colonies per plate. All colonies were swabbed from the agar and DNA extracted using the

359    Qiagen Blood & Tissue Kit. From this pooled DNA *fumC* and *fimH* genes were amplified by

360    PCR by using the same primers and conditions as described above for control samples.

361    Amplicons were then purified and pooled by sample using the Qiagen Gel Extraction kit, then

362    prepared for sequencing using Nextera XT DNA library prep kit using standard protocol except

363    for usage of 52.5ul of RSB in the final magnetic bead cleanup step. The resulting library was

364    sequenced on the Illumina MiSeq (v3 kit). Sequencing data was analyzed using a Python

365    program of our construction, Population-Level Allele Profiler (PLAP), and has been made

366    available for public use on GitHub: github.com/marade/PLAP. The process is described below

367    (see also Suppl. Fig. 9).

368    For each sample, adapter sequences were removed using Trim-Galore, and resulting trimmed

369    reads were aligned to a list of all known *fumC* and *fimH* alleles using KMA with strict 99.99%

370    identity matching[50, 51]. For each KMA-detected allele per sample, trimmed reads were again

17

371    aligned to the sequence using Minimap2 and SAMtools[52, 53]. Any candidate allele which had at

372    least 1 base supported by <0.8% of reads was removed from consideration. False positives were

373    filtered using a moving 10bp window for each allele as follows. Reads of ≥100bp with 100%

374    identity within the window were counted. Alleles with low initial coverage, unstable coverage

375    (high average deviation from the mean), and high similarity in coverage pattern to an allele with

376    more stable coverage were removed from consideration. If >3 alleles were left for consideration

377    for a gene, 10bp moving window analysis was repeated with ≥200bp reads. If for any interval in

378    this second analysis, >60% of coverage was lost compared to the first moving window coverage,

379    the allele was discarded. Heterogeneity at any positions that remained undescribed by surviving

380    alleles was recorded. Relative abundance of all alleles was determined using the minimum

381    coverage found during first moving window analysis. In samples found by PLAP to be ≥50%

382    made up of <100bp reads (overtagmented samples), allele prevalence was calculated manually

383    by ascertaining base(s) unique to each allele and using the coverage of these base(s) to calculate

384    prevalence.

385    Out of the 28 total sets of fecal and urine samples chosen for this study, at least one sample failed

386    PCR amplification or sequencing library prep in 4 sets and therefore all samples from these sets

387    were dropped. From the remaining 24 sets we were able to sequence *fumC* and *fimH* in all three

388    samples. Out of those, 67 (89%) samples – 22 first fecal, 24 urine, and 21 second fecal – reached

389    ≥9,000X coverage per gene and were included in the analysis.

390    **Determining within-sample clonal group breakdown**

391    Identity of strains present in a sample was determined by combining *fumC* and *fimH* allele

392    numbers and determining the ST type using Enterobase. In uniclonal and unambiguous samples,

393    every allele had one match supported by the equation for maximum acceptable difference

18

394   between same-strain *fumC* and *fimH*. Therefore, these alleles formed a CH type based on which

395   ST type was determined.

396   For ambiguous-simple samples, the most prevalent *fumC* and *fimH* alleles formed an equation-

397   supported CH type. Any alleles that also had a single equation-supported match were assigned to

398   form a CH type. For all other alleles, Enterobase was consulted to determine which allele

399   combinations have been observed. If the CH type(s) produced was between alleles that had

400   different prevalences according to the equation, the "remaining" prevalence was calculated for

401   the allele with the greater prevalence. This allele was then paired with allele(s) for which an

402   Enterobase-logged CH type was not available and/or any novel alleles until the "remaining"

403   prevalence was consumed. If there were any allele(s) that remained after this step, they were

404   paired with the major allele of the opposite gene.

405   For ambiguous-complex samples, the most prevalent *fumC* and most prevalent *fimH* allele were

406   assigned to the same CH type. The "remaining" prevalence was calculated for the allele with the

407   greater prevalence and treated as an unmatched allele. From this step, we proceeded as with

408   ambiguous-simple samples.

409   **Determining prevalence of clonal groups by culturing**

410   Prevalence of ciprofloxacin-resistant clones in each sample was determined by diluting ~1ul of

411   sample with ≥300ul of $H_2O$, plating 40ul of this dilution on MacConkey agar, picking >130

412   single *E. coli* colonies, patching on Hardy-Chrom UTI agar to verify *E. coli* identity, then

413   patching colonies on LB-ciprofloxacin. Prevalence of other clonal groups was validated by

414   plating on MacConkey agar and subsequent patching of single colonies onto Hardy-Chrom UTI

19

415    agar to distinguish *E. coli*. *fumC* and *fimH* alleles of these colonies were then determined by 7-

416    SNP clonotyping and Sanger sequencing[54].

417    **Statistical and phylogenetic analysis**

418    To determine the 99% confidence interval (CI) for the prevalence of ciprofloxacin-resistant

419    strains, the number of resistant colonies was treated as number of successes and the total number

420    of picked colonies was treated as the total. To determine the 99% CI for the prevalence of

421    ciprofloxacin-sensitive strains, the number of colonies of that strain was treated as number of

422    successes and the total number of picked colonies was treated as the total. Confidence intervals

423    were calculated using Stata[55]. All t-tests were run using GraphPad

424    (http://www.graphpad.com/quickcalcs/ConfInterval1.cfm).

425    Phylogenetic trees were constructed using MEGA7[56]. Erroneous base coverage graph was

426    generated using seaborn[57]. *Escherichia coli fumC* alleles were downloaded from Enterobase

427    MLST allele data. *Escherichia coli fimH* alleles used are publicly available[58]. *Escherichia*

428    *fergusonii* and *albertii fumC* alleles were downloaded from NCBI. *Klebsiella pneumonia* and

429    *Enterobacter aerogenes* alleles of *fimH* were downloaded from the PATRIC database

430    (www.patricbrc.org).

431    **ACKNOWLEDGEMENTS**

436    E.V.S. conceived the project and designed the experiments. D.K. performed control sample

437    sequencing and analysis. All other sequencing, validation, and analysis was performed by S.G.S.

438    V.T. provided study data and samples. M.R. programmed the algorithm; M.R. and S.G.S. tested

439    and calibrated it. S.G.S. and E.V.S. wrote the manuscript with input from all authors.

