

36 **Abstract**

37

38 **Background.** Rising antibiotic resistance increasingly compromises empiric treatment. Knowing
39 the antibiotic susceptibility of a pathogen's close genetic relative(s) may improve empiric
40 antibiotic selection.

41

42 **Methods.** Using genomic and phenotypic data from three separate clinically-derived databases
43 of *Escherichia coli* isolates, we evaluated multiple genomic methods and statistical models for
44 predicting antibiotic susceptibility, focusing on potentially rapidly available information such as
45 lineage or genetic distance from archived isolates. We applied these methods to derive and
46 validate prediction of antibiotic susceptibility to common antibiotics.

47

48 **Results.** We evaluated 968 separate episodes of suspected and confirmed infection with
49 *Escherichia coli* from three geographically and temporally separated databases in Ontario,
50 Canada, from 2010-2018. Across all approaches, model performance (AUC) ranges for
51 predicting antibiotic susceptibility were greatest for ciprofloxacin (0.76-0.97), and lowest for
52 trimethoprim-sulfamethoxazole (0.51-0.80). When a model predicted a susceptible isolate, the
53 resulting (post-test) probabilities of susceptibility were sufficient to warrant empiric therapy for
54 most antibiotics (mean 92%). An approach combining multiple models could permit the use of
55 narrower-spectrum oral agents in 2 out of every 3 patients while maintaining high treatment
56 adequacy (~90%).

57

58 **Conclusions.** Methods based on genetic relatedness to archived samples in *E. coli* could be
59 used to predict antibiotic resistance and improve antibiotic selection.

60

61 **Keywords:** Empiric antibiotics; antibiotics; genomics; Gram-negative; antibiotic resistant
62 organisms; rapid diagnostics.

63 **Background**

64
65 Antibiotic resistance is a major global threat to public health (1). Antibiotic resistant organisms
66 (AROs) and mechanisms of antibiotic resistance are selected through the use of antibiotics in
67 humans, animals, and environments (2). The prevalence of AROs in human infections has been
68 increasing in many regions and across many different bacterial species (1,3). As a result,
69 empiric antibiotic therapy (administered prior to knowledge of the organism's antibiotic
70 susceptibility phenotype) has become increasingly challenging for both community- and
71 hospital-acquired infectious syndromes. Inadequate empiric therapy, i.e., treatment that does
72 not include an agent to which the etiologic pathogen is susceptible, has been associated with
73 worse patient outcomes (4-6). Moreover, increasing antibiotic resistance in common infections
74 leads to more frequent use of broader spectrum antibiotic agents, with added toxicity and
75 enhanced selection of antibiotic resistance by targeting a wider array of pathogenic,
76 opportunistic, and commensal bacteria.

77
78 Reducing the time from presentation and sample collection to reporting of antibiotic
79 susceptibility has long been touted as a potential means to improve early adequate therapy and
80 reduce the use of unnecessarily broad antibiotic agents (7,8). Development of rapid diagnostic
81 tests that can narrow these 'windows' of empiric antibiotic therapy are the focus of active
82 research, but translation of these tests to clinical practice has been slow, due to inconsistency
83 between individual or combined genetic loci and expected phenotype, specialized equipment,
84 cost, challenges of commercialization, and poor integration into clinical workflow (7). Recently,
85 genomic approaches have been identified as rapid diagnostic tests, offering the promise of
86 culture-independent (and -dependent) identification of: (1) species; (2) relationship(s) to genetic
87 neighbors/groups/clusters from databases of known isolates; and (3) prediction of antibiotic
88 resistance based on (2). However, the traditional approach to predicting antibiotic resistance

89 rests on identification of individual resistance loci to predict phenotype. This requires a high-
90 quality database of resistance-causative elements and is further complicated by significant cost,
91 large physical space requirements, complicated workflow, limited expertise, and long
92 sequencing/bioinformatic processing times even with real-time sequencing technologies (9). On
93 the other hand, a recently introduced alternative approach called genomic neighbor typing infers
94 antibiotic resistance and susceptibility by identifying sample's closest relatives in a database of
95 genomes with known phenotypes (10). This relies on strong correlation between phylogenetic
96 group and resistance phenotype, which is observed for many bacteria (11-13).

