

1 **Predicting Variable Gene Content in *Escherichia coli* using Conserved**  
2 **Genes**

3  
4 Marcus Nguyen<sup>a,b</sup>, Zachary Elmore<sup>c</sup>, Clay Ihle<sup>c</sup>, Francesco S. Moen<sup>c</sup>, Adam D. Slater<sup>c</sup>,  
5 Benjamin N. Turner<sup>c</sup>, Bruce Parrello<sup>b,d</sup>, Aaron A. Best<sup>c</sup>, James J. Davis<sup>a,b#</sup>

6  
7 <sup>a</sup>Data Science and Learning Division, Argonne National Laboratory, Lemont, IL, 60439, USA

8 <sup>b</sup>Consortium for Advanced Science and Engineering, University of Chicago, Chicago, IL,  
9 60637, USA

10 <sup>c</sup>Biology Department, Hope College, Holland, MI, 49423, USA

11 <sup>d</sup>Fellowship for Interpretation of Genomes, Burr Ridge, IL, 60527, USA

12  
13  
14 #Corresponding Author

15 email: [jjdavis@anl.gov](mailto:jjdavis@anl.gov)

16 Phone: +1-630-252-1190

17 Fax: +1-630-252-6333

18  
19 **Running Title**

20 Predicting *E. coli* Gene Content

21  
22 **Abbreviations**

23 AMR: Antimicrobial Resistance

24 BV-BRC: Bacterial and Viral Bioinformatics Resource Center

25 MAG: Metagenome assembled genome

26 MLST: Multi Locus Sequence Type

27 PATRIC: PATHosystems Resource Integration Center

28 RAST: Rapid Annotation Subsystem Technology

29 XGBoost: Extreme Gradient Boosting

30  
31 Abstract: 250 words

32 Text: 5,238 words

33  
34 **Keywords**

35 machine learning, horizontal gene transfer, antimicrobial resistance, bacterial virulence,  
36 phylogeny

## 39 Abstract

40 Having the ability to predict the protein-encoding gene content of an incomplete genome or  
41 metagenome assembled genome is important for a variety of bioinformatic tasks. In this  
42 study, as a proof of concept, we built machine learning classifiers for predicting variable  
43 gene content in *Escherichia coli* genomes using only the nucleotide k-mers from a set of 100  
44 conserved genes as features. Protein families were used to define orthologs, and a single  
45 classifier was built for predicting the presence or absence of each protein family occurring  
46 in 10-90% of all *E. coli* genomes. The resulting set of 3,259 XGBoost classifiers had a per-  
47 genome average macro F1 score of 0.944 [0.943-0.945, 95% CI]. We show that the F1  
48 scores are stable across MLSTs, and that the trend can be recapitulated by sampling a  
49 smaller number of core genes or diverse input genomes. Surprisingly, the presence or  
50 absence of poorly annotated proteins, including “hypothetical proteins” were accurately  
51 predicted (F1 = 0.902 [0.898-0.906, 95% CI]). Models for proteins with horizontal gene  
52 transfer-related functions had slightly lower F1 scores, but were still accurate (F1s = 0.895,  
53 0.872, 0.824, and 0.841 for transposon, phage, plasmid, and antimicrobial resistance-  
54 related functions, respectively). Finally, using a holdout set of 419 diverse *E. coli* genomes  
55 that were isolated from freshwater environmental sources, we observed an average per-  
56 genome F1 score of 0.880 [0.876-0.883, 95% CI], demonstrating the extensibility of the  
57 models. Overall, this study provides a framework for predicting variable gene content  
58 using a limited amount of input sequence data.

59

60

## 61 **Importance**

62 Having the ability to predict the protein-encoding gene content of a genome is important  
63 for assessing genome quality, binning genomes from shotgun metagenomic assemblies, and  
64 assessing risk due to the presence of antimicrobial resistance (AMR) and other virulence  
65 genes. In this study, we built a set of binary classifiers for predicting the presence or  
66 absence of variable genes occurring in 10-90% of all publicly available *E. coli* genomes.  
67 Overall, the results show that a large portion of the *E. coli* variable gene content can be  
68 predicted with high accuracy, including genes with functions relating to horizontal gene  
69 transfer. This study offers a strategy for predicting gene content using limited input  
70 sequence data.

71

## 72 Introduction

73 In genomic and metagenomic sequencing studies it is common to encounter  
74 incomplete genomes. Having the ability to predict the protein-encoding gene content for  
75 an incomplete genome or metagenome assembled genome (MAG) based on limited data is  
76 important for a multitude of bioinformatic tasks. These include estimating genome quality  
77 and completeness, assessing the metabolic capabilities of the organism or community, and  
78 understanding the potential for a genome to encode antimicrobial resistance (AMR) and  
79 other virulence factors. Over the years, a variety of approaches for predicting protein-  
80 encoding gene content and their associated functions from limited source data have been  
81 devised, and they vary based on the study design and the experimental needs. Some  
82 approaches are as straight forward as searching for a known set of conserved genes, while  
83 others use more complex algorithmic and artificial intelligence-based approaches to  
84 predict the presence or absence of key genes.

85  
86 Perhaps the largest body of work predicting gene content and protein functions using  
87 limited input data has come from the field of metagenomics. Metagenomic studies  
88 routinely perform amplicon sequencing of the 16S rRNA gene in order to determine the  
89 microbial diversity of an environment. However, since 16S sequencing does not provide  
90 information about the gene content of the sample, many methods have been developed to  
91 infer gene content and functions from the 16S sequence. These include popular tools such  
92 as, PICRUSt (1), Piphillin (2), Tax4Fun (3), PAPRICA (4), PanFP (5), and MicFunPred (6).  
93 Although the algorithmic steps vary, in essence, these tools utilize the relationship between  
94 phylogeny and gene content and make their predictions based on a set of closely related  
95 reference genomes (7). These methods are useful for inferring the metabolic capabilities of  
96 the constituents of a sample, especially when deep shotgun sequencing is unavailable.  
97 However, they come with understandable limitations due to the limited information  
98 encoded in 16S amplicon sequences (particularly single variable regions), the size and  
99 scope of the reference databases, and the ability of horizontal gene transfer to disrupt gene  
100 content between close relatives.

101  
102 Another related bioinformatic application, prediction of genome completeness—i.e.,  
103 inferring that all genes are present in a given genome—is crucial for assessing genome  
104 quality and the reliability of downstream comparative analyses (8). This is important  
105 when working with genomes that are assembled from short reads, genomes from single  
106 cell sequencing data, and MAGs. Many studies have established measures to predict  
107 genome completeness by searching for sets of universal or lineage-specific genes from  
108 close relatives including CheckM (9), BUSCO (10), CEGMA (11), mOTU (12), Anvi'o (13),  
109 and CheckV (14). Parrello and colleagues recently extended this concept by building a tool  
110 that predicts genome completeness for bacterial and archaeal genomes from the protein  
111 annotations in the PATRIC database (15, 16). Using a set of approximately 2,000 well  
112 annotated “roles,” which are the individual atomic functions of a protein in the SEED  
113 annotation schema (17), they built a set of machine learning classifiers that predicted the  
114 presence or absence of each role based on the presence or absence of the other roles in the  
115 set. This enabled them to both quantify the completeness of the genome and provide an  
116 estimate for the expected number of occurrences of each role per genome. In most cases,

117 when genome completeness scores deviate from expectation, the genome is reliably  
118 incomplete or contaminated with sequences from another organism. Similarly, recent  
119 updates to the CheckM algorithm, CheckM2, incorporate the use of machine learning  
120 models, which include the KEGG protein annotations as part of the feature vector (18, 19).  
121 Another recent tool called MetaPredict uses a set of classifiers to predict the presence or  
122 absence of KEGG modules in a MAG, based on the presence or absence of the existing  
123 annotations in the MAG (20). Although all of these methods have proven to be useful for  
124 predicting genome completeness, a potential downside is that they are designed to predict  
125 the presence or absence of well characterized genes, which may not fully capture patterns  
126 in the variable strain-specific gene content across a species.

127  
128 AMR genes and other virulence factors are often among the set of strain specific genes that  
129 vary between the members of a species. Many of these genes are found on mobile genetic  
130 elements, so their occurrences sometimes do not match the phylogeny of a given taxon.  
131 Many bioinformatic tools have been developed to search for AMR and virulence genes  
132 within a genome or metagenome using both sequence similarity (21-28) and machine  
133 learning techniques (29, 30). Since shotgun metagenomic studies sample multiple  
134 genomes, and their assemblies are often incomplete, methods that identify AMR and  
135 virulence genes are not always able to identify the source genome for a given AMR gene.  
136 To this end, some studies have attempted to predict the source genomes for the AMR genes  
137 in a sample using either statistical (31) or machine learning methods (32, 33).

138  
139 Many studies have also been designed to predict AMR phenotypes from genome sequences  
140 by training machine learning models using the genomes and laboratory-derived  
141 antimicrobial susceptibility test data (34). Importantly, several of these studies have  
142 demonstrated that AMR phenotypes can be predicted using the phylogeny of the strains,  
143 either by learning the tree structure, mapping phenotypes from close relatives, or building  
144 machine learning models from conserved parts of the genome (35-38). This has been  
145 demonstrated even in cases where the AMR phenotype is the result of a horizontal gene  
146 acquisition and is presumably due to the machine learning models learning non-linear  
147 relationships in the input data. Although there is a clear link between phenotype and  
148 genotype, to date we still lack well-developed tools for predicting whether an AMR or  
149 virulence gene should or should not be present in a genome given a set of existing  
150 sequences from a contig or MAG.