440    **REFERENCES**

441    1.  Heintz-Buschart A, Wilmes P. 2018. Human gut microbiome: Function matters. Trends

442        Microbiol. 26(7):563-574.

443    2.  Caputi V, Giron MC. 2018. Microbiome-gut-brain axis and Toll-like receptors in

444        Parkinson's Disease. Int J Mol Sci 19(6):1689.

445    3.  Perez-Pardo P, Hartog M, Garssen J, Kraneveld AD. 2017. Microbes tickling your

446        tummy: the importance of the gut-brain axis in Parkinson's Disease. Curr Behav

447        Neurosci Rep 4(4):361-368.

448    4.  Sanmiguel C, Gupta A, Mayer EA. 2015. Gut Microbiome and obesity: A plausible

449        explanation for obesity. Curr Obes Rep 4(2):250-261.

450    5.  De la Cuesta-Zuluaga J, Corrales-Agudelo V, Velásquez-Mejía EP, Carmona JA, Abad

451        JM, Escobar JS. 2018. Gut microbiota is associated with obesity and cardiometabolic

452        disease in a population in the midst of Westernization. Sci Rep 8:11356.

453    6.  Roszyk E, Puszczewicz M. 2017. Role of human microbiome and selected bacterial

454        infections in the pathogenesis of rheumatoid arthritis. Reumatologia 55(5):242-250.

455    7.  Bu J, Wang Z. 2018. Cross-talk between gut microbiota and heart via the routes of

456        metabolite and immunity. Gastroenterol Res Pract 2018:6458094.

457    8.  Dzidic M, Boix-Amorós A, Selma-Royo M, Mira A, Collado MC. 2018. Gut microbiota

458        and mucosal immunity in the neonate. Med Sci Basel. 6(3): E56.

459    9.  Nunez G. 2017. Linking pathogen virulence, host immunity and the microbiota at the

460        intestinal barrier. Keio J Med 66(1):14.

461    10. Tenaillon O, Skurnik D, Picard B, Denamur E. 2010. The population genetics of

462        commensal *Escherichia coli*. Nature Reviews. 8(3):207-217.

463    11. Gordon DM, O'Brien CL, Pavli P. 2015. *Escherichia coli* diversity in the lower intestinal

464        tract of humans. Environ Microbiol Rep. 7(4):642-648.

465    12. Costea PI, Coelho LP, Sunagwa S, Much R, Huerta-Cepas J, Forslund K, Hildebrand F,

466        Kushugulova A, Zeller G, Bork P. 2017. Subspecies in the global human gut microbiome.

467        Mol Sys Biol. 13(12):960.

468    13. Metwaly A, Haller D. 2019. Strain-level diversity in the gut: the *P. copri* case. Cell Host

469        Microbe. 25(3):349-350.

470    14. Zhang C, Zhao L. 2016. Strain-level dissection of the contribution of the gut microbiome

471        to human metabolic disease. Genome Med. 8(1):41.

472    15. Leatham MP, Banerjee S, Autieri SM, Mercado-Lubo R, Conway T, Cohen PS. 2009.

473        Precolonized human commensal *Escherichia coli* clones serve as a barrier to *E. coli*

474        O157:H7 growth in the streptomycin-treated mouse intestine. Infect Immun. 77(7):2876-

475        86.

476    16. Hecht AL, Casterline BW, Earley ZM, Goo YA, Goodlett DR, Bubeck Wardenburg J.

477        2016. Clone competition restricts colonization of an enteric pathogen and prevents colitis.

478        EMBO Rep. 17(9):1281-91.

479    17. Lam LH, Monack DM. 2014. Intraspecies competition for niches in the distal gut dictate

480        transmission during persistent Salmonella infection. PLoS Pathog. 10(12):e1004527.

481    18. Sassone-Corsi M, Nuccio SP, Liu H, Hernandez D, Vu CT, Takahashi AA, Edwards RA,

482         Raffatellu M. 2016. Microcins mediate competition among Enterobacteriaceae in the

483         inflamed gut. Nature. 540(7632):280-283.

484    19. Moreno E, Johnson JR, Perez T, Prats G, Kuskowski MA, Andreu A. 2009. Structure and

485         urovirulence characteristics of the fecal *Escherichia coli* population among healthy

486         women. Microbes Infect. 11(2):274-280.

487    20. Bailey JK, Pinyon JL, Anantham S, Hall RM. 2010. Commensal *Escherichia coli* of

488         healthy humans: a reservoir for antibiotic-resistance determinants. J Med Microb.

489         59:1331-1339.

490    21. Gorrie CL, Mirceta M, Wick RR, Judd LM, Wyres KL, Thomson NR, Strugnell RA,

491         Pratt NF, Garlick JS, Watson KM, Hunter PC, McGloughlin SA, Spelman DW, Jenney

492         AWJ, Holt KE. 2018. Antimicrobial-resistant *Klebsiella pneumoniae* carriage and

493         infection in specialized geriatric care wards linked to acquisition in the referring hospital.

494         Clin Infect Dis. 67(2):161-170.

495    22. Li H, Zhu J. 2017. Targeted metabolic profiling rapidly differentiates *Escherichia coli*

496         and *Staphylococcus aureus* at species and strain level. Rapid Commun Mass Spectrom.

497         31(19):1669-1676.

498    23. Galardini M, Koumoutsi A, Herrera-Dominguez L, Cordero Varela JA, Telzerow A,

499         Wagih O, Wartel M, Clermont O, Denamur E, Typas A, Beltrao P. 2017. Phenotype

500         inference in an *Escherichia coli* strain panel. Elife. 6:e31035.

501    24. Bevan ER, McNally A, Thomas CM, Piddock LJV, Hawkey PM. 2018. Acquisition and

502         loss of CTX-M-producing and non-producing *Escherichia coli* in the fecal microbiome of

503         travelers to South Asia. mBio. 9(6):e02408-18.

504    25. Robin F1,2, Beyrouthy R3,2, Bonacorsi S4,5, Aissa N6, Bret L7, Brieu N8, Cattoir V9,

505        Chapuis A10, Chardon H8, Degand N11, Doucet-Populaire F12, Dubois V13, Fortineau

506        N14, Grillon A15, Lanotte P16, Leyssene D17, Patry I18, Podglajen I19, Recule C20,

507        Ros A21, Colomb-Cotinat M22, Ponties V22, Ploy MC23, Bonnet R3,2. 2017. Inventory

508        of extended-spectrum-β-lactamase-producing Enterobacteriaceae in France as assessed

509        by a multicenter study. Antimicrob Agents Chemother. 61(3): pii: e01911-16.