97

98 As neighbor typing uses all genomic data available from a given set of reads, identification of a
99 best match isolate or lineage (i.e. a genetically related cluster or group such as multi-locus
100 sequence type) can occur within minutes (as opposed to hours or days for loci-based
101 approaches depending upon the sequencing technology used). Proof of principle has been
102 demonstrated for *Streptococcus pneumoniae* and *Neisseria gonorrhoeae*, with determination of
103 resistance or susceptibility within ten minutes of Oxford Nanopore Technologies' © MinION
104 sequencing of cultured isolates and respiratory metagenomic samples (10). Limited data also
105 suggest that the association between antibiotic susceptibility phenotype and lineage may also
106 hold true for *Enterobacteriaceae* (14,15), but this approach requires further validation. It is also
107 unknown whether predicting antibiotic susceptibility phenotype based on the phenotype of the
108 nearest genetic neighbor provides advantages over using the average phenotype of a broader
109 (or higher level) lineage (e.g. ST, clonal complex, or cluster). In order to understand the
110 potential clinical application of these techniques, we sought to validate the association between
111 genetic relatedness (using nearest neighbor and lineage markers) and antibiotic susceptibility
112 phenotype in the most common Gram-negative pathogen in humans, *Escherichia coli*.

113

114 **Methods**

115

116 *Study Design*

117

118 We performed a retrospective study to evaluate whether genetic relatedness can predict
119 antibiotic susceptibility in *E. coli* isolates from 968 episodes of suspected and confirmed human
120 infection. Three separate datasets were combined for this analysis and include: (Dataset 1) 411
121 *E. coli* isolates from bloodstream infections at Sunnybrook Health Sciences Centre (SHSC), a
122 single tertiary care medical centre in Toronto, Canada, collected over the years 2010-2015;
123 (Dataset 2) 177 *E. coli* from suspected urinary tract infections from SHSC for the year 2018; and
124 (Dataset 3) 380 multi-drug resistant (MDR) *E. coli* isolates from urinary sources from the
125 Canadian province of Ontario collected in 2010 and 2015, where MDR was defined as
126 resistance to at least three different classes of routinely tested antibiotics.

127

128 *Resistance Phenotype*

129

130 Antibiotic susceptibility phenotypes for ciprofloxacin (fluoroquinolones), trimethoprim-
131 sulfamethoxazole (sulfonamides), ceftriaxone (3rd generation cephalosporins), gentamicin
132 (aminoglycosides), and ertapenem (carbapenems), were determined for each isolate using
133 Vitek 2 AST cards. Clinical laboratory standards institute (CLSI) 2015 breakpoints were
134 employed for determining susceptible and non-susceptible phenotypes for all datasets
135 (Supplementary Table 1). For Dataset 2, only formal ESBL testing was available (not
136 ceftriaxone MICs reported) and as such we classified all non-ESBL producing *E. coli* as
137 susceptible to ceftriaxone. Given that we calculated susceptibility using MIC's and static
138 breakpoints for all datasets, there were no temporal changes in the interpretation of
139 susceptibility. We have considered all non-susceptible isolates as resistant throughout this
140 paper.

141

142 *Whole Genome Sequencing*

143

144 Genomes for each dataset were sequenced separately using a NextSeq High Output platform
145 with Nextera Library Preparation with mean coverages of 134X, 90X, and 81X for dataset 1, 2,
146 and 3 respectively. Further sequencing details can be found in the Supplementary Methods.

147

148 *Overall Prediction Approach*

149 As overfitting could be a major limitation of our approach, and to simulate potential clinical
150 implementation strategies, we externally validated previously collected derivation datasets
151 (Dataset 1 and 3) for predicting the susceptibility of isolates from the most recent dataset
152 (Dataset 2). Where applicable, we also evaluated the sets internally with bootstrapping to adjust
153 for optimism (overfitting). We followed the general principles of the TRIPOD statement for
154 reporting of multivariable prediction models (16).

155

156 Four prediction model approaches were employed and are described in detail in the
157 Supplemental material. Briefly the first model is an *ST Parametric Model Approach* which uses
158 ST (marker of lineage) as a categorical predictor within a logistic regression model to predict the
159 probability of susceptibility. The second model is a *Cluster Parametric Model Approach* that
160 uses labelled genetic clusters (marker of lineage) as categorical predictors in a logistic
161 regression model to predict the probability of susceptibility. The third model is an *ST Reference*
162 *Database Approach* that uses the average prevalence of susceptibility of an antibiotic for a
163 given ST (marker of lineage) as the predicted probability of susceptibility. And the fourth model
164 is a *Best Genetic Match Reference Database Approach* that uses the susceptibility phenotype
165 of the best genetic match (nearest neighbor) in a reference database as the predicted
166 susceptibility. The above analyses were performed separately for each antibiotic.