151  
152 *E. coli* is the most widely studied bacterial species, and there are currently well over 30,000  
153 sequenced *E. coli* genomes in the public domain. All of the genes of the species can be  
154 thought of as a pan-genome consisting of a conserved set of core genes that are held in  
155 common among all members of the species, plus tens of thousands of accessory genes that  
156 are often strain-specific and vary in their frequency of occurrence (39-43). These variable  
157 genes encode a variety of known and unknown functions including AMR and virulence  
158 genes. Since pre-existing tools are mainly designed to predict the presence or absence  
159 genes with well annotated pprotein functions, in this study, as a proof of concept, we  
160 wanted to see the extent to which it is possible to predict the presence or absence of  
161 variable genes in *E. coli* using only the nucleotide sequences of a set of universal genes to  
162 make the predictions.

163  
164  
165

## Materials and Methods

### 166 Genomes and datasets

167 A high-quality, diverse set of publicly available *Escherichia coli* genomes was selected for  
168 building the models. All *E. coli* genomes were downloaded from the Bacterial and Viral  
169 Bioinformatics Resource Center (BV-BRC) FTP site (<ftp.bvbrc.org>) on March 28, 2022  
170 (**Figure 1**). The BV-BRC is a large bioinformatics resource center that maintains the  
171 PATHosystems Resource Integration Center (PATRIC)(15, 44). Each bacterial genome in  
172 the BV-BRC has been uniformly annotated using the Rapid Annotation Using Subsystem  
173 Technology (RAST) pipeline (45), and the analysis includes computation of genome quality  
174 (16), protein encoding gene annotations (45), protein family assignments, etc. (46). All *E.*  
175 *coli* genomes lacking “Complete” or “WGS” designations (sourced from GenBank) (47), and  
176 those that were listed as being poor quality (16), were excluded from consideration. Any  
177 genome with less than half of the average number of genes per genome was also excluded.  
178 The genome set was further filtered to ensure that all of the core genes that were used for  
179 generating features for the models (described below) were present, and that each gene  
180 with a given function was within 50-200% of the median gene length. This resulted in a set  
181 of 34,527 *E. coli* genomes that were available for modeling.  
182

183 In order to reduce the size of the set of genomes for computing efficiency, while  
184 maintaining genomic diversity, genomes were clustered based on nucleotide k-mer  
185 similarity. A set of 100 core genes, defined as those corresponding to the protein families  
186 that were most highly conserved across the entire set of *E. coli* genomes, was computed  
187 (**Table S1**). Nucleotide 7-mer counts were computed for the core genes of each genome  
188 using KMC 2.3.0 (48) and the genomes were clustered based on their 7-mer distances using  
189 the agglomerative clustering function in scikit-learn (version 0.20.3) (49) using the  
190 parameters: `n_clusters='4000'`, `affinity='l1'`, and `linkage='average'`. From this, we selected a  
191 final set of 4,000 diverse *E. coli* genomes representing each cluster that was used for  
192 training and testing the models in this study (**Table S2**).  
193

194 Since the goal of this study is to predict whether a protein-encoding gene is present or  
195 absent within a given *E. coli* genome, we chose to use the PATRIC local protein families  
196 (PATtyFams) to describe this set (46). Genes encoding proteins that were members of the  
197 same protein family were considered to be orthologous. We computed the frequency of  
198 occurrence for each protein family across the entire starting set of 34,527 *E. coli* genomes  
199 and chose to model the protein families occurring in 10-90% of the genomes. This resulted  
200 in a total of 3,259 *E. coli* protein families that were modeled (**Table S3**). We chose not to  
201 model protein families occurring in less than 10% of the genomes to limit class imbalance  
202 and to keep the number of models tractable.  
203

### 204 Model generation

205 The set of 100 nearly universal core genes (described above) was chosen for generating the  
206 k-mer based feature sets for the models (**Table S1**). The genes were found in each *E. coli*

207 genome, and the nucleotide sequences were subdivided into canonical 7-mers using KMC  
208 version 2.3.0 (48). 7-mers were chosen because they train rapidly while retaining accuracy  
209 (**Table S4**). A matrix was created where the columns were the k-mers, the rows were the  
210 genomes, and each cell contained the counts of each k-mer. K-mers containing ambiguous  
211 nucleotides were not considered. A binary classifier was computed for each of the 3,259  
212 protein families described above, where the labels were the presence or absence of the  
213 family in each genome.

214  
215 Models were built using Extreme Gradient Boosting (XGBoost) version 0.81 (50) as  
216 described previously (36, 51). Unless otherwise stated, all models were evaluated using a  
217 10-fold cross validation, where 80% of the data were used for training, 10% for testing, and  
218 10% as a holdout set to monitor for overfitting in each fold. Model parameters were  
219 chosen based on tuning experiments for conserved gene models that were previously  
220 published (36). These included a maximum tree depth of 16 and a learning rate of 0.0625.  
221 Due to the high computing volume, unless otherwise stated, results are shown for the first  
222 five of ten folds.

223

## 224 **Environmental genomes**

225 A holdout set of genomes from 419 environmental *E. coli* isolates was used to evaluate the  
226 models that were trained on the public genomes (**Table S5**). *E. coli* isolates were collected  
227 from freshwater samples in rivers, streams, and Lake Macatawa in the Macatawa  
228 Watershed (Holland, Michigan, USA) between 2012 and 2019 as part of year-round water  
229 quality monitoring efforts. EPA Method 1603 (52) was used to monitor *E. coli* levels in the  
230 watershed and served as the basis for strain collection. Isolated colonies displaying  
231 morphology consistent with *E. coli* on mTEC plates were streaked for isolation on nutrient  
232 agar plates to obtain pure cultures of putative *Escherichia* strains. Purified isolates were  
233 archived as glycerol stocks and stored at -80°C for downstream genome sequencing. All  
234 strains were screened via standard biochemical identification tests to ensure consistency  
235 with *E. coli* phenotypes prior to sequencing. Genomic DNA extraction was performed with  
236 the DNeasy® PowerLyzer® Microbial Kit (Qiagen). Sequencing library preparation was  
237 performed with the Nextera XT DNA Library Prep kit (Illumina). Library QC was  
238 performed with the Qubit™ dsDNA HS Assay Kit (Invitrogen) and an Agilent 2200  
239 TapeStation system, using the High Sensitivity D5000 ScreenTape System (Agilent). Pooled  
240 libraries (24 per run) were sequenced on an Illumina MiSeq using the MiSeq Reagent Kit V2  
241 (500 cycle, PE 2x250), according to manufacturer instructions. Genomes were assembled  
242 and annotated using the BV-BRC assembly and annotation services (44).

243

## 244 **Subset analyses**

245 Several experiments were conducted to determine how models performed with less data.  
246 In order to evaluate model performances on a smaller number of genomes, clustering (as  
247 described above) was performed to generate sets that were 500, 1000, 2000, and 4000  
248 genomes in size, and modeling was subsequently performed on these representative  
249 genome sets. To evaluate model performances using fewer conserved protein families, the  
250 top 25, 50, and 75 most conserved genes were selected from the original set of 100

251 conserved genes (**Table S1**) and models were trained on each respective set. Model  
252 performances were then recorded as described above.

253

254

## 255 **Genomic comparisons**

256 Multi Locus Sequence Types (MLST) were computed for all genomes using the MLST tool  
257 version 2.21.0 developed by Torsten Seemann (<https://github.com/tseemann/mlst>), which  
258 uses the PubMLST database (53). The phylogenetic tree was computed based on a  
259 concatenated nucleotide sequence alignment of the genes corresponding to the five most  
260 conserved protein families in **Table S1**. Genes were aligned using MAFFT v7.130b (54).  
261 The alignment was curated by removing all inserts occurring in less than 5% of the genes,  
262 and poor quality sequences were removed by hand using the alignment editor JalView  
263 version 2.11.2.0(55). A tree was generated with FastTree version 2.1.7 using the  
264 generalized time reversible model for nucleotide sequences (56). Trees were rendered in  
265 iTOL (57). *Salmonella enterica* serovar Typhimurium LT2 was used as an outgroup for the  
266 tree (GenBank ID: AE006468.2).

267

## 268 **Data availability**

269 Genomes for environmental isolates have been deposited at SRA under bioprojects  
270 PRJNA923802 and PRJNA918992 . Modeling software is available on github  
271 <https://github.com/BV-BRC-dependencies/EColiVariableGeneModels>.

272

273

## 274 **Results**

### 275 **Predicting variable gene content**

276 In order to predict the presence or absence of variable genes across the set of *E. coli*  
277 genomes, we first determined a set of genes that were amenable to modeling. To do this,  
278 we defined orthologous genes as those that belong to the same PATRIC local protein family  
279 (46) across the *E. coli* genomes in the BV-BRC database (**Figure 1**) (44). The local protein  
280 families are restricted to each genus, and they are computed using the same set of  
281 signature amino acid k-mers that are used by the RAST annotation system to project  
282 protein functions (46). Overall, many of the variable genes encode proteins with  
283 uncharacterized functions. Since the objective was to be able to predict the presence or  
284 absence of variable genes regardless of their annotation status, the PATRIC protein family  
285 algorithm worked well because it places all proteins with “hypothetical” and  
286 uncharacterized functions into families using either signature k-mers or sequence  
287 similarity with BLAST (58), enabling the tracking of the these poorly annotated genes. We  
288 chose to exclude highly conserved genes occurring in greater than 90% of the genomes and  
289 rare genes occurring in less than 10% of the genomes in order to build better balanced  
290 models for predicting presence or absence. This resulted in a final set of 3,259 variable  
291 genes occurring in 10-90% of the *E. coli* genomes that were modeled in this study (**Figure**  
292 **2**).