510    26. Gupta M, Didwal G, Bansal S, Kaushal K, Batra N, Gautam V, Ray P. 2019. Antibiotic-

511        resistant Enterobacteriaceae in healthy gut flora: A report from north Indian semiurban

512        community. Indian J Med Res. 149(2):276-280.

513    27. Johnson JR, Johnston B, Clabots C, Kuskowski MA, Castanheira M. 2010. *Escherichia

514        coli* sequence type ST131 as the major cause of serious multidrug-resistant *E. coli*

515        infections in the United States. Clin Infect Dis. 51(3):286-294.

516    28. Johnson JR, Tchesnokova V, Johnston B, Clabots C, Roberts PL, Billig M, Riddell K,

517        Rogers P, Qin X, Butler-Wu S, Price LB, Aziz M, Nicolas-Chanoine MH, Debroy C,

518        Robicsek A, Hansen G, Urban C, Platell J, Trott DJ, Zhanel G, Weissman SJ, Cookson

519        BT, Fang FC, Limaye AP, Scholes D, Chattopadhyay S, Hooper DC, Sokurenko EV.

520        2013. Abrupt emergence of a single dominant multidrug-resistant clone of *Escherichia

521        coli*. J Infect Dis. 207(6):919-928.

522    29. Burgess MJ, Johnson JR, Porter SB, Johnston B, Clabots C, Lahr BD, Uhl JR, Banerjee

523        R. 2015. Long-term care facilities are reservoirs for antimicrobial-resistant sequence type

524        131 *Escherichia coli*. Open Forum Infect Dis. 2(1):ofv011.

525    30. Johnson JR, Porter S, Thuras P, Castanheira M. 2017. The pandemic *H*30 subclone of

526        sequence type 131 (ST131) as the leading cause of multidrug-resistant *Escherichia coli*

527        infections in the United States (2011–2012). Open Forum Infect Dis. 4(2):ofx089.

528    31. Tchesnokova V, Rechkina E, Chan D, Haile HG, Larson L, Schroeder DW, Solyanik T,

529        Shibuya S, Hansen KE, Ralston JD, Riddell K, Scholes D, Sokurenko EV. 2019.

530        Pandemic uropathogenic fluoroquinolone-resistant *Escherichia coli* have enhanced

531        ability to persist in the gut and cause bacteriuria in healthy women. Clin Inf Dis.

532        (accepted)

533    32. Ong SH, Kukkillaya VU, Wilm A, Lay C, Ho EX, Low L, Hibberd ML, Nagarajan N.

534        2013. Species Identification and Profiling of Complex Microbial Communities Using

535        Shotgun Illumina Sequencing of 16S rRNA Amplicon Sequences. Parkinson J, ed. PLoS

536        One 8(4):e60811.

537    33. Chen W, Zhang CK, Cheng Y, Zhang S, Zhao H. 2013. A Comparison of Methods for

538        Clustering 16S rRNA Sequences into OTUs. Casiraghi M, ed. PLoS One. 8(8):e70837.

539    34. Zolfo M, Tett A, Jousson O, Donati C, Segata N. 2017. MetaMLST: multi-locus clone-

540        level bacterial typing from metagenomic samples. Nucleic Acids Res. 45(2):e7.

541    35. Scholz M, Ward DV, Pasolli E, Tolio T, Zolfo M, Asincar F, Truong DT, Tett A,

542        Morrow AL, Segata N. 2016. Clone-level microbial epidemiology and population

543        genomics from shotgun metagenomics. Nature Methods. 13:435-438.

544    36. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. 2016. An integrated

545        metagenomics pipeline for strain profiling reveals novel patterns of bacterial transmission

546        and biogeography. Genome Res. 26(11):1612-1625.

25

547    37. Fischer M, Strauch B, Renard BY. 2017. Abundance estimation and differential testing

548        on strain level in metagenomics data. Bioinformatics. 33(14):i124-i132.

549    38. Weissman SJ, Johnson JR, Tchesnokova V, Billig M, Dykhuizen D, Riddell K, Rogers P,

550        Qin X, Butler-Wu S, Cookson BT, Fang FC, Scholes D, Chattopadhyay S, Sokurenko

551        EV. 2012. High-resolution two-locus clonal typing of extraintestinal pathogenic

552        *Escherichia coli*. Appl Environ Microbiol. 78(5):1353-1360.

553    39. National Center for Emerging and Zoonotic Infectious Diseases, Division of Healthcare

554        Quality Promotion. "Biggest Threats and Data". Centers for Disease Control and

555        Prevention. www.cdc.gov/drugresistance/biggest_threats.html

556    40. Eckburg PB, Bik EM, Bernstein CN, Purdom E, Dethlefsen L, Sargent M, Gill SR,

557        Nelson KE, Relman DA. 2005. Diversity of the human intestinal microbial flora. Science.

558        308(5728):1635-8.

559    41. Anderson MA, Whitlock JE, Harwood VJ. 2006. Diversity and distribution of

560        *Escherichia coli* genotypes and antibiotic resistance phenotypes in feces of humans,

561        cattle, and horses. App Environ Microbiol. 72(11):6914-22.

562    42. Richter TKS, Hazen TH, Lam D, Coles CL, Seidman JC, You Y, Silbergeld EK, Fraser

563        CM, Rasko DA. 2018. Temporal variability of *Escherichia coli* diversity in the

564        gastrointestinal tracts of Tanzanian children with and without exposure to antibiotics.

565        mSphere. 3(6):e00558-18.

566    43. Diard M, Garry L, Selva M, Mosser T, Denamur E, Matic I. 2010. Pathogenicity-

567        associated islands in extraintestinal pathogenic *Escherichia coli* are fitness elements

568        involved in intestinal colonization. J Bacteriol. 192(19):4885-93.

569     44. Le Gall T, Clermont O, Gouriou S, Picard B, Nassif X, Denamur E, Tenaillon O. 2007.

570         Extraintestinal virulence is a coincidental by-product of commensalism in B2

571         phylogenetic group *Escherichia coli* strains. Mol Biol Evol. 24(11):2373-84.

572     45. Nielsen KL, Stegger M, Godfrey PA, Feldgarden M, Andersen PS, Frimodt-Moller N.