167

168 To further simulate the antibiotic decision-making process, we explored two different scenarios
169 where an antibiotic was selected based upon sequential model outputs, either favouring narrow
170 spectrum agents or favoring high likelihood of adequacy of coverage. We called these
171 *Sequential Decision-making Models*. Further details on these methods are described in the
172 Supplemental material. Institutional research ethics board approval was obtained for this study.
173

174 **Results**

175

176 *Description of the Datasets*

177

178 We collected and sequenced the genomes of 968 unique *E. coli* isolates from separate clinical
179 episodes of suspected infection, across three datasets. The collection and sequencing details
180 are outlined in the methods and supplementary materials sections. The characteristics of each
181 dataset are shown in Table 1 below.

182

183 *Genetic Relatedness of Different Datasets and Antibiotic Resistance Phenotype*

184

185 To illustrate the linkage between relatedness and resistance, a genetic tree was constructed
186 using Mash distances, associated ST, antibiotic susceptibility phenotype, and source dataset
187 (Figure 1). Broad genetic clusters tended to match or be nested within STs, and closely related
188 isolates have similar antibiograms.

189

190 From the Mash tree illustrating high level relationships between the genomes (Figure 1), there
191 are clear genetic clusters that emerge. These phylogenetic groups are generally nested within
192 specific STs, with ST131 being the most prevalent in each dataset (Table 1). Unsurprisingly,
193 there is a higher prevalence of resistance for almost all antibiotic groups in the MDR dataset
194 (Dataset 3) compared to Datasets 1 and 2. Despite the fact that the three datasets are

195 separated on different scales (temporally, geographically, anatomically), they show genetic
196 diversity that is well distributed across different phylogroups (Figure 1).

197

198 *ST Parametric Model Approach (Lineage-based)*

199 We calculated the AUCs and test characteristics for predicting antibiotic susceptibility of Dataset
200 2 isolates using a parametric logistic regression model with STs as categorical predictors,
201 across a variety of derivation datasets (Supplementary Table 3). AUCs ranged from 0.89-0.91
202 for ciprofloxacin, 0.77-0.80 for ceftriaxone, 0.68 to 0.75 for gentamicin, and 0.6-0.73 for
203 trimethoprim-sulfamethoxazole. For ertapenem, we performed internal derivation for Dataset 1,
204 Dataset 3, and Datasets 1 and 3 combined (Dataset 2 could not be used as there was 100%
205 susceptibility to ertapenem), and found apparent and optimism adjusted AUCs ranging from 0.7-
206 0.99 and 0.67-0.99 respectively (Supplementary Table 2).

207

208 *Cluster Parametric Model Approach (Lineage-based)*

209

210 We calculated the AUC and test characteristics for predicting antibiotic susceptibility (internally)
211 for each dataset with a parametric logistic regression model with clusters as categorical
212 predictors, using a variety of derivation datasets (Supplementary Table 4). AUCs ranged from
213 0.76-0.9 for ciprofloxacin, 0.69-0.82 for ceftriaxone, 0.66-0.77 for gentamicin, and 0.65-0.75 for
214 trimethoprim-sulfamethoxazole. For ertapenem, we performed internal derivation for Dataset 1
215 and Dataset 3, and found apparent and optimism adjusted AUCs ranging from 0.74-0.98 and
216 0.7-0.98 respectively (Supplementary Table 2).

217

218 *ST Reference Database Approach (Lineage-based)*

219

220 We calculated the AUCs and test characteristics for predicting antibiotic susceptibility for
221 isolates in Dataset 2 with an *ST Reference Database*, using a variety of derivation datasets

222 (Supplementary Table 5). AUCs ranged from 0.85-0.95 for ciprofloxacin, 0.67-0.85 for
223 ceftriaxone, 0.73-0.83 for gentamicin, and 0.56-0.8 for trimethoprim-sulfamethoxazole.

