# Machine learning to predict the source of campylobacteriosis using whole genome data

Nicolas Arning*[1], Samuel K. Sheppard[2], David A. Clifton[3] and Daniel J. Wilson[1]

[1] Big Data institute, Nuffield Department of Population Health, University of Oxford, Li Ka Shing Centre for Health Information and Discovery, Old Road Campus, Oxford, OX3 7LF, UK

[2] The Milner Centre of Evolution, Department of Biology & Biochemistry, University of Bath, Claverton Down, Bath, BA2 7AZ, UK

[3] Institute of Biomedical Engineering, Department of Engineering Science, University of Oxford, OX3 7DQ, UK

## Abstract

Campylobacteriosis is among the world's most common foodborne illnesses, caused predominantly by the bacterium *Campylobacter jejuni*. Effective interventions require determination of the infection source which is challenging as transmission occurs via multiple sources such as contaminated meat, poultry, and drinking water. Strain variation has allowed source tracking based upon allelic variation in multi-locus sequence typing (MLST) genes allowing isolates from infected individuals to be attributed to specific animal or environmental reservoirs. However, the accuracy of probabilistic attribution models has been limited by the ability to differentiate isolates based upon just 7 MLST genes. Here, we broaden the input data spectrum to include core genome MLST (cgMLST) and whole genome sequences (WGS), and implement multiple machine learning algorithms, allowing more accurate source attribution. We increase attribution accuracy from 64% using the standard iSource population genetic approach to 71% for MLST, 85% for cgMLST and 78% for kmerized WGS data using machine learning. To gain insight beyond the source model prediction, we use Bayesian inference to analyse the relative affinity of *C. jejuni* strains to infect humans and identified potential differences, in source-human transmission ability among clonally related isolates in the most common disease causing lineage (ST-21 clonal complex). Providing generalizable computationally efficient methods, based upon machine learning and population genetics, we provide a scalable approach to global disease surveillance that can continuously incorporate novel samples for source attribution and identify fine-scale variation in transmission potential.

## Author summary

*C. jejuni* are the most common cause of food-borne bacterial gastroenteritis but the relative contribution of different sources are incompletely understood. We traced the origin of human *C. jejuni* infections using machine learning algorithms that compare the DNA sequences of bacteria sampled from infected people, contaminated chickens, cattle, sheep, wild birds and the environment. This approach achieved improvement in accuracy of source attribution by 33% over existing methods that use only a subset of genes within the genome and provided evidence for the relative contribution of different infection sources. Sometimes even very similar bacteria showed differences, demonstrating the value of basing analyses on the entire genome when developing this algorithm that can be used for understanding the global epidemiology and other important bacterial infections.

# Introduction

2   *Campylobacter jejuni* and *Campylobacter coli* are among the most common causes of gastroenteritis

3   globally and are responsible for approximately nine million annual cases in the European Union (1,2).

4   These zoonotic bacteria are a common commensal constituent of the gut microbiota of bird and

5   animal species (3,4) but cause serious infections in humans. Symptoms include nausea, fever,

6   abdominal pain, and severe diarrhoea, with potential for the development of debilitating, and

7   sometimes fatal, sequelae (5,6). Various infection sources have been identified including animal

8   faeces, contaminated drinking water and especially raw or under-cooked poultry and other meats (7).

9   However, effectively combating disease requires a detailed understanding of the relative contribution

10  of different sources to human infection.

11

12  As in many other bacterial species, *Campylobacter* populations represent diverse assemblages of

13  strains (3,8–10). Within this structured population, some lineages are more commonly observed in

14  particular host species (3,4,11). Because of this host association, DNA sequence comparisons of

15  bacteria from human gastroenteritis and potential reservoir populations have potential to reveal the

16  infection source. This has identified contaminated poultry as a major source of human infection

17  (12,13). Based on the body of evidence including DNA sequence analysis (14), targeted interventions

18  have been implemented, including improved biosecurity measures on poultry farms, which have

19  halved recorded campylobacteriosis cases in New Zealand (15,16).