573         2016. Adaptation of *Escherichia coli* traversing from the faecal environment to the

574         urinary tract. Int J Med Microbiol. 306(8):595-603.

575     46. Moreno E, Andreu A, Perez T, Sabate M, Johnsom JR, Prats G. 2005. Relationship

576         between *Escherichia coli* strains causing urinary tract infection in women and the

577         dominant faecal flora of the same hosts. Epidemiol Infect. 134:1015-1023.

578     47. Smati M, Clermont O, Le Gal F, Schichmanoff O, Jauréguy F, Eddi A, Denamur E,

579         Picard B. 2013. Real-time PCR for quantitative analysis of human commensal

580         Escherichia coli populations reveals a high frequency of subdominant phylogroups. Appl

581         Environ Microbiol. 79(16):5005-12.

582     48. Krueger F. 2016. Trim Galore. https://github.com/FelixKrueger/TrimGalore. [Online;

583         accessed 2018-11-28]

584     49. Philip TLC, Clausen F, Aarestrup M, Lund O. 2018. Rapid and precise alignment of raw

585         reads against redundant databases with KMA", BMC Bioinformatics. 19:307.

586     50. Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics.

587         34:3094-3100.

588     51. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,

589         Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The Sequence

590         alignment/map (SAM) format and SAMtools, Bioinformatics. 25(16) 2078-9.

27

591    52. Tchesnokova V, Avagyan H, Billig M, Chattopadhyay S, Aprikian P, Chan D, Pseunova

592        J, Rechkina E, Riddell K, Scholes D, Fang FC, Johnson JR, Sokurenko EV. 2016. A

593        Novel 7-Single Nucleotide Polymorphism-Based Clonotyping Test Allows Rapid

594        Prediction of Antimicrobial Susceptibility of Extraintestinal *Escherichia coli* Directly

595        From Urine Specimens. Open Forum Infect Dis 3(1):ofw002.

596    53. StataCorp. 2019. Stata Statistical Software: Release 16. College Station, TX: StataCorp

597        LLC.

598    54. Kumar S, Stecher G, Tamura K. 2016. MEGA7: Molecular Evolutionary Genetics

599        Analysis version 7.0. Mol Biol Evol. 33(7):1870-1874.

600    55. Waskom M, Botvinnik O, O'Kane D, Hobson P, Lukauskas S, Gemperline DC,

601        Augspurger T, Halchenko Y, Cole JB, Warmenhoven J, de Ruiter J, Pye C, Hoyer S,

602        Vanderplas J, Villalba S, Kunter G, Quintero E, Bachant P, Martin M, Meyer K, Miles A,

603        Ram Y, Yarkoni T, Williams ML, Evans C, Fitzgerald C, Fonnesback C, Lee A, Qalieh

604        A. 2017. Seaborn: statistical data visualization. http://seaborn.pydata.org. [Online;

605        accessed 2019-02-05].

606    56. Roer L, Tchesnokova V, Allesoe R, Muradova M, Chattopadhyay S, Ahrenfeldt J,

607        Thomsen MCF, Lund O, Hansen F, Hammerum AM, Sokurenko E, Hasman H. 2017.

608        Development of a Web Tool for *Escherichia coli* Subtyping Based on *fimH* Alleles. J

609        Clin Microbiol. 55:2538–2543.

610    **FIGURE LEGENDS**

611    **Figure 1. Congruency of *fumC* and *fimH* allele counts in fecal and urine samples.** Size of

612    bubbles corresponds to number of samples with designated *fumC*/*fimH* allele counts (i.e. 1

28

613    sample with one *fumC* allele and three *fimH* alleles). Linear fit with Pearson square correlation

614    index shown.

615    **Figure 2**. **Validation of predicted *H*30 allele prevalence.** PLAP-predicted prevalence of *H*30

616    alleles vs actual *H*30 load in *H*30-containing fecal samples. Prevalence of predicted *fumC*40 (**A**)

617    and predicted *fimH*30 (**B**). Predicted prevalence of *fumC*40 and *fimH*30 is expressed as

618    percentage of all *E. coli* in each sample. Experimentally confirmed *H*30 load is expressed as

619    percent of *H*30 (ciprofloxacin-resistant) single colonies to all plated *E. coli* single colonies in

620    percent. At least 130 colonies were tested per sample. Outliers, marked in open circles, were

621    outside the 99% confidence interval of the number of colonies tested.

622    **Figure 3. Validation of predicted *fumC/fimH* allele prevalence. A.** PLAP-predicted vs

623    experimental within-sample *fumC/fimH* allele prevalence in 18 samples. Experimental allele

624    prevalence was determined by CH typing of at least 40 single bacterial colonies per sample.

625    Outliers (open circles) were outside the 99% confidence interval of the number of colonies

626    sampled. **B.** Predicted prevalence of *fumC* vs *fimH* alleles from the same CH type in 11 samples

627    where no sharing of alleles between strains was present.

628    **Figure 4.** Difference in predicted prevalence between *fumC* and *fimH* alleles from the same *E.*

629    *coli* strain. Deviation in absolute numbers is shown on the top. Deviation as a percentage of the

630    prevalence of the allele is shown on the bottom. Open circles indicate *fimH* data points. Shaded

631    circles indicate *fumC* data points. Trend lines and equations were used to determine intervals for

632    matching (i.e. belonging to the same CH type) *fumC* and *fimH* alleles.

29

633  **Figure 5**. Representative examples of each sample category defined by within-sample

634  breakdown of prevalence for *fumC* and *fimH* alleles. Number of fecal and urine samples

635  belonging in each category is listed below.

636  **Figure 6**. **Diversity of *E. coli* in individual fecal/urine samples.** *H*30 content was determined

637  by PLAP and/or culturing.

638  **Figure 7. Counts of *E. coli* strains in fecal and urine samples**. Number of strains detected by

639  PLAP in (**A**) first fecal vs urine, (**B**) second fecal vs urine, and (**C**) first fecal vs second fecal

640  samples. Each bubble indicates participants with the corresponding number of *E. coli* strains in

641  the designated sample. The bubble size indicates number of participants with the determined

642  number of strains. Linear fit with Pearson square correlation index shown.