224

225 *Best Genetic Match Reference Database Approach (Nearest neighbor)*

226

227 We calculated the AUCs and test characteristics for predicting antibiotic susceptibility for
228 isolates in Dataset 2 with a *Best Genetic Match Reference Database*, using a variety of
229 derivation datasets (Supplementary Table 6). For all isolates, AUCs ranged from 0.83-0.92 for
230 ciprofloxacin, 0.58-0.72 for ceftriaxone, 0.65-0.66 for gentamicin, and 0.53-0.63 for
231 trimethoprim-sulfamethoxazole. For the top 25%'tile of best match datasets, AUCs ranged from
232 0.85-0.97 for ciprofloxacin, 0.57-0.88 for ceftriaxone, 0.76-0.86 for gentamicin, and 0.54-0.73 for
233 trimethoprim-sulfamethoxazole. In any method based on comparing new samples with an
234 existing database, it is important to investigate how large the database needs to be to permit
235 accurate prediction. Thus, we evaluated the impact of varying reference database sizes on the
236 performance of the genetic distance approach, with results shown in Figure S1.

237

238 *Summary of Test Characteristics Across Models*

239 In Figures 2a-d, we summarize the post-test probabilities of susceptibility for: (1) positive model
240 predictions indicating a likely susceptible isolate (PPV [positive predictive value]); and (2) for
241 negative model predictions indicating a likely resistant isolate (1-NPV [Negative predictive
242 value]). Here we see that the models that provide the best positive and negative predictive
243 values are the *ST Reference Database*, the *Best Genetic Match Reference Database*, and
244 combinations of the two (Supplementary Table 7). We also see that for most antibiotics and
245 models, that the post-test probabilities for results indicating a susceptible result are sufficiently
246 high (relative to syndromic thresholds) to support recommendation of therapy. Similarly, for
247 most antibiotics and models the post-test probabilities of results indicating resistance is

248 sufficiently low to support withholding therapy for a particular agent. Lastly, combined derivation
249 datasets tend to have the most consistent model performance.

250

251 *Sequential Antibiotic Decision-Making Based Upon a Prespecified Antibiotic Cascade*

252 When using a sequential selection model favouring narrow spectrum antibiotics, adequate
253 therapy was achieved for 89-98% of recommendations, with 47-70% of recommendations being
254 a narrower spectrum agent with oral formulation (i.e. ciprofloxacin or trimethoprim-
255 sulfamethoxazole). Spectrum scores for the narrow spectrum model cascade were consistently
256 lower than those attained with the employment of typical empiric agents (Supplementary Table
257 7). Using an adequacy focused cascade, adequate therapy was achieved for 98-100% of
258 recommendations, with 0-5% of recommendations being narrower spectrum agents with oral
259 formulation options. Spectrum scores were generally high, but were lower than the most broad
260 spectrum agent ertapenem.

261

262 To summarize, when using a narrow-spectrum focused cascade, the models yield excellent
263 adequacy while enabling over 2/3 of recommendations to be narrow spectrum oral agents. When
264 using an adequacy focused cascade they yield close to perfect adequacy but at the expense of
265 using broader spectrum agents (Supplementary Table 7).

266

267

268 **Discussion**

269

270 In this study, we demonstrate that we can predict the resistance phenotype of *E. coli* by rapidly
271 determining the genetic relatedness of the infecting pathogen to a database of sequenced
272 isolates with known resistance phenotypes. We show that these relationships can be used to
273 generate post-test probabilities of susceptibility in excess of 0.8 or 0.9 which render antibiotics
274 with a high prevalence of resistance to be empirically useful (e.g. ciprofloxacin or trimethoprim-

275 sulfamethoxazole for urinary tract infection). In essence, adoption of this approach would modify
276 the current stages of empiric therapy to introduce a new window that is informed by genetic
277 relatedness, and information that is available in advance of standard phenotypic testing. This
278 approach is a supplement to, not a replacement for, gold standard phenotype. The contribution
279 of such a system could be to improve the quality of antimicrobial prescribing in the window
280 between culture positivity and phenotypic testing results, by 1) reducing the expected time from
281 culture positivity to adequate treatment and 2) reducing the duration of use of broad-spectrum
282 agents to treat infections that could be adequately treated with a narrow-spectrum agent.

283

284 Looking at the distribution of STs across the datasets (Figure 1, Table 1) we see that Dataset 3
285 is enriched for ST131, and this is presumably due to the intentionally biased sampling approach
286 toward MDR. However, it means that lineage and nearest neighbor approaches will be able to
287 draw sample phenotype predictions from different datasets based on genetic proximity, not
288 simply ones most closely related on temporal, geographic, or anatomic scales.