20

21  Extending the principal of linking source-sink populations using genotype data, methods have been

22  developed to attribute *C. jejuni* to the likely source based on bacterial gene frequencies in potential

23  reservoir populations (17,18). Among the most common genotyping approaches for *C. jejuni* has been

24    multi-locus sequence typing (MLST) that catalogues DNA sequence variation across seven

25    housekeeping genes that are common to all strains (19,20). Isolates with identical alleles at all loci are

26    assigned to the same sequence type (ST) and those with identical sequences at most or all loci are

27    grouped within the same clonal complex (CC). Using these data, and allele frequencies, it has been

28    possible to probabilistically assign clinical isolates (STs and CCs) to host source using source attribution

29    models such as the asymmetric island model implemented in *iSource* (17) and the Bayesian population

30    assignment model STRUCTURE (18,21). Both methods have been instructive in estimating the relative

31    contribution of a range of domestic and wild animal hosts to human infection, with poultry often

32    identified as the principal source of human campylobacteriosis across different regions and countries

33    (17,18,22–25).

34

35    There are two main limitations when using genotype data to for bacterial source attribution. The first

36    is that the ability to attribute is only as good as the degree of genotype segregation. For example, in

37    *C. jejuni* there are host restricted genotypes (3,26) that can be readily attributed to a given host source

38    when observed in human infections, as well as ecological generalists (27,28) that have relatively

39    recently transitioned between hosts and cannot therefore be attributed with confidence (29). While

40    host switching potentially imposes a biological constraint on quantitative attribution models, the

41    second limitation is far more tractable. Specifically, most current source attribution methods are

42    subject to limitations imposed by the underlying data. Reflecting the technology of the time, MLST-

43    based source attribution is based only on a small fraction of the genome (approximately 0.2% for *C.*

44    *jejuni* (25)) and there is considerable potential for better strain differentiation using current

45    techniques.

46

47    The increasing availability of large whole genome sequence (WGS) datasets has greatly enhanced

48    analyses of bacterial population structure and diversity (30). However, exploiting the full information

49    can be challenging due to variable gene content and the complexity of interpreting the short reads

50    produced by next generation sequencing. Notwithstanding this, some studies have attempted to

51    overcome the limited discriminatory power of MLST in attribution studies by screening WGS data to

52    identify elements (SNPs and genes) that segregate by host (31–33). Using these host segregating

53    markers as input data has improved the resolution of existing attribution models, including

54    STRUCTURE, and provided information about potential infection reservoirs and the UK and France.

55    However, using bespoke marker selection approaches with software designed for MLST data does not

56    maximize the potential of WGS data for source attribution.

57

58    Here, we present a machine learning approach using WGS data to predict the source of human *C. jejuni*

59    infection. This has two principal advantages over existing techniques. First, building on WGS-based

60    machine learning source attribution approaches applied to *Salmonella enterica* and *Escherichia coli*

61    (34,35), we take an agnostic approach to identify which machine learning tool performs best from a

62    broad range of available algorithms. Second, we use a WGS input capture approach using data types

63    deposited in public databases allowing the analysis of existing MLST, core-genome MLST and WGS

64    datasets and the reuse of data for continuous updatable monitoring in a generalizable framework.

65    Thus, we aimed to overcome limitations of the currently available methods and use the output to

66    investigate the infective potential of *C. jejuni* strains.

67

## Methods

### Dataset acquisition

70    A total of 5,798 *C. jejuni* and *C. coli* genomes isolated from various sources and host species were

71   available on the public database for molecular typing and microbial genome diversity: PubMLST

72   (https://pubmlst.org/) (S1 Table). WGS data corresponded to MLST ST and CC designations as well as

73   core genome (cg) MLST classes. The dataset was divided into training (75%) and testing (25%) sets

74   using phylogeny-aware sorting, wherein all members of one ST were sorted entirely into either training

75   or testing sets (S1 Table). The ST based sorting accounts for the phylogenetic non-independence of

76   samples (36). To allow for sufficient sample sizes per reservoir population (hereafter "class"), only the

77   five most prevalent classes for MLST and cgMLST were used (chicken, cattle, sheep, wild bird and

78   environment). For farm animals the classes "chicken" and "chicken offal or meat" were combined to

79   "chicken" (likewise for sheep and cattle), whilst "environment", "sand" and "river water" were

80   combined into "environment", consistent with previous studies (18,37).