643  **Figure 8. Persistence of *E. coli* strains in fecal samples.** (**A**) Prevalence of criterion fecal

644  strains in first vs second fecal samples. White data points indicate *H*30 strains while shaded data

645  points indicate non-*H*30 strains. Circled cluster represents 4 strains present at 100% prevalence

646  in both samples. Dotted lines indicate the mean prevalence for strains in first and second fecal

647  samples. Distribution of prevalences in both first and second fecal samples is not significantly

648  different from random (t-test, p>0.05). (**B**) Prevalence of non-criterion fecal strains in first vs

649  second fecal samples. Dotted lines indicate the mean prevalence for transient strains in first and

650  second fecal samples. Transient strains are defined as strains that are present in only one of the

651  two fecal samples from the same participant. Distribution of prevalences in both first and second

652  fecal samples is significantly skewed towards lower prevalences (t-test, p<0.01).

653  **Supplemental Figure 1**. Coverage of erroneous bases in *H*30-only, ST101-only, and mix sample

654  sequencing. Coverage is expressed in percentage of total reads aligned to each gene.

30

655     **Supplemental Figure 2**. Correlation between input and PLAP-derived (deep seq) prevalences of

656     *fumC* and *fimH* alleles of *H*30 and ST101 in 1:1, 1:4, and 1:100 mixes.

657     **Supplemental Figure 3**. Phylogenetic relationships between predicted novel *fumC* alleles and

658     known *E. coli fumC* alleles. *Escherichia fergusonii* and *albertii fumC* alleles also presented for

659     outgroup reference. Alleles not labelled with a species are known *E. coli* alleles or putative novel

660     alleles. Alleles found in the sample as the novel allele are highlighted in the same color as the

661     novel allele to show distance between predicted novel alleles and other *fumC* alleles present in

662     the sample. Alleles present in multiple different samples are marked with the appropriate colors

663     next to the allele name.

664     **Supplemental Figure 4**. Phylogenetic relationships between predicted novel *fimH* alleles and

665     known *E. coli fimH* alleles. *Klebsiella pneumoniae* and *Enterobacter aerogenes fimH* alleles also

666     presented for outgroup reference. Alleles not labelled with a species are known *E. coli* alleles or

667     putative novel alleles. Alleles found in the sample as the novel allele are highlighted in the same

668     color as the novel allele to show distance between predicted novel alleles and other *fimH* alleles

669     present in the sample. Alleles present in multiple different samples are marked with the

670     appropriate colors next to the allele name.

671     **Supplemental Figure 5**. **A**. Comparison of actual *H*30 load in *H*30-containing fecal samples to

672     PLAP-predicted *fumC*-40/*fimH*-30 prevalences with minority rule correction (i.e. the smaller

673     prevalence of the two was used). Prevalence of *fumC*-40/*fimH*-30 is expressed as percentage of

674     all *E. coli* in each sample. *H*30 load is expressed as ratio of *H*30 (ciprofloxacin-resistant) single

675     colonies to all plated *E. coli* single colonies in percent. **B**. PLAP-predicted allele prevalence

676     (with minority rule correction) compared to experimental allele prevalence as determined by

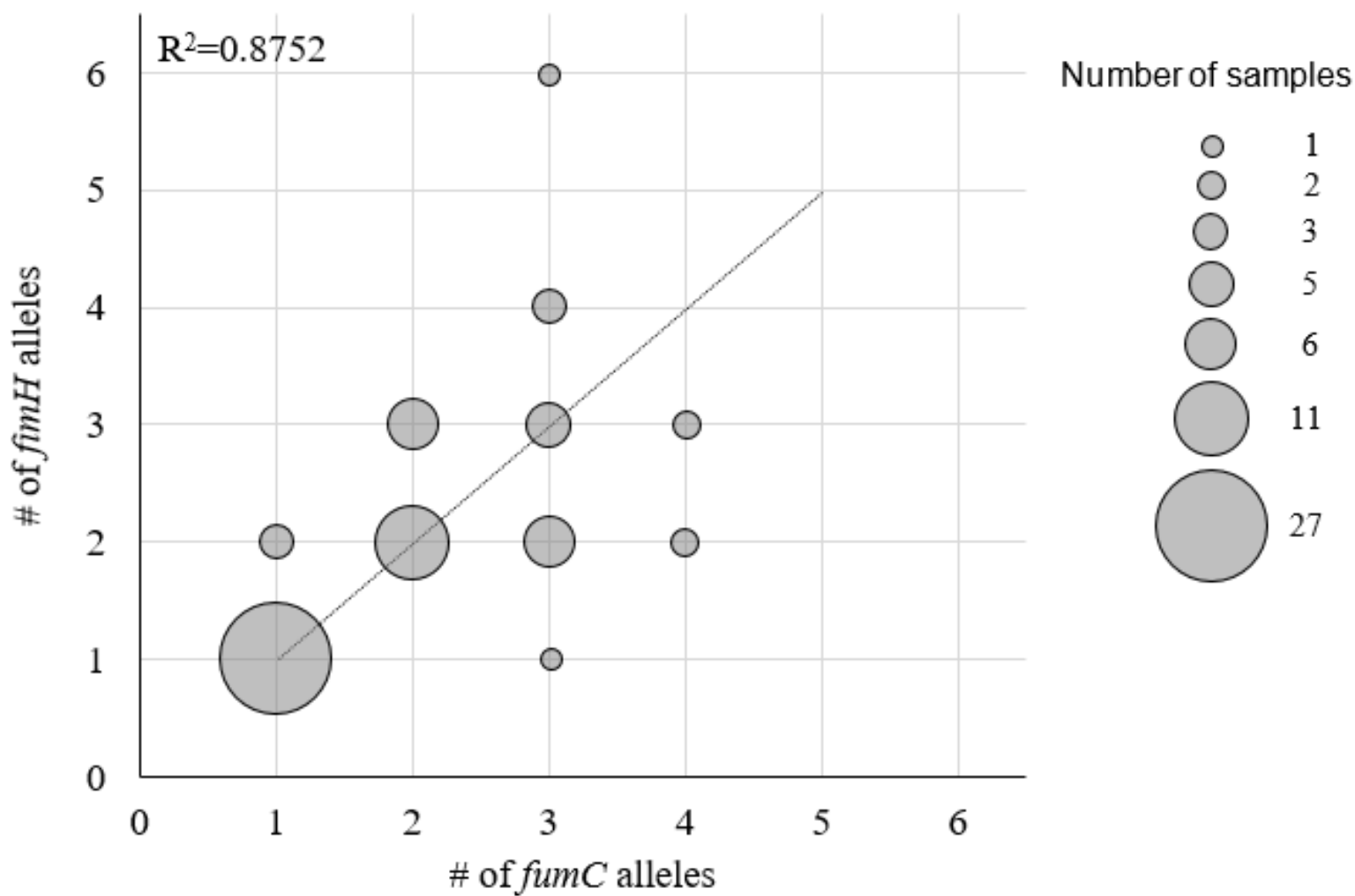677     surveying at least 40 single colonies per sample.