289

290 When predicting antibiotic resistance based on lineage, *ST Parametric Model* and *Cluster*
291 *Parametric Model Approaches* (Supplementary Tables 3 and 4) provide reasonable
292 discrimination for ciprofloxacin susceptibility (AUCs 0.76 - 0.91), in keeping with existing
293 literature (15). This emphasizes the strong association between lineage and ciprofloxacin
294 susceptibility. By contrast there were only modest associations for the other antibiotic classes
295 (AUCs 0.6-0.82). For all antibiotic classes, the use of combined derivation datasets (1 and 2 and
296 3 or 1 and 3) seemed to perform well most consistently, and this supports the use of an
297 aggregated derivation dataset across time, geography, and anatomic location.

298

299 Using an *ST Reference Database* approach has the appeal of providing improved predictions

300 for less common STs. When considering only those isolates with a matching ST in the reference
301 database, the discrimination of this approach paralleled and sometimes exceeded that of *ST*
302 and *Cluster Parametric Models* (Supplementary Table 5). The notable downside to this
303 approach is the inability to provide predictions for sequence types outside of the reference
304 database, though the proportion of these were small and decreased with increasing reference
305 database size.

306

307 *A Best Genetic Match Reference Database Approach* (nearest neighbor) offers potential
308 improvement over the *ST Reference Database Approach* (lineage), in that it might improve
309 predictive performance for those classes that have a weaker association with specific lineage
310 (e.g. ceftriaxone, gentamicin, trimethoprim-sulfamethoxazole). The genetic distance approach
311 seemed to operate best under two circumstances: (1) when only the 'top matches' were
312 considered and (2) when a combined derivation dataset was used. This 'top match' nearest
313 neighbor approach is potentially implemented using a predefined threshold of Mash distance (or
314 other genetic distance measure), but suffers from having a significant number of samples for
315 which predictions may not be offered. Interestingly, the 'top match' approach can improve AUCs
316 compared to other approaches for the antibiotics that are less strongly associated with
317 phenotype (ceftriaxone, gentamicin, trimethoprim-sulfamethoxazole).

318

319 One important consideration with a genetic distance based model, is the consideration of how
320 large the reference database should be. Using repeated sampling methods, we found that
321 optimal performance was achieved with comparatively small reference database sizes,
322 consisting of 100-200 samples. This was consistent across the different classes of antibiotics
323 tested, with plateaus in performance after ~200 samples (Supplemental Figure 1), and
324 consistent with previous work in other pathogens (10). However, the necessary size will likely
325 depend on the diversity of the population being evaluated, with more diverse populations

326 requiring larger databases. This is reflected in the performance seen with combined derivation
327 datasets.

328

329 There are limitations to our study. First, we were only able to consider the construction of
330 reference databases confined to the geographic region of the province of Ontario. However, this
331 region is geographically large and contains a population of over 14 million people. As such, our
332 results still support construction of regional databases at least at this jurisdictional level, which is
333 generalizable to many areas globally. Secondly, we did not evaluate the utility of this approach
334 for other bacterial species, however *E. coli* is the most common Gram-negative pathogen in the
335 hospital and community. Thirdly, we did not examine in detail the reasons for failure to
336 accurately predict a susceptibility phenotype. The potential reasons for imperfect performance
337 are numerous, and include comprehensiveness of reference database, the acquisition of mobile
338 genetic elements or new resistance mutations, human/labelling error, and imperfection in
339 phenotypic testing methodologies. We did not seek to explore all of these reasons, but instead
340 we sought to quantify the overall additional benefit these approaches could add. However,
341 future work specifically exploring and characterizing the modes of failure is warranted. Other
342 future work will aim to prospectively evaluate these techniques in a clinical setting using rapid
343 sequencing approaches across geographic scales and additional pathogens. Our recent work
344 suggests we can predict susceptibility within minutes, and when this is combined with rapid
345 DNA extraction kits (<30 mins) and rapid library preparation kits (<15 min), then it is currently
346 feasible to go from clinical collection to result in under 60 minutes (10). This time frame will likely
347 shrink as DNA extraction and library preparation steps are further improved and simplified. In
348 summary, our results suggest that rapidly obtainable genomic information from clinical isolates
349 can support intelligent choices that improve empiric antibiotic therapy, both by rescuing narrow
350 spectrum agents for therapeutic use and by better selecting broader spectrum agents.