293

294 Importantly, because this study is designed to detect the presence or absence of variable  
295 genes, the set of protein families modeled in this study differs considerably from previous  
296 work. For instance, there are only 179 protein families that are modeled in this study that  
297 are also used to predict genome completeness by the BV-BRC genome quality tool (16)  
298 (**Figure S1**). The set of genes used in this study encode a diverse set of proteins with a  
299 variety of strain-specific functions (**Table S3**). Overall, 679 of the genes encode proteins  
300 with functions that exist in a curated SEED annotation subsystem (17). Some of the more  
301 common functions in the set of 3,259 modeled families include components of secretion  
302 systems, fimbriae and flagella, toxins and antitoxins, and genes involved in transcriptional  
303 control. Many have annotations relating to horizontal gene transfer (e.g., phage,  
304 transposition, and plasmid conjugation related functions). Over 40% of the genes encode  
305 proteins with poorly annotated functions containing the terms, “hypothetical,”  
306 “uncharacterized,” “putative,” or “mobile element protein” (we note that in the SEED  
307 annotation schema the term “mobile element protein” is an outdated term that is more  
308 often synonymous with “hypothetical protein,” rather than a function demonstrated to be  
309 involved in horizontal gene transfer).

310  
311 In order to reliably predict the presence or absence of each variable gene in each *E. coli*  
312 genome, a set of 100 highly conserved genes, present in all of the *E. coli* genomes (**Table**  
313 **S1**) was used to generate a feature set of nucleotide 7-mer counts (**Figure 1**). One XGBoost  
314 classifier was built for each of the 3,259 variable genes to predict its presence or absence.  
315 The models were trained and tested on a high-quality set of 4,000 *E. coli* genomes that was  
316 down sampled from all of the *E. coli* genomes in the BV-BRC. The training set includes 534  
317 distinct MLSTs and 133 genomes that are untyped (53) (**Figure 3A, Table S2, Figure S2**).  
318 The F1 scores averaged across all 3,259 protein families was  $0.912 \pm 0.910-0.914$  ( $\pm 95\%$   
319 confidence interval over 5 folds), with median F1 score of 0.926 (**Table S3, Figure S3**).  
320 When the F1 scores are averaged by genome or MLST, rather than by protein family, we  
321 observe similar results with F1 scores equal to 0.944 [0.943-0.945] per genome and 0.918  
322 [0.913-0.923] per MLST (**Table 1, Figure 3C**).

323  
324 Although the high F1 scores with cross validation indicate that the models are robust,  
325 models built for longer genes could have higher accuracies than shorter genes because the  
326 *ab initio* gene callers have difficulty accurately predicting shorter open reading frames  
327 (59). Likewise, genes that occur more frequently across the training set may have  
328 distribution patterns that are more consistent with the phylogeny of the conserved genes  
329 that were used as features, making their models more accurate. This might explain why the  
330 F1 scores are slightly higher when they are averaged by genome or MLST, because the  
331 more commonly occurring families are contributing more to these averages. To assess  
332 these potential sources of error, we plotted the average F1 scores for each protein family  
333 versus the median protein length for the protein family members and observe a weak  
334 correlation between gene length and accuracy  $PCC = 0.173$  (**Figure 4A**). Similarly, when  
335 we plot F1 versus the occurrence of each family across the training set of 4,000 genomes,  
336 we observe a slightly upward trend in the average F1 scores with a  $PCC$  of 0.612. This  
337 trend is not dramatic, and the genes occurring least frequently, in 10-11% of the genomes,  
338 still have an average F1 score of 0.885 [0.866-0.904] (**Figure 4B**). Although these data  
339 indicate weak trends in model accuracy relating to protein length and abundance in the

340 training set, this does not appear to be a major source of bias that could explain the high F1  
341 scores that we observe.

342

### 343 **Models built with less data retain accuracy**

344 In order to understand how using less data influences the performance of the models, we  
345 first built models using 7-mers from the top 25, 50, and 75 core genes. As expected, the  
346 models that were based on 25 core genes performed slightly worse ( $F1 = 0.886 \pm 0.882-$   
347  $0.889$ , averaged by protein family) because they contain less information and gradually  
348 improved as the number of core genes was increased (**Figure 5A**). Likewise, we built  
349 models using the original set of 100 core genes as features, and gradually increased the size  
350 of the training set from 500 to 4,000 diverse *E. coli* genomes. The models trained on 500  
351 genomes had an F1 score of  $0.863 \pm 0.859-0.867$  (averaged by protein family), and the F1  
352 scores gradually increased beyond 0.9 as the models were trained with 4,000 genomes  
353 (**Figure 5B**). This improvement is likely due to the better representation of the variable  
354 genes across the training set. Overall, the data suggest that reliable models can be built  
355 with fewer conserved genes or training set genomes with a correspondingly modest  
356 decrease in performance. Unless otherwise stated, results reported in this study are for  
357 models built from a feature set of 100 core genes and a training set of 4,000 *E. coli*  
358 genomes.

359

### 360 **Horizontally transferred genes can be predicted**

361 Since the feature set for the models are based on conserved genes, it is possible that models  
362 for certain protein families outperform others due to their tight coupling to the phylogeny  
363 or may underperform due to the effects of horizontal gene transfer. When we examine the  
364 F1 scores based on the protein functions encoded by the variable genes, we find that the  
365 accuracy of the models is typically higher in genes that are well annotated (**Table 2, Table**  
366 **S3**). For instance, models for variable genes with functions occurring in subsystems ( $F1 =$   
367  $0.935 \pm 0.931-0.940$ ), or which have full Enzyme Commission (EC) numbers ( $F1 = 0.945 \pm$   
368  $0.937-0.952$ ) have significantly higher F1 scores than those that do not. Conversely, genes  
369 that are annotated with functions involved in horizontal gene transfer, including those  
370 encoding functions relating to transposable elements ( $F1 = 0.895 \pm 0.882-0.907$ ), phage  
371 elements ( $F1 = 0.872 \pm 0.868-0.876$ ), or conjugation and other plasmid-related functions  
372 ( $F1 = 0.824 \pm 0.814-0.834$ ) all had had significantly lower F1 scores than the genes that did  
373 not (**Table 2**). A total of 14 AMR-related protein families were modeled, and they have an  
374 average F1 score of 0.841 [0.814-0.869]. The average F1 scores for the AMR genes, and  
375 their non-uniform distributions over the genomes used in the study (**Figure S4**), indicate  
376 that these have similar characteristics to the other horizontally transferred genes that were  
377 modeled (**Table 3, Table S6**). We note that in most cases the pattern of occurrence for  
378 each AMR protein family does not tend to cluster with the clades of the phylogenetic tree.  
379 These results indicate that although protein families with horizontal gene transfer-related  
380 functions do have lower F1 scores than other variable genes, their presence or absence can  
381 still be reliably predicted ( $F1 > 0.8$ ).

382

## 383 **Model performance on an environmental holdout set**

384 Although the collection of public *E. coli* genomes is large, it is biased toward laboratory,  
385 surveillance, and clinical strains. We wanted to observe how well the models trained on  
386 the public genomes would extend to novel genomes. To do this, we sequenced a collection  
387 of 419 environmental *E. coli* isolates that were collected from freshwater environments.  
388 Importantly, none of these genomes previously existed in the public archives. Overall, the  
389 collection is comprised of 136 distinct MLSTs and 37 untyped genomes, and the  
390 distribution of MLSTs differs from that of the public collection (**Figure 3B, Figure S2**).  
391 When the models that were trained on the public data are applied to these genomes, we  
392 observe F1 scores of 0.880 [0.876-0.882] averaged by genome, 0.867 [0.862-0.873]  
393 averaged by MLST (**Figure 3D**), and 0.718 [0.712-0.724] averaged by protein family (**Table**  
394 **1**). The diverse genomes lacking an MLST designation, have a lower average F1 score of  
395 0.700 [0.693-0.706], and is likely due to their genetic diversity which has not been learned  
396 by the models. Likewise, the lower F1 scores averaged by protein family are likely due to  
397 the differences in distribution of the protein families across these genomes. For instance,  
398 approximately 38% of the protein families that existed in 10-90% of the public genomes  
399 occur in less than 10% of the genomes in the environmental collection (**Figure S5**). In  
400 other words, the distribution of these protein families deviates from the expectation of the  
401 models trained on the public data. However, since these these families are rare within the  
402 environmental collection, they are insufficient to dramatically alter the results when  
403 averaged by genome or MLST (**Figure S6**), both of which remain greater than 0.86. Overall,  
404 these results indicate that the models are robust for predicting variable gene content in  
405 holdout set of diverse environmental *E. coli* genomes.

## 408 **Discussion**

409 The public repositories contain an abundance of incomplete genomes and MAGs, but we  
410 currently lack tools for predicting the additional genes that they should encode. In this  
411 study, as a proof of concept, we used the nucleotide sequences of core genes to predict the  
412 variable gene content across *E. coli*. Overall, the average F1 scores were greater than 0.9  
413 over the training set of 4,000 diverse genomes, indicating that the data from the core genes  
414 is sufficient for predicting the presence or absence of many of the variable genes. When we  
415 looked at how the accuracy relates to protein functions, we found that genes that were well  
416 annotated, either belonging to a SEED subsystem or annotated as having complete EC  
417 numbers were more easily predicted, since they had significantly higher F1 scores than  
418 those that did not. Conversely, models for genes with functions associated with horizontal  
419 gene transfer had significantly lower F1 scores. Although this is unsurprising given that  
420 horizontal gene transfer moves these genes in patterns that do not necessarily match the  
421 phylogeny, it is noteworthy that genes with annotations containing the terms “plasmid”  
422 and “conjugation,” which was the category with the lowest average F1 score in our analysis,  
423 still had a remarkably high average F1 score of 0.824. This is likely due to the ability of  
424 XGBoost to track non-linear relationships. Another surprise was that the models for genes  
425 with protein functions containing the terms “hypothetical,” “uncharacterized,” and  
426 “putative” had average F1 scores of 0.902 and were only slightly lower than the F1 scores

427 for the set of families with curated annotations. This suggests that despite being poorly  
428 characterized, the occurrence of these sequences is rather easily predicted, implying that  
429 there is much more to learn about their distribution patterns and value to be added by  
430 elucidating their functions.