678     **Supplemental Figure 6.** Putative rare novel *fumC* alleles identified by lowering the error

679     threshold from 0.8% to 0.5%, marked in open shapes. Known alleles from the same sample as

680     the rare novel allele are marked in filled-in shapes of the same type and color. FumC-40 was

681     present in 3 different samples and therefore is marked by 3 different shapes.
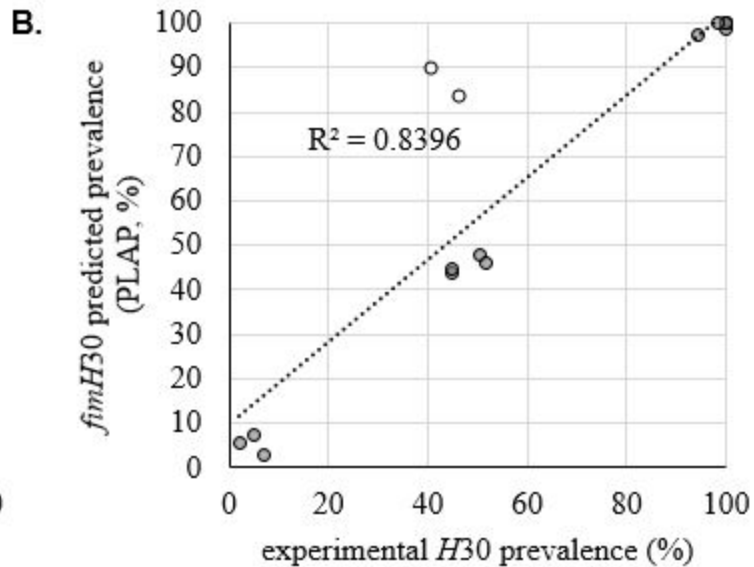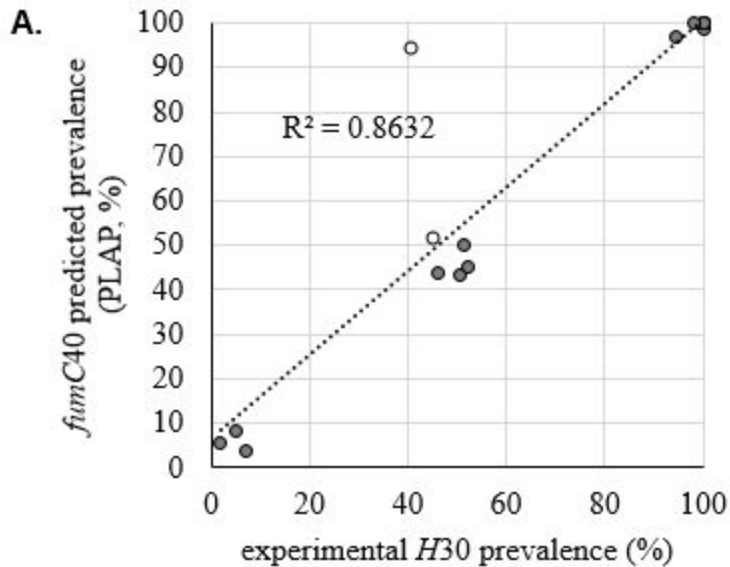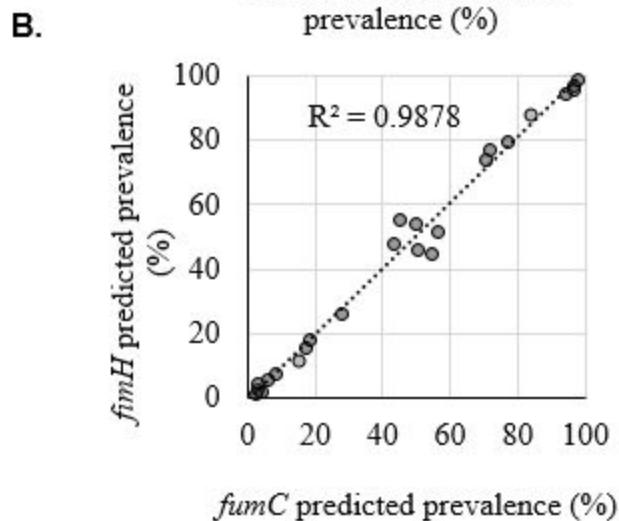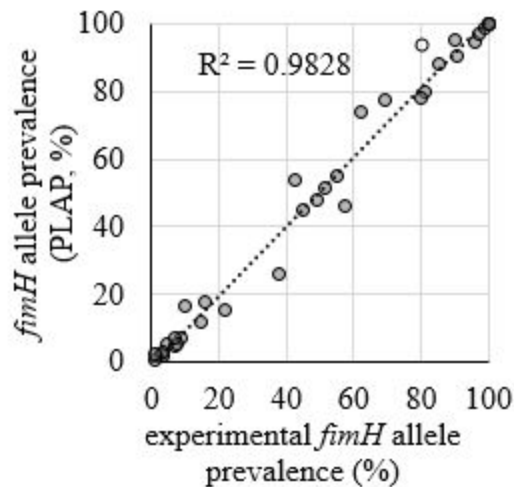
682     **Supplemental Figure 7.** Putative rare novel *fimH* alleles identified by lowering the error

683     threshold from 0.8% to 0.5%, marked in open shapes. Known alleles from the same sample as

684     the rare novel allele are marked in filled-in shapes of the same type and color. FimH-30 was