351

352 **Funding**

353

354 We would like to acknowledge the funding support provided through the McLaughlin Accelerator

355 Grant at the University of Toronto. DRM was supported through the Clinician Scientist Training

356 Program at the University of Toronto, and a Canadian Institutes of Health Research Training

357 Award. The project described was also supported by Cooperative Agreement Number

358 U54GM088558 from the National Institute Of General Medical Sciences. The content is solely

359 the responsibility of the authors and does not necessarily represent the official views of the

360 National Institute Of General Medical Sciences or the National Institutes of Health.

361

362 **Conflicts of Interest**

363

364 ML has received honoraria/consulting income from Merck, Pfizer, Antigen Discovery, and

365 Affinivax, and research support through his institution from Pfizer.

366 **Acknowledgements**

367

368 There are no additional acknowledgements.

369

370

371 **References**

- 372 1. Review on Antimicrobial Resistance. 2016. Tackling Drug-resistant Infections Globally: Final
373 Report and Recommendations.
- 374 2. Lipsitch M, Samore MH. 2002. Antimicrobial Use and Antimicrobial Resistance: A Population
375 Perspective. *Emerging Infectious Diseases*.
- 376 3. Frieden T. 2013. Antibiotic Resistance Threats in the United States 2013.
- 377 4. Kumar A, Roberts D, Wood KE, Light B, Parrillo JE, Sharma S, Suppes R, Feinstein D,
378 Zanotti S, Taiberg L, Gurka D, Kumar A, Cheang M. 2006. Duration of hypotension before
379 initiation of effective antimicrobial therapy is the critical determinant of survival in human
380 septic shock. *Crit Care Med* 34:1589–1596.
- 381 5. Bodilsen J, Dalager-Pedersen M, Schønheyder HC, Nielsen H. 2016. Time to antibiotic
382 therapy and outcome in bacterial meningitis: a Danish population-based cohort study. *BMC*
383 *Infect Dis* 16:392.
- 384 6. Doganis D, Sifas K, Mavrikou M, Issaris G, Martirosova A, Perperidis G, Konstantopoulos A,
385 Sinaniotis K. 2007. Does early treatment of urinary tract infection prevent renal damage?
386 *Pediatrics* 120:e922–8.
- 387 7. MacFadden DR, Leis JA, Mubareka S, Daneman N. 2014. The opening and closing of
388 empiric windows: the impact of rapid microbiologic diagnostics. *Clin Infect Dis* 59:1199–1200.
- 389 8. Timbrook TT, Morton JB, McConeghy KW, Caffrey AR, Mylonakis E, LaPlante KL. 2017. The
390 Effect of Molecular Rapid Diagnostic Testing on Clinical Outcomes in Bloodstream Infections:
391 A Systematic Review and Meta-analysis. *Clin Infect Dis* 64:15–23.
- 392 9. Boolchandani M, D’Souza AW, Dantas G. 2019. Sequencing-based methods and resources
393 to study antimicrobial resistance. *Nat Rev Genet* 20:356–370.
- 394 10. Břinda K, Callendrello A, Ma KC, MacFadden DR, Charalampous T, Lee RS, Cowley L,
395 Wadsworth CB, Grad YH, Kucherov G, O’Grady J, Baym M, Hanage WP. 2020. Rapid

- 396 inference of antibiotic resistance and susceptibility by genomic neighbour typing. *Nature*
397 *Microbiology*.
- 398 11. Nazir H, Cao S, Hasan F, Hughes D. 2011. Can phylogenetic type predict resistance
399 development? *J Antimicrob Chemother* 66:778–787.
- 400 12. Massot M, Daubié A-S, Clermont O, Jaureguy F, Couffignal C, Dahbi G, Mora A, Blanco J,
401 Branger C, Mentré F, Eddi A, Picard B, Denamur E. 2016. Phylogenetic, virulence and
402 antibiotic resistance characteristics of commensal strain populations of *Escherichia coli* from
403 community subjects in the Paris area in 2010 and evolution over 30 years. *Microbiology*
404 162:642.
- 405 13. Hussain A, Ranjan A, Nandanwar N, Babbar A, Jadhav S, Ahmed N. 2014. Genotypic and
406 Phenotypic Profiles of *Escherichia coli* Isolates Belonging to Clinical Sequence Type 131
407 (ST131), Clinical Non-ST131, and Fecal Non-ST131 Lineages from India. *Antimicrob Agents*
408 *Chemother* 58:7240.
- 409 14. Tchesnokova V, Billig M, Chattopadhyay S, Linardopoulou E, Aprikian P, Roberts PL,
410 Skrivankova V, Johnston B, Gileva A, Igusheva I, Toland A, Riddell K, Rogers P, Qin X,
411 Butler-Wu S, Cookson BT, Fang FC, Kahl B, Price LB, Weissman SJ, Limaye A, Scholes D,
412 Johnson JR, Sokurenko EV. 2013. Predictive diagnostics for *Escherichia coli* infections
413 based on the clonal association of antimicrobial resistance and clinical outcome. *J Clin*
414 *Microbiol* 51:2991–2999.
- 415 15. MacFadden DR, Melano RG, Coburn B, Tijet N, Hanage WP, Daneman N. 2019. Comparing
416 Patient Risk Factors, Sequence Type, and Resistance Loci Identification Approaches for
417 Predicting Antibiotic Resistance in Bloodstream Infections. *J Clin Microbiol*.
- 418 16. Collins GS, Reitsma JB, Altman DG, Moons KGM. 2015. Transparent Reporting of a
419 multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD): the TRIPOD
420 Statement. *Br J Surg* 102:148–158.
- 421