431  
432 Using a holdout set of 419 *E. coli* genomes from freshwater environmental isolates, the  
433 models retained extensibility with F1 scores of 0.880 and 0.867 averaged by genome and  
434 MLST respectively, indicating that the models work well even in diverse genomes. In both  
435 the training set and the holdout set, we observed slight correlations in the accuracy of each  
436 model and the underlying protein length and occurrence of each family. However, these  
437 trends were insufficient to explain the high F1 scores for the models. The influence of rare  
438 families was more dramatic in the holdout set lowering the F1 score averaged by protein  
439 family to 0.718. However, since almost 40% of the families occurred in less than 10% of  
440 the genomes in the holdout set, their per-genome effect was considerably smaller. Adding  
441 diverse genomes to the training set as they become available would eventually correct this  
442 issue.

443  
444 One limitation of this study is that by focusing on the set of variable genes occurring in 10-  
445 90% of the genomes, many of the rarely occurring genes were omitted. Although  
446 predicting the presence or absence of this massive and enigmatic set of genes is obviously  
447 desirable, this was done to control the study size and because these rare genes often lacked  
448 sufficient numbers to provide balanced sets for modeling. As long as the *E. coli* pan genome  
449 remains open and we continue to observe new genes with each new genome (41, 42), this  
450 will always be a problem, so predicting the presence or absence of the rarest genes may  
451 require a different modeling strategy. However, we expect that as the number of diverse  
452 genomes increases, the number of protein families that can be used to build balanced  
453 classifiers in the way that we did in this study will also continue to increase.

454  
455 Our highest quality set of models was generated using 100 core genes as features on a  
456 training set of 4,000 *E. coli* genomes and covered the set of protein families occurring in 10-  
457 90% of the genomes. This resulted in a collection of over 3,259 XGBoost models. This  
458 approach represented a rather significant outlay of computing resources, with each model  
459 taking approximately 4 minutes on an Intel Xeon Gold 6148 machine utilizing 128 cores,  
460 for a total of 7.6 days for computing the entire set. Although this experimental design is  
461 admittedly brute force, it is nevertheless tractable and could be extended to other well  
462 sequenced species. Indeed, unlike previous models that we have built for predicting AMR  
463 phenotypes using larger k-mer sizes and more complex matrices (36, 60, 61), these models  
464 are simple binary classifiers and have small memory footprints, and thus could be  
465 computed in parallel on a cluster with a modest amount of memory per node, rather than a  
466 high memory server. In designing this study, we attempted several other matrix designs  
467 and algorithms, including several deep learning approaches which had the potential to  
468 make the task more succinct. However, these attempts have been unsuccessful in our  
469 hands due to the size of the dataset, and ultimately the strategy of computing one classifier  
470 per family was successful. One way to reduce the computational burden might be to use  
471 fewer core genes or training set genomes. We showed that systematically reducing the size  
472 of the training set, while maintaining diversity, or using a smaller number of core genes for

473 the feature set resulted in modest losses in accuracy. These tradeoffs may be deemed  
474 acceptable in certain circumstances. Using this study as a proof of concept, we expect that  
475 future studies will find more elegant modeling solutions.

476  
477 In conclusion, we have found that it is possible to predict the presence or absence of a large  
478 number of the *E. coli* variable genes by building classifiers that use k-mers from a set of  
479 conserved genes. These models were highly accurate and worked even for families with  
480 hypothetical and unknown functions. This study provides a potential framework for  
481 predicting whether an incomplete genome or MAG should or should not be expected to  
482 contain a given gene, and has implications for the estimation of genome quality, the  
483 assessment of risk due to AMR and other virulence genes, and the ability to predict the  
484 presence of other important genes.

485

## 486 **Acknowledgements**

487 We thank Emily Dietrich for her careful editing and Bob Olson for technical assistance. This  
488 work was funded in part by the United States National Institute of Allergy and Infectious  
489 Diseases Bacterial and Viral Bioinformatics resource center award [Contract No.  
490 75N93019C00076] to PI Rick Stevens, and by the United States Defense Advanced  
491 Research Projects Agency iSENTRY Friend or Foe program award [Contract No.  
492 HR0011150042] to JJD, the National Science Foundation Awards [MCB-1616737 and DBI-  
493 1229585] to AAB, and the Herbert H. and Grace A. Dow Foundation. The funders had no  
494 role in study design, data collection and interpretation, or the decision to submit the work  
495 for publication.

496

497

## 498 **Competing Interests**

499 The authors declare no competing interests.

500

501

## 502 **References**

- 503 1. Douglas GM, Maffei VJ, Zaneveld JR, Yurgel SN, Brown JR, Taylor CM, Huttenhower  
504 C, Langille MG. 2020. PICRUSt2 for prediction of metagenome functions. *Nature*  
505 *Biotechnology* 38:685-688.
- 506 2. Iwai S, Weinmaier T, Schmidt BL, Albertson DG, Poloso NJ, Dabbagh K, DeSantis TZ.  
507 2016. Piphillin: improved prediction of metagenomic content by direct inference from  
508 human microbiomes. *PloS one* 11:e0166104.
- 509 3. Wemheuer F, Taylor JA, Daniel R, Johnston E, Meinicke P, Thomas T, Wemheuer B.  
510 2018. Tax4Fun2: a R-based tool for the rapid prediction of habitat-specific functional  
511 profiles and functional redundancy based on 16S rRNA gene marker gene sequences.  
512 *BioRxiv*:490037.

- 513 4. Bowman JS, Ducklow HW. 2015. Microbial communities can be described by metabolic  
514 structure: a general framework and application to a seasonally variable, depth-stratified  
515 microbial community from the coastal West Antarctic Peninsula. *PLoS one* 10:e0135868.
- 516 5. Jun S-R, Robeson MS, Hauser LJ, Schadt CW, Gorin AA. 2015. PanFP: pangenome-  
517 based functional profiles for microbial communities. *BMC research notes* 8:1-7.
- 518 6. Mongad DS, Chavan NS, Narwade NP, Dixit K, Shouche YS, Dhotre DP. 2021.  
519 MicFunPred: A conserved approach to predict functional profiles from 16S rRNA gene  
520 sequence data. *Genomics* 113:3635-3643.
- 521 7. Djemiel C, Maron P-A, Terrat S, Dequiedt S, Cottin A, Ranjard L. 2022. Inferring  
522 microbiota functions from taxonomic genes: a review. *GigaScience* 11.
- 523 8. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy T, Schulz  
524 F, Jarett J, Rivers AR, Eloe-Fadrosh EA. 2017. Minimum information about a single  
525 amplified genome (MISAG) and a metagenome-assembled genome (MIMAG) of  
526 bacteria and archaea. *Nature biotechnology* 35:725-731.
- 527 9. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. 2015. CheckM:  
528 assessing the quality of microbial genomes recovered from isolates, single cells, and  
529 metagenomes. *Genome research* 25:1043-1055.
- 530 10. Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov G,  
531 Kriventseva EV, Zdobnov EM. 2018. BUSCO applications from quality assessments to  
532 gene prediction and phylogenomics. *Molecular biology and evolution* 35:543-548.
- 533 11. Parra G, Bradnam K, Ning Z, Keane T, Korf I. 2009. Assessing the gene space in draft  
534 genomes. *Nucleic acids research* 37:289-297.
- 535 12. Sunagawa S, Mende DR, Zeller G, Izquierdo-Carrasco F, Berger SA, Kultima JR, Coelho  
536 LP, Arumugam M, Tap J, Nielsen HB. 2013. Metagenomic species profiling using  
537 universal phylogenetic marker genes. *Nature methods* 10:1196-1199.
- 538 13. Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, Delmont TO. 2015.  
539 Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* 3:e1319.
- 540 14. Nayfach S, Camargo AP, Schulz F, Eloe-Fadrosh E, Roux S, Kyrpides NC. 2021.  
541 CheckV assesses the quality and completeness of metagenome-assembled viral genomes.  
542 *Nature biotechnology* 39:578-585.
- 543 15. Davis JJ, Wattam AR, Aziz RK, Brettin T, Butler R, Butler RM, Chlenski P, Conrad N,  
544 Dickerman A, Dietrich EM. 2020. The PATRIC Bioinformatics Resource Center:  
545 expanding data and analysis capabilities. *Nucleic acids research* 48:D606-D612.
- 546 16. Parrello B, Butler R, Chlenski P, Olson R, Overbeek J, Pusch GD, Vonstein V, Overbeek  
547 R. 2019. A machine learning-based service for estimating quality of genomes using  
548 PATRIC. *BMC bioinformatics* 20:1-9.
- 549 17. Overbeek R, Olson R, Pusch GD, Olsen GJ, Davis JJ, Disz T, Edwards RA, Gerdes S,  
550 Parrello B, Shukla M. 2014. The SEED and the Rapid Annotation of microbial genomes  
551 using Subsystems Technology (RAST). *Nucleic acids research* 42:D206-D214.
- 552 18. Chklovski A, Parks DH, Woodcroft BJ, Tyson GW. 2022. CheckM2: a rapid, scalable  
553 and accurate tool for assessing microbial genome quality using machine learning.  
554 *bioRxiv*.
- 555 19. Kanehisa M, Goto S. 2000. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic  
556 acids research* 28:27-30.
- 557 20. Geller-McGrath D, Konwar KM, Edgcomb VP, Pachiadaki M, Roddy JW, Wheeler TJ,  
558 McDermott JE. 2022. MetaPredict: A machine learning-based tool for predicting

559 metabolic modules in incomplete bacterial genomes. bioRxiv  
560 doi:10.1101/2022.12.21.521254:2022.12.21.521254.