685     present in 3 different samples and therefore is marked by 3 different shapes.
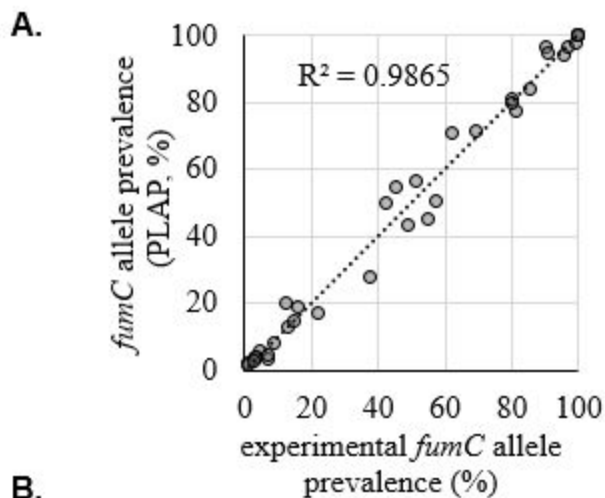
686     **Supplemental Figure 8.** Sampling of volunteer sample sets. Length of segments is proportional

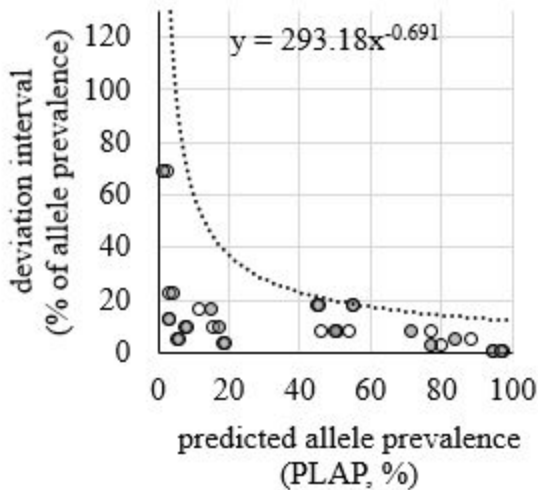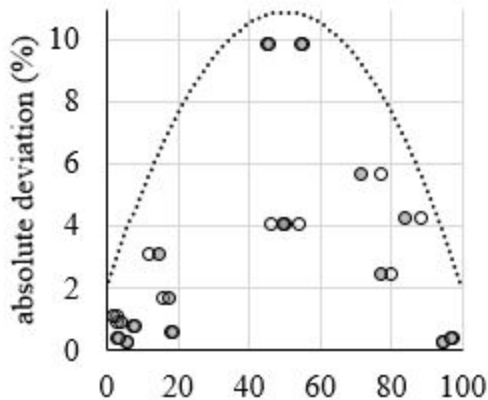687     to number of days between samples.

688     **Supplemental Figure 9. PLAP algorithm workflow.** Algorithms previously developed by

689     other groups include Trim-Galore, KMA, Minimap2. Not pictured but used during windowed