422 **Tables**

423

424 **Table 1. Characteristics of datasets.**

	Dataset 1	Dataset 2	Dataset 3
Number of Isolates (n=968)	411	177	380
Collection Period	2010-2015	2018	2010 and 2015
Location	Toronto (City)	Toronto (City)	Southeastern Ontario
Location Type	Hospital Lab	Hospital Lab	Hospital Lab
Inpatient/Outpatient	Inpatient	In/Outpatient	In/Outpatient
Anatomic Site	Blood	Urine	Variable
Sampling Bias	None	None	Multi-drug Resistant (MDR)
Antibiotic Susceptibility*			
<i>Ciprofloxacin</i>	297 (72)	120 (68)	118 (31)
<i>Ceftriaxone</i>	357 (87)	155 (88)	236 (62)
<i>Gentamicin</i>	355 (86)	156 (88)	229 (60)
<i>Trimethoprim-Sulfamethoxazole</i>	292 (71)	125 (71)	79 (21)
<i>Ertapenem</i>	410 (99)	177 (100)	338 (89)
Predominant ST*			
1193	14 (3.4)	15 (8.5)	21 (5.5)
127	15 (3.6)	7 (4.0)	3 (0.8)
131	87 (21)	36 (20)	170 (45)
38	7 (1.7)	2 (1.1)	12 (3.2)
405	9 (2.2)	1 (0.6)	11 (2.9)
648	6 (1.5)	8 (4.5)	17 (4.5)
69	22 (5.4)	13 (7.3)	23 (6.1)
73	57 (14)	20 (11)	20 (5.3)
95	58 (14)	19 (11)	12 (3.2)
Other	136 (33)	56 (32)	91 (24)

425

426 *N(%)

427

428

429 **Figure Legends**

430 **Figure 1.** Mash tree (left), ST, phenotypic susceptibility by antibiotic, and dataset, by
431 individual isolate. Antibiotic susceptibility is denoted in green and resistance in red, with
432 antibiotics abbreviated as: CIP (ciprofloxacin), CRO (ceftriaxone), GEN (gentamicin), SXT
433 (trimethoprim-sulfamethoxazole), ETP (ertapenem). ST 9999 represents all remaining or
434 unknown STs.

435
436 **Figure 2a.** Selected post-test probabilities of ciprofloxacin susceptibility (in Dataset 2)
437 based on Model Predictions of Resistant or Susceptible, by model type and derivation
438 dataset.

439
440 **Figure 2b.** Selected post-test probabilities of ceftriaxone susceptibility (in Dataset 2)
441 based on Model Predictions of Resistant or Susceptible, by model type and derivation
442 dataset.

443
444 **Figure 2c.** Selected post-test probabilities of gentamicin susceptibility (in Dataset 2)
445 based on Model Predictions of Resistant or Susceptible, by model type and derivation
446 dataset.

447
448 **Figure 2d.** Selected post-test probabilities of trimethoprim-sulfamethoxazole
449 susceptibility (in Dataset 2) based on Model Predictions of Resistant or Susceptible, by
450 model type and derivation dataset.

451
452 **The aim of the models is to shift green data-points to the right (indicating that the model*
453 *has correctly classified certain isolates as susceptible), and red data points to the left*
454 *(indicating that the model has correctly classified certain isolates as resistant).*

455
456
457