561 21. Antonopoulos DA, Assaf R, Aziz RK, Brettin T, Bun C, Conrad N, Davis JJ, Dietrich  
562 EM, Disz T, Gerdes S. 2019. PATRIC as a unique resource for studying antimicrobial  
563 resistance. *Briefings in bioinformatics* 20:1094-1102.

564 22. Alcock BP, Raphenya AR, Lau TT, Tsang KK, Boucharde M, Edalatmand A, Huynh W,  
565 Nguyen A-LV, Cheng AA, Liu S. 2020. CARD 2020: antibiotic resistome surveillance  
566 with the comprehensive antibiotic resistance database. *Nucleic acids research* 48:D517-  
567 D525.

568 23. Yin X, Jiang X-T, Chai B, Li L, Yang Y, Cole JR, Tiedje JM, Zhang T. 2018. ARGs-  
569 OAP v2. 0 with an expanded SARG database and Hidden Markov Models for  
570 enhancement characterization and quantification of antibiotic resistance genes in  
571 environmental metagenomes. *Bioinformatics* 34:2263-2270.

572 24. Bortolaia V, Kaas RS, Ruppe E, Roberts MC, Schwarz S, Cattoir V, Philippon A, Allesoe  
573 RL, Rebelo AR, Florensa AF. 2020. ResFinder 4.0 for predictions of phenotypes from  
574 genotypes. *Journal of Antimicrobial Chemotherapy* 75:3491-3500.

575 25. Feldgarden M, Brover V, Haft DH, Prasad AB, Slotta DJ, Tolstoy I, Tyson GH, Zhao S,  
576 Hsu C-H, McDermott PF. 2019. Validating the AMRFinder tool and resistance gene  
577 database by using antimicrobial resistance genotype-phenotype correlations in a  
578 collection of isolates. *Antimicrobial agents and chemotherapy* 63:e00483-19.

579 26. Hunt M, Mather AE, Sánchez-Busó L, Page AJ, Parkhill J, Keane JA, Harris SR. 2017.  
580 ARIBA: rapid antimicrobial resistance genotyping directly from sequencing reads.  
581 *Microbial genomics* 3.

582 27. Hunt M, Bradley P, Lapierre SG, Heys S, Thomsit M, Hall MB, Malone KM, Wintringer  
583 P, Walker TM, Cirillo DM. 2019. Antibiotic resistance prediction for *Mycobacterium*  
584 tuberculosis from genome sequence data with Mykrobe. *Wellcome open research* 4.

585 28. Liu B, Zheng D, Jin Q, Chen L, Yang J. 2019. VFDB 2019: a comparative pathogenomic  
586 platform with an interactive web interface. *Nucleic acids research* 47:D687-D692.

587 29. de Nies L, Lopes S, Busi SB, Galata V, Heintz-Buschart A, Laczny CC, May P, Wilmes  
588 P. 2021. PathoFact: a pipeline for the prediction of virulence factors and antimicrobial  
589 resistance genes in metagenomic data. *Microbiome* 9:1-14.

590 30. Arango-Argoty G, Garner E, Pruden A, Heath LS, Vikesland P, Zhang L. 2018.  
591 DeepARG: a deep learning approach for predicting antibiotic resistance genes from  
592 metagenomic data. *Microbiome* 6:1-15.

593 31. Rice EW, Wang P, Smith AL, Stadler LB. 2020. Determining hosts of antibiotic  
594 resistance genes: a review of methodological advances. *Environmental Science &*  
595 *Technology Letters* 7:282-291.

596 32. Haffiez N, Chung TH, Zakaria BS, Shahidi M, Mezbahuddin S, Maal-Bared R, Dhar BR.  
597 2022. Exploration of machine learning algorithms for predicting the changes in  
598 abundance of antibiotic resistance genes in anaerobic digestion. *Science of The Total*  
599 *Environment*:156211.

600 33. Sun Y, Clarke B, Clarke J, Li X. 2021. Predicting antibiotic resistance gene abundance in  
601 activated sludge using shotgun metagenomics and machine learning. *Water Research*  
602 202:117384.

- 603 34. McDermott PF, Davis JJ. 2021. Predicting antimicrobial susceptibility from the bacterial  
604 genome: a new paradigm for one health resistance monitoring. *Journal of Veterinary*  
605 *Pharmacology and Therapeutics* 44:223-237.
- 606 35. Moradigaravand D, Palm M, Farewell A, Mustonen V, Warringer J, Parts L. 2018.  
607 Prediction of antibiotic resistance in *Escherichia coli* from large-scale pan-genome data.  
608 *PLoS computational biology* 14:e1006258.
- 609 36. Nguyen M, Olson R, Shukla M, VanOeffelen M, Davis JJ. 2020. Predicting antimicrobial  
610 resistance using conserved genes. *PLoS computational biology* 16:e1008319.
- 611 37. Aytan-Aktug D, Nguyen M, Clausen PTLC, Stevens R, Aarestrup FM, Lund O, Davis J.  
612 2021. Predicting antimicrobial resistance using partial genome alignments. *Msystems*  
613 6:e00185-21.
- 614 38. Břinda K, Callendrello A, Ma KC, MacFadden DR, Charalampous T, Lee RS, Cowley L,  
615 Wadsworth CB, Grad YH, Kucherov G. 2020. Rapid inference of antibiotic resistance  
616 and susceptibility by genomic neighbour typing. *Nature microbiology* 5:455-464.
- 617 39. Her H-L, Wu Y-W. 2018. A pan-genome-based machine learning approach for predicting  
618 antimicrobial resistance activities of the *Escherichia coli* strains. *Bioinformatics* 34:i89-  
619 i95.
- 620 40. Ding W, Baumdicker F, Neher RA. 2017. panX: pan-genome analysis and exploration.  
621 *Nucleic Acids Research* 46:e5-e5.
- 622 41. Rasko DA, Rosovitz MJ, Myers GS, Mongodin EF, Fricke WF, Gajer P, Crabtree J,  
623 Sebaihia M, Thomson NR, Chaudhuri R, Henderson IR, Sperandio V, Ravel J. 2008. The  
624 pangenome structure of *Escherichia coli*: comparative genomic analysis of *E. coli*  
625 commensal and pathogenic isolates. *J Bacteriol* 190:6881-93.
- 626 42. Tettelin H, Riley D, Cattuto C, Medini D. 2008. Comparative genomics: the bacterial  
627 pan-genome. *Current opinion in microbiology* 11:472-477.
- 628 43. Touchon M, Hoede C, Tenaillon O, Barbe V, Baeriswyl S, Bidet P, Bingen E, Bonacorsi  
629 S, Bouchier C, Bouvet O. 2009. Organised genome dynamics in the *Escherichia coli*  
630 species results in highly diverse adaptive paths. *PLoS genetics* 5:e1000344.
- 631 44. Olson RD, Assaf R, Brettin T, Conrad N, Cucinell C, Davis James J, Dempsey  
632 Donald M, Dickerman A, Dietrich Emily M, Kenyon Ronald W, Kuscuglu M,  
633 Lefkowitz Elliot J, Lu J, Machi D, Macken C, Mao C, Niewiadomska A, Nguyen M,  
634 Olsen Gary J, Overbeek Jamie C, Parrello B, Parrello V, Porter Jacob S, Pusch  
635 Gordon D, Shukla M, Singh I, Stewart L, Tan G, Thomas C, VanOeffelen M, Vonstein  
636 V, Wallace Zachary S, Warren Andrew S, Wattam Alice R, Xia F, Yoo H, Zhang Y,  
637 Zmasek Christian M, Scheuermann Richard H, Stevens Rick L. 2022. Introducing the  
638 Bacterial and Viral Bioinformatics Resource Center (BV-BRC): a resource combining  
639 PATRIC, IRD and ViPR. *Nucleic Acids Research* doi:10.1093/nar/gkac1003.
- 640 45. Brettin T, Davis JJ, Disz T, Edwards RA, Gerdes S, Olsen GJ, Olson R, Overbeek R,  
641 Parrello B, Pusch GD. 2015. RASTtk: a modular and extensible implementation of the  
642 RAST algorithm for building custom annotation pipelines and annotating batches of  
643 genomes. *Scientific reports* 5:1-6.
- 644 46. Davis JJ, Gerdes S, Olsen GJ, Olson R, Pusch GD, Shukla M, Vonstein V, Wattam AR,  
645 Yoo H. 2016. PATtyFams: protein families for the microbial genomes in the PATRIC  
646 database. *Frontiers in microbiology* 7:118.
- 647 47. Sayers EW, Cavanaugh M, Clark K, Pruitt KD, Schoch CL, Sherry ST, Karsch-Mizrachi  
648 I. 2021. GenBank. *Nucleic acids research* 49:D92-D96.