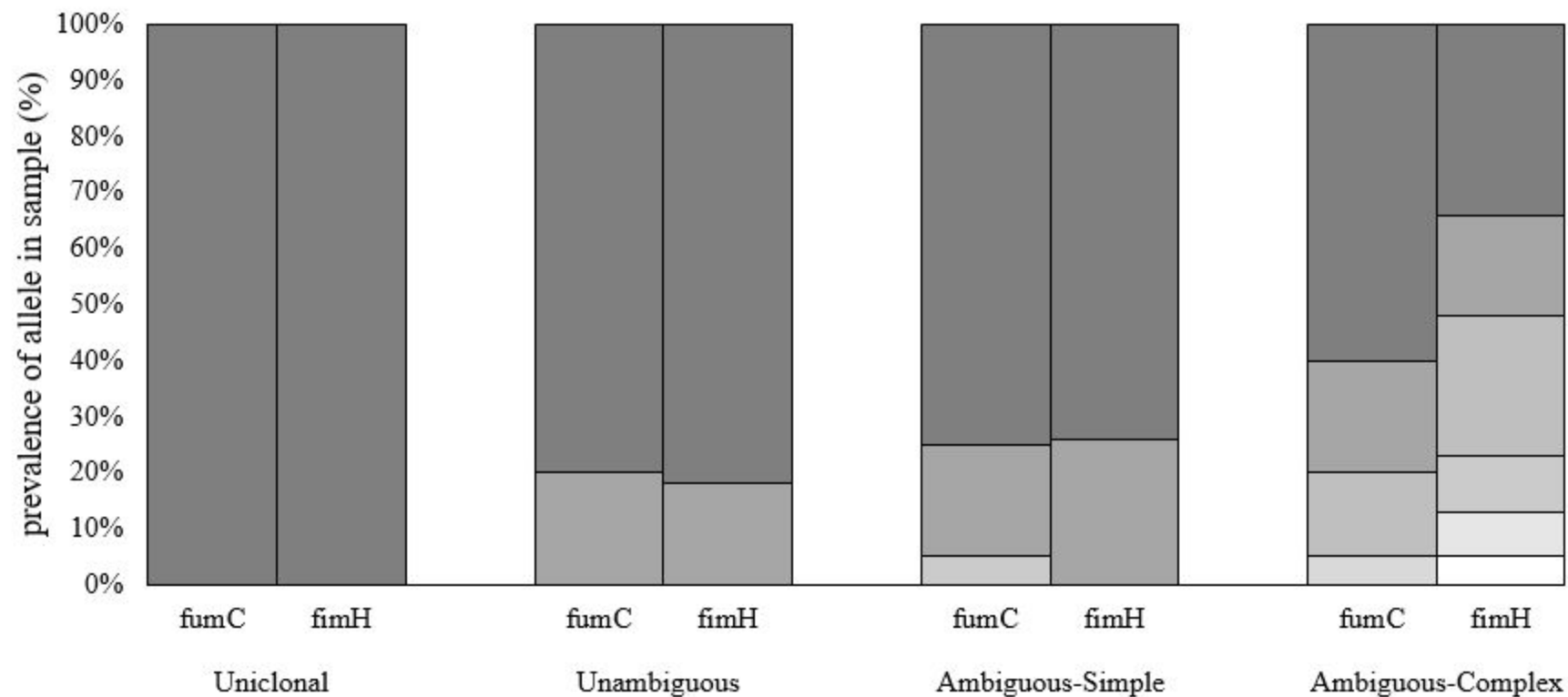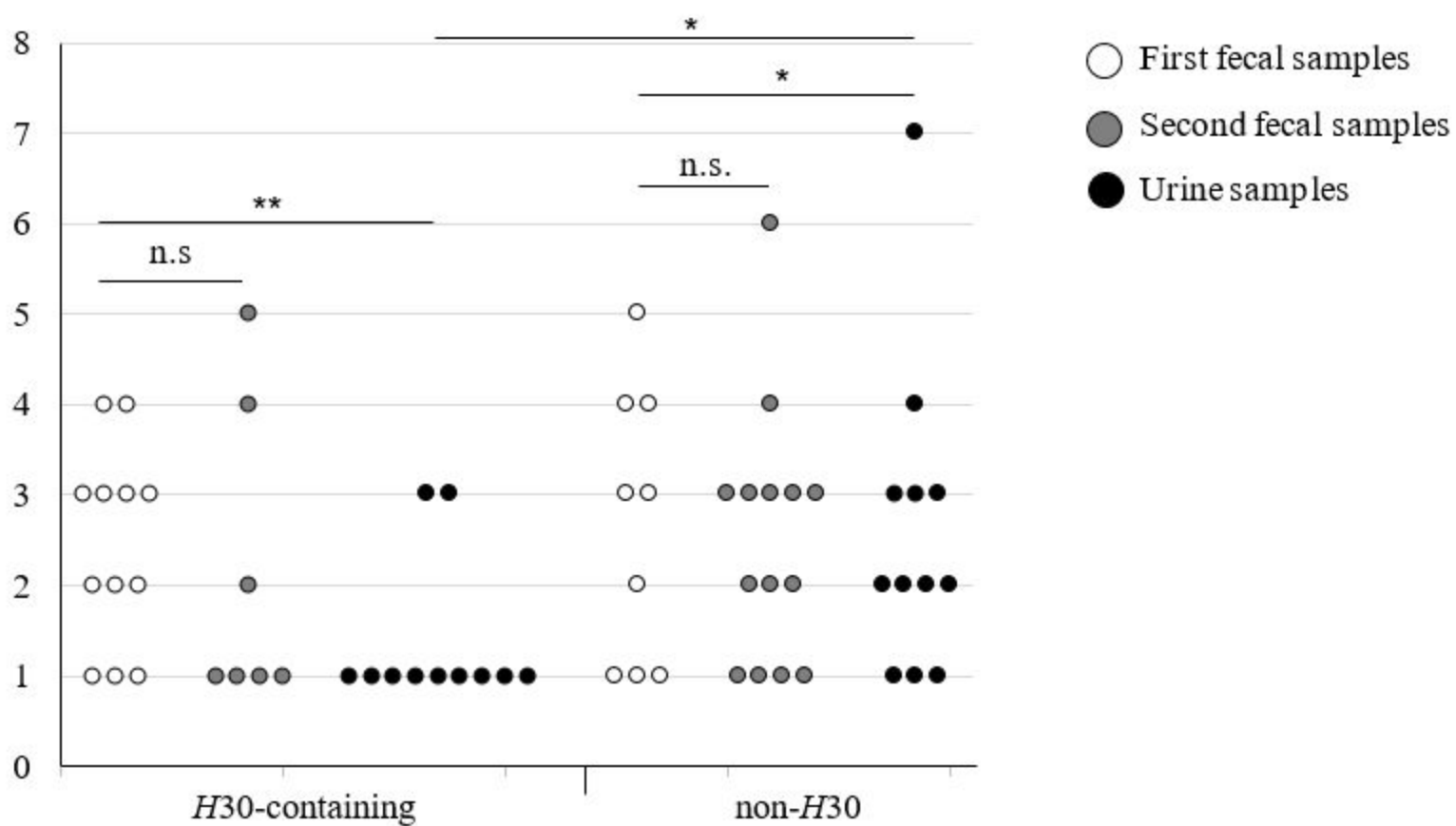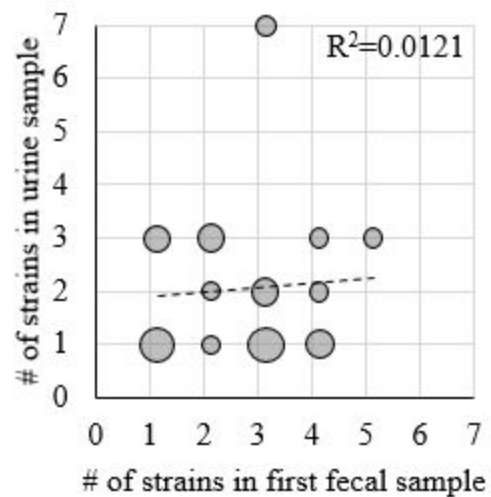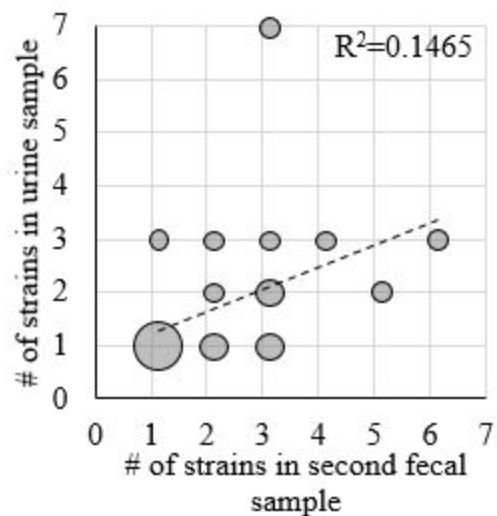690     coverage checks is SAMtools.

691

**A.** *fumC40* predicted prevalence (PLAP, %) vs experimental *H30* prevalence (%), R² = 0.8632

**B.** *fimH30* predicted prevalence (PLAP, %) vs experimental *H30* prevalence (%), R² = 0.8396

**A.**



**B.**

$$y = 293.18x^{-0.691}$$

| | Uniclonal | Unambiguous | Ambiguous-Simple | Ambiguous-Complex |
|---|---|---|---|---|
| Fecal samples | 14 | 10 | 9 | 10 |
| Urine samples | 13 | 2 | 6 | 3 |
| Total | **27** | **12** | **15** | **13** |

**A.**

# of strains in urine sample (y-axis), # of strains in first fecal sample (x-axis), $R^2=0.0121$

**B.**

# of strains in urine sample (y-axis), # of strains in second fecal sample (x-axis), $R^2=0.1465$

**C.**

# of strains in second fecal sample (y-axis), # of strains in first fecal sample (x-axis), $R^2=0.0645$

Number of samples

1   2   4   7

**A.** prevalence in second fecal sample (%) vs prevalence in first fecal sample (%)

**B.** prevalence in second fecal sample (%) vs prevalence in first fecal sample (%)