- 649 48. Kokot M, Długosz M, Deorowicz S. 2017. KMC 3: counting and manipulating k-mer  
650 statistics. *Bioinformatics* 33:2759-2761.
- 651 49. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M,  
652 Prettenhofer P, Weiss R, Dubourg V. 2011. Scikit-learn: Machine learning in Python. *the*  
653 *Journal of machine Learning research* 12:2825-2830.
- 654 50. Chen T, Guestrin C. Xgboost: A scalable tree boosting system, p 785-794. *In* (ed),  
655 51. VanOeffelen M, Nguyen M, Aytan-Aktug D, Brettin T, Dietrich EM, Kenyon RW,  
656 Machi D, Mao C, Olson R, Pusch GD. 2021. A genomic data resource for predicting  
657 antimicrobial resistance from laboratory-derived antimicrobial susceptibility phenotypes.  
658 *Briefings in Bioinformatics* 22:bbab313.
- 659 52. United States Environmental Protection Agency. 2014. Method 1603: *Escherichia coli*  
660 (*E. coli*) in Water by Membrane Filtration Using Modified membrane-Thermotolerant  
661 *Escherichia coli* Agar (Modified mTEC). United States Environmental Protection  
662 Agency, Washington, DC, USA.
- 663 53. Jolley KA, Maiden MC. 2010. BIGSdb: scalable analysis of bacterial genome variation at  
664 the population level. *BMC bioinformatics* 11:1-11.
- 665 54. Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:  
666 improvements in performance and usability. *Molecular biology and evolution* 30:772-  
667 780.
- 668 55. Waterhouse AM, Procter JB, Martin DM, Clamp M, Barton GJ. 2009. Jalview Version  
669 2—a multiple sequence alignment editor and analysis workbench. *Bioinformatics*  
670 25:1189-1191.
- 671 56. Price MN, Dehal PS, Arkin AP. 2010. FastTree 2—approximately maximum-likelihood  
672 trees for large alignments. *PloS one* 5:e9490.
- 673 57. Letunic I, Bork P. 2007. Interactive Tree Of Life (iTOL): an online tool for phylogenetic  
674 tree display and annotation. *Bioinformatics* 23:127-128.
- 675 58. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL.  
676 2009. BLAST+: architecture and applications. *BMC bioinformatics* 10:1-9.
- 677 59. Tripp HJ, Sutton G, White O, Wortman J, Pati A, Mikhailova N, Ovchinnikova G, Payne  
678 SH, Kyrpides NC, Ivanova N. 2015. Toward a standard in structural genome annotation  
679 for prokaryotes. *Standards in Genomic Sciences* 10:1-9.
- 680 60. Nguyen M, Brettin T, Long S, Musser JM, Olsen RJ, Olson R, Shukla M, Stevens RL,  
681 Xia F, Yoo H. 2018. Developing an in silico minimum inhibitory concentration panel test  
682 for *Klebsiella pneumoniae*. *Scientific reports* 8:1-11.
- 683 61. Nguyen M, Long SW, McDermott PF, Olsen RJ, Olson R, Stevens RL, Tyson GH, Zhao  
684 S, Davis JJ. 2019. Using machine learning to predict antimicrobial MICs and associated  
685 genomic features for nontyphoidal *Salmonella*. *Journal of Clinical Microbiology*  
686 57:e01260-18.
- 687  
688

689

690 **Tables**691 **Table 1.** Macro F1 scores averaged by genome, MLST, and protein family.

<b>Averaged by:</b>	<b>Training set (4,000 genomes)</b>	<b>Hold-out set (419 genomes)</b>
Genome	0.944 [0.943-0.945]	0.880 [0.876-0.882]
MLST	0.918 [0.913-0.923]	0.867 [0.862-0.873]
Protein family	0.912 [0.910-0.914]	0.718 [0.712-0.724]

692

693

694

695 **Table 2.** Commonly occurring protein functions in the set of 3,259 modeled protein  
696 families with their average F1 scores<sup>1</sup>.

<b>Annotations</b>	<b>With the annotation</b>		<b>Without the annotation</b>	
	<b>number</b>	<b>avg F1</b>	<b>number</b>	<b>avg F1</b>
Hypothetical etc. proteins <sup>2</sup>	1343	0.902 [0.898-0.906]	1916	0.918 [0.915-0.921]
Occurring in subsystems	679	0.935 [0.931-0.940]	2580	0.905 [0.903-0.908]
With the term "phage"	513	0.872 [0.868-0.876]	2746	0.919 [0.916-0.922]
With complete EC numbers	245	0.945 [0.937-0.952]	3014	0.909 [0.906-0.912]
Transporters	166	0.949 [0.940-0.959]	3093	0.910 [0.907-0.912]
Membrane proteins	133	0.954 [0.945-0.962]	3126	0.910 [0.907-0.912]
With the term "secretion"	119	0.965 [0.959-0.971]	3140	0.910 [0.907-0.912]
Transcriptional regulation <sup>2</sup>	97	0.945 [0.934-0.956]	3162	0.911 [0.908-0.913]
With the term "fimbriae"	78	0.974 [0.964-0.984]	3181	0.910 [0.908-0.913]
Toxins and anti-toxins	70	0.932 [0.917-0.946]	3189	0.911 [0.909-0.914]
With the terms "transposase" or "transposon"	67	0.895 [0.882-0.907]	3192	0.912 [0.910-0.915]
With the terms "conjugation" or "plasmid"	59	0.824 [0.814-0.834]	3200	0.913 [0.911-0.916]
Relating to flagellar function	39	0.948 [0.943-0.952]	3220	0.911 [0.909-0.914]

697 <sup>1</sup>The average is reported with the 95% confidence interval698 <sup>2</sup>Containing the terms "hypothetical", "uncharacterized," "putative," or "mobile element"699 <sup>3</sup>Containing the terms, transcriptional "activator," "repressor," "regulator," or  
700 "antiterminator"

701

702

703

704

**Table 3.** AMR protein families modeld in this study with their respective F1 scores.

Protein Family	F1 Score	Frac. Genomes with Protein	BV-BRC Annotation
PLF_561_00005992	0.794 [0.756-0.833]	0.238	Aminoglycoside 3''-nucleotidyltransferase (EC 2.7.7.-) => ANT(3'')-Ia (AadA family)
PLF_561_00057308	0.798 [0.753-0.842]	0.153	Aminoglycoside 3''-nucleotidyltransferase (EC 2.7.7.-) => ANT(3'')-Ia (AadA family)
PLF_561_00005448	0.817 [0.782-0.853]	0.346	Aminoglycoside 3''-phosphotransferase (EC 2.7.1.87) => APH(3'')-I
PLF_561_00005227	0.812 [0.779-0.844]	0.350	Aminoglycoside 6-phosphotransferase (EC 2.7.1.72) => APH(6)-Ic/APH(6)-Id
PLF_561_00009406	0.791 [0.759-0.824]	0.137	Aminoglycoside N(3)-acetyltransferase (EC 2.3.1.81) => AAC(3)-II,III,IV,VI,VIII,IX,X
PLF_561_00009579	0.836 [0.809-0.862]	0.160	Chloramphenicol/florfenicol resistance, MFS efflux pump => FloR family
PLF_561_00013716	0.853 [0.834-0.872]	0.189	Class A beta-lactamase (EC 3.5.2.6) => CTX-M family, extended-spectrum
PLF_561_00004782	0.831 [0.813-0.850]	0.405	Class A beta-lactamase (EC 3.5.2.6) => TEM family
PLF_561_00004401	0.989 [0.983-0.996]	0.629	Colicin E2 tolerance protein CbrC-like protein => CbrC
PLF_561_00011342	0.850 [0.819-0.881]	0.220	Macrolide 2'-phosphotransferase => Mph(A) family
PLF_561_00003770	0.915 [0.881-0.948]	0.651	Small multidrug resistance (SMR) efflux transporter => EmrE, broad substrate specificity
PLF_561_00013078	0.841 [0.823-0.859]	0.284	Small multidrug resistance (SMR) efflux transporter => QacE delta 1, quaternary ammonium compounds
PLF_561_00006144	0.814 [0.791-0.837]	0.353	Tetracycline resistance, MFS efflux pump => Tet(A)
PLF_561_00006969	0.837 [0.803-0.871]	0.184	Tetracycline resistance, MFS efflux pump => Tet(B)

705

706

707 **Figures**

708 **Figure 1.** Workflow used in this study. All *E. coli* genomes were downloaded from BV-BRC,  
709 filtered for genome quality, and down selected using hierarchical clustering. All *E. coli*  
710 genus-level protein families were also taken from BV-BRC and downselected for those that  
711 occur in 10-90% of the genomes in order to enable building balanced models. Then one  
712 matrix was built per family using 7-mer nucleotide frequencies from a set of 100 core  
713 genes held in common by all of the genomes. Finally, one binary XGBoost classifier was  
714 built to predict the presence or absence of each protein family a given *E. coli* genome.  
715

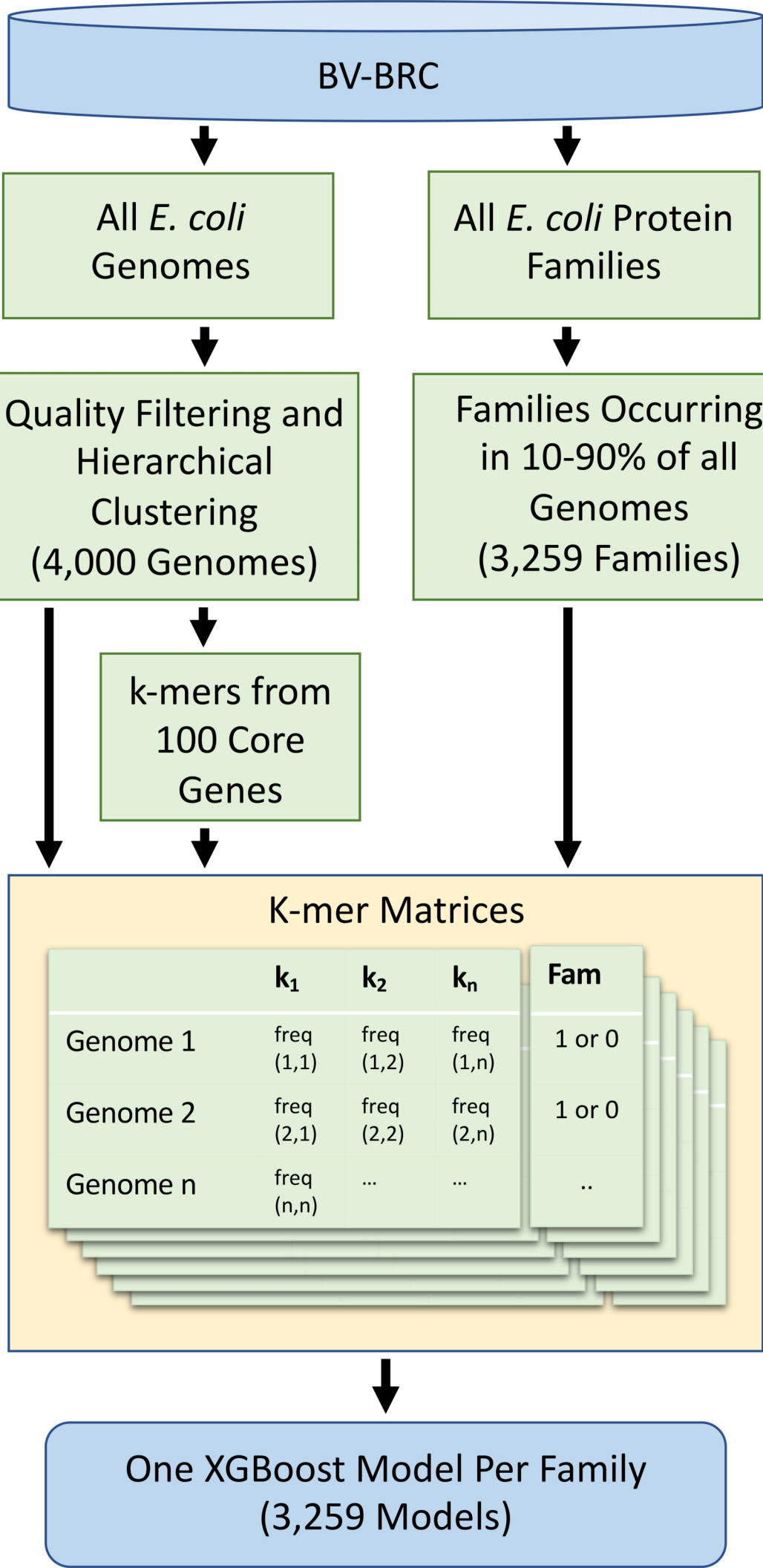
716  
717 **Figure 2.** Histogram of the frequency of occurrence for all protein families found in the *E.*  
718 *coli* genomes in the BV-BRC. Models in this study were built to predict the presence or  
719 absence of protein families occurring in 10-90% of all genomes, shown in blue.  
720

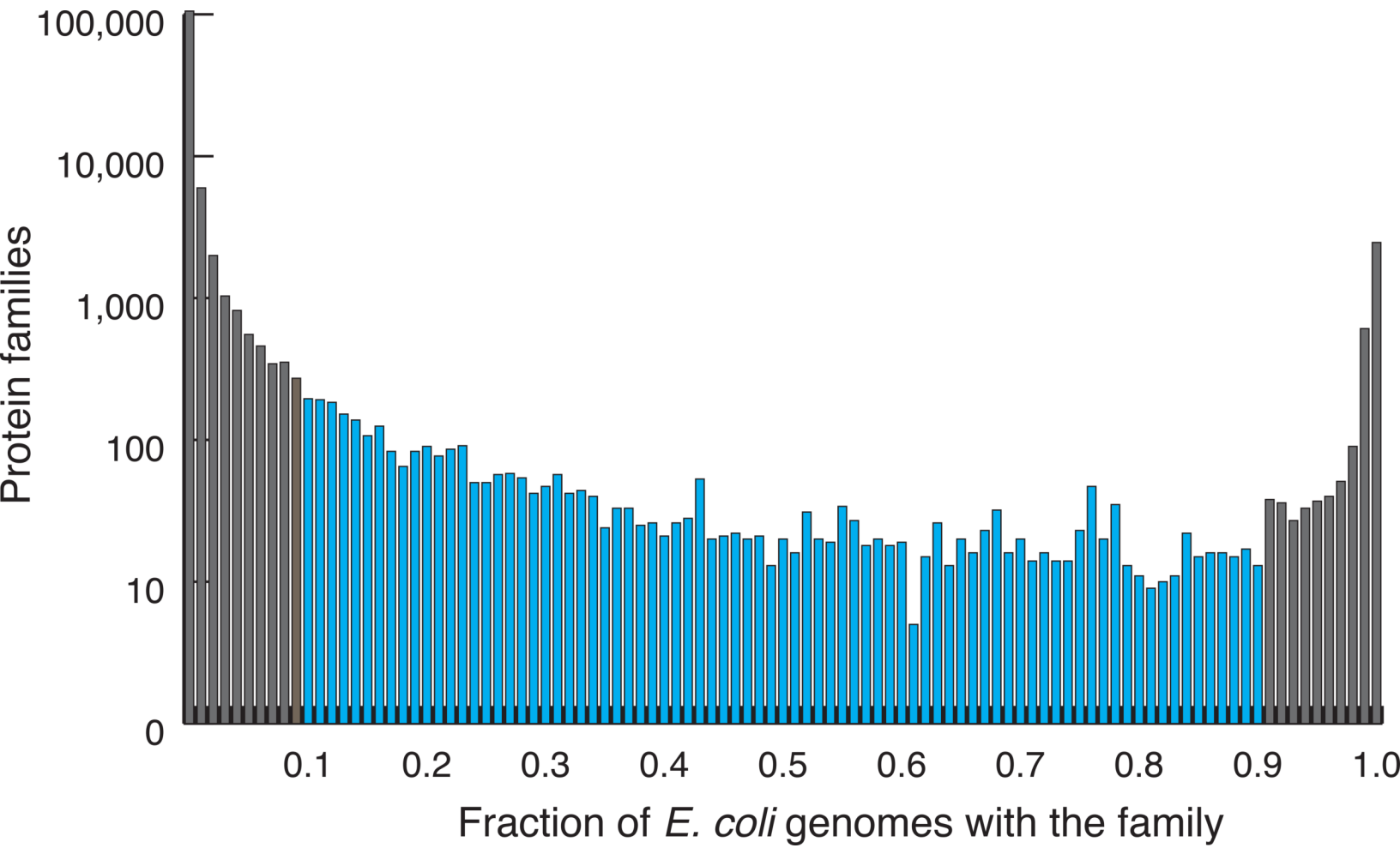
721  
722 **Figure 3.** MLST distributions and F1 scores averaged by MLST. A) Histogram of the 20  
723 most frequently occurring MLSTs in the training set of 4000 diverse genomes from the BV-  
724 BRC; B) Histogram of the 20 most frequently occurring MLSTs in the holdout set of 419  
725 environmental genomes; C) F1 scores averaged by MLST for the set of 4,000 BV-BRC  
726 genomes; D) F1 scores averaged by MLST for the holdout set of environmental genomes.  
727 Error bars depict the 95% confidence intervals. The MLST labeled with a dash represent  
728 all genomes with undetermined MLSTs in each set.  
729

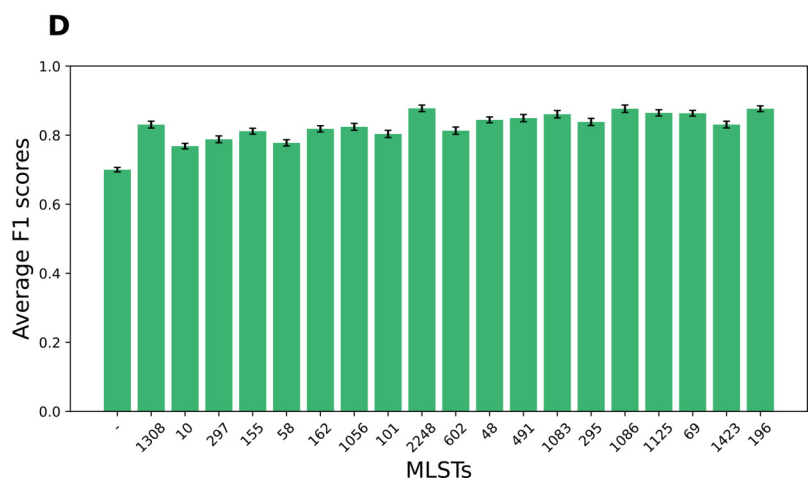
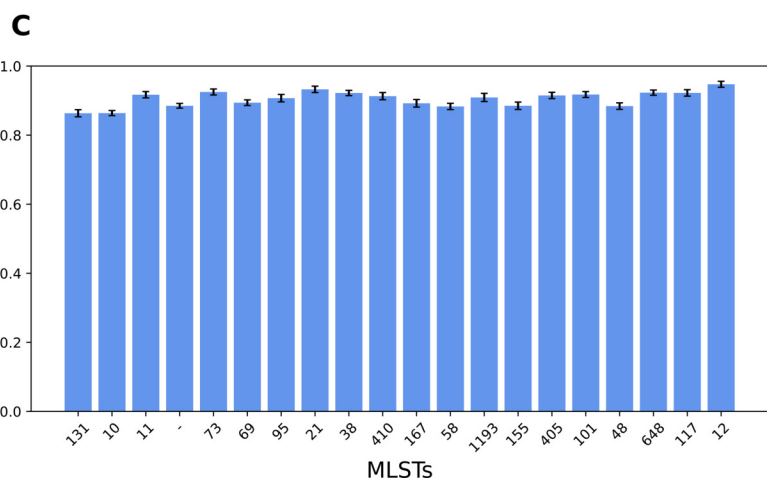
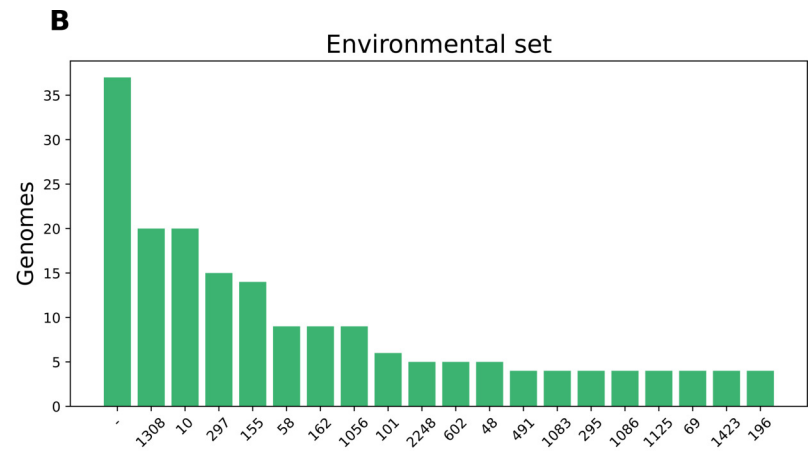
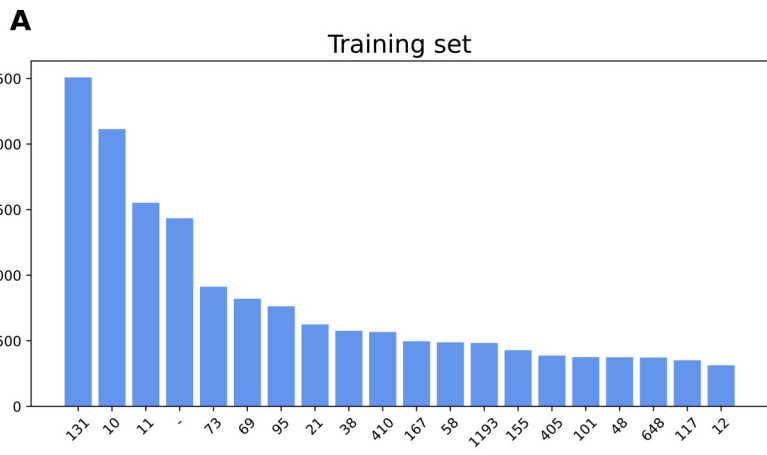
730  
731 **Figure 4.** Average F1 scores versus protein length and protein family occurrence. A) F1  
732 scores averaged by protein family plotted by the median protein length for all family  
733 members; B) F1 scores averaged by protein family versus the fraction of *E. coli* genomes in  
734 the training set containing a member of the given protein family. Gray bars depict the 95%  
735 confidence intervals.  
736

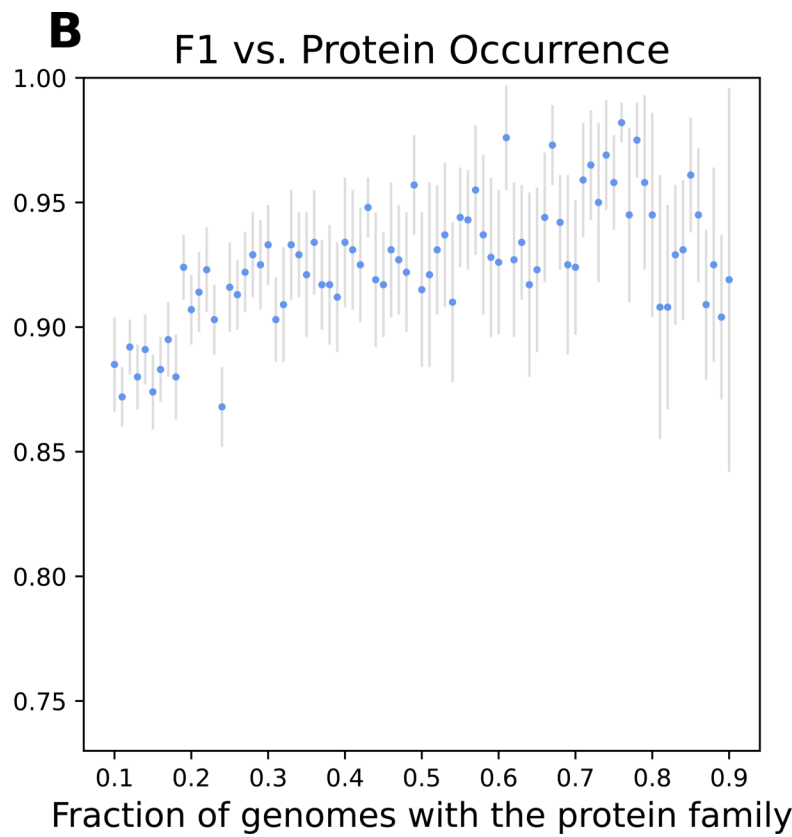
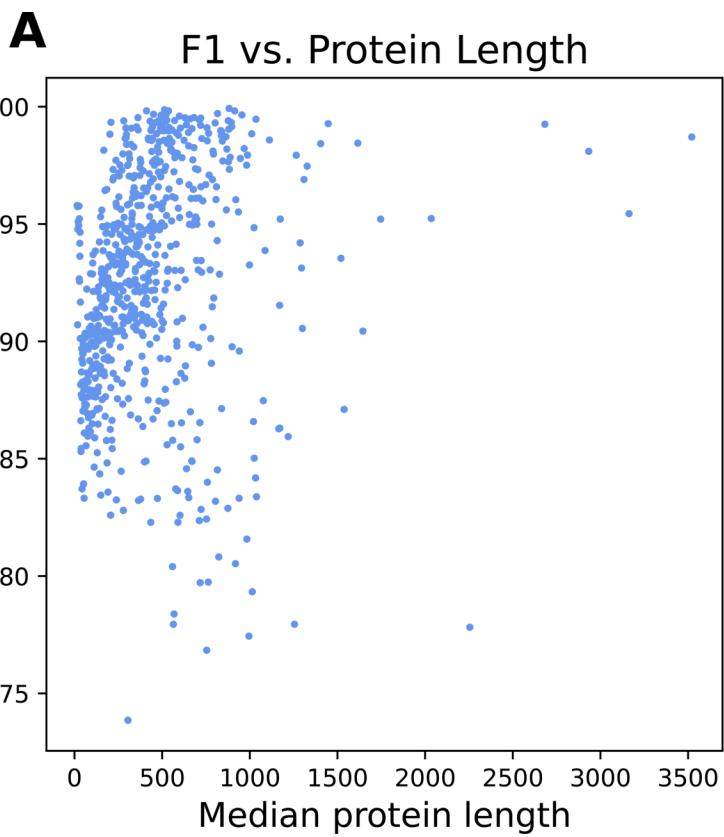
737  
738 **Figure 5.** F1 scores versus number of core genes and number of genomes used to train the  
739 models. A) F1 scores averaged by protein family versus the number of core genes used to  
740 train the model; B) F1 scores averaged by protein family versus the number of diverse *E.*  
741 *coli* genomes used to train the models. Gray bars depict the 95% confidence intervals over  
742 5 folds.  
743

744  
745  
746  
747  
748

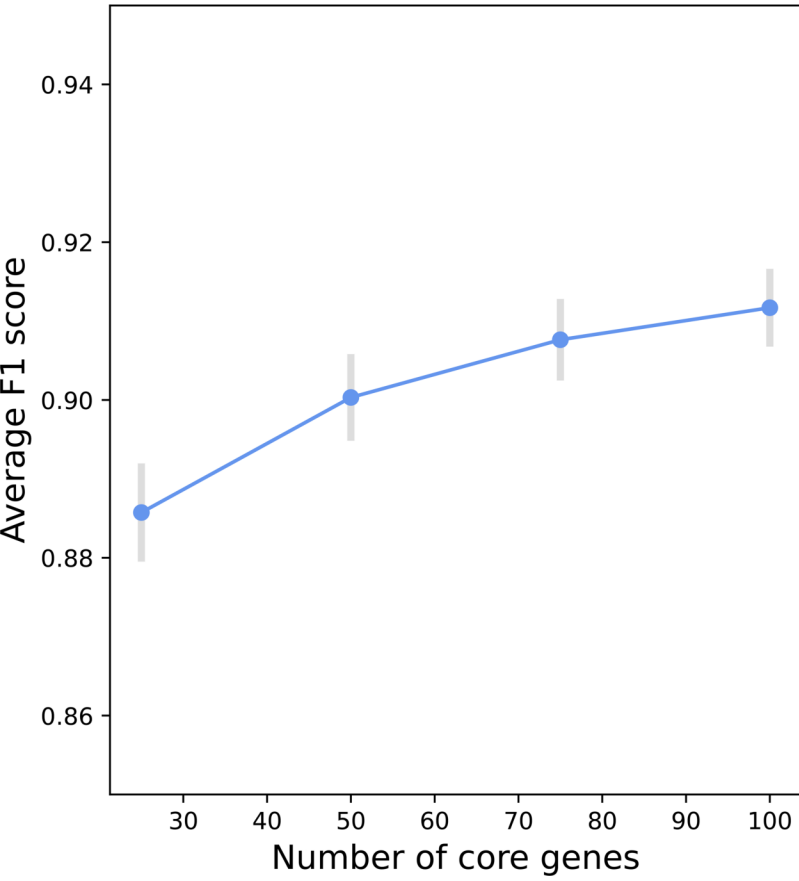








**A** F1 vs. Number of Core Genes



**B** F1 vs. Number of Genomes