81

## Feature engineering

83   The allelic profiles of MLST and cgMLST were used directly. To potentially exploit the gradient of

84   separation encoded in the sequences underlying the MLST allelic profiles, we downloaded the

85   underlying allele sequences and encoded the nucleotides as dummy variables and k-mers (k=21) using

86   DSK (38). DSK was also used for encoding the WGS as k-mers. Using k=21 led to a prohibitively large

87   input vector due to the number of unique k-mers found in all genomes (109,675,176). We reduced the

88   number of k-mers by applying a variance threshold where k-mers which were present or absent in

89   more than 99% of the samples were discarded, reducing the numbers of unique k-mers to 7,285,583.

90   Furthermore, we performed feature selection by testing the dependence of the source labels on every

91   individual k-mer using the Chi-Square statistic. To avoid data-leakage we only performed the feature

92   selection using the training data and labels to select the 100,000 k-mers with the highest score.

93

5

## Algorithm training

94

95   All machine learning and deep learning was performed in Python (for a list of all algorithms see Figure

96   1). The xgboost library (39) was used for the gradient boosting classifiers with all other machine

97   learners implemented in scikit-learn (40). The hyper-parameters for each classifier were chosen using

98   Cartesian grid search on five-fold cross-validation of the training set. The Keras library

99   (https://github.com/keras-team/keras) was used to construct deep learning algorithms aimed at

100  supplying a wide range of commonly used architectures. We found this to work best, empirically, given

101  that there is no principled means of architecture selection for such models. Specifically: (i) A recurrent

102  neural network consisting of a layer with 64 gated recurrent units, a 50% dropout layer and Rectified

103  Linear Unit (ReLU) activation layer; (ii) A 1-dimensional convolutional network with two convolutional

104  layers of kernel size 3 and 5 respectively and 30 filters, both followed by 50% dropout layers and a

105  ReLU layer; (iii) A Long short-term memory network consisting of one LSTM layer with 64 units and a

106  50% dropout layer; (iv) A Shallow dense network with one dense layer with 64 units followed by a 50%

107  dropout layer and a ReLU activation layer; (v) A Deep dense network with 6 dense layers starting with

108  128 units and halving units with each successive layer. All individual dense layers are followed by a

109  50% dropout layer and a ReLU layer.

110

111  To all deep learning architectures, we added an output layer comprising a dense layer with soft-max

112  activation with one unit for every class. We encoded the labels as dummy variables and used

113  categorical cross-entropy as a loss function together with the Adam optimiser (41). Cyclical learning

114  rates were used with a maximum learning rate of 0.1 and a minimum learning rate of 0.0001 to

115  overcome local minima. The accuracy on the test set was measured at every epoch and the overall

116  best performing weights were stored as a checkpoint. The data was deployed in batches of 128

117  samples with every batch randomly undersampled so that each class was represented in equal

6

118    proportions. The training was run for 500 generations with early stopping after 50 generations.

119

## Algorithm testing

121    Both machine learning and deep learning were tested on the same 25% test set. The original data were

122    skewed in source composition by ratios which did not necessarily reflect source origin of infection. We

123    therefore used two methods to rebalance the classes in testing. The first test set featured an even

124    distribution of classes, whereas the second undersampled the over-abundant chicken-origin genomes

125    to emulate relative contribution to human disease. We used the ratios predicted by Wilson et al. (12),

126    where *Campylobacter* genomes from chickens were 1.61 times more common than those from cattle.

127    In both methods, rebalancing the classes was achieved by undersampling, which we repeated 200

128    times with replacement and averaged the accuracy over all iterations whilst also recording the

129    variance. For performance metrics we registered accuracy, precision (positive predictive value), recall

130    (sensitivity), F1, negative predictive value, specificity and speed. Speed was measured relative to other

131    classifiers where a scale was defined with 0 being the slowest classifier and 1 being the quickest and

132    all intermediate values being normalised within these confines. For comparison to previous methods,

133    iSource was applied to the test dataset (17). Having established that XGBoost on cgMLST was the best

134    performing source attribution method, we retrained the classifier with both training and testing data

135    and applied it to all 15,988 human cgMLST samples available on the PubMLST database. The prediction

136    took 892 milliseconds on a Dell OptiPlex 7060 desktop using ten threads on an Intel Core i7-8700 CPU

137    and 16 GB RAM.

138

## Phylogenetic analysis

140    We defined the generalist index as the number of sources the ST was found in across all isolates in the

141  dataset, which included additional samples for which only MLST data was available (S1 Table). We

142  built a phylogeny of CC21 genomes from both source-associated and human isolates using Neighbour

143  Joining, based on pairwise hamming distances of k-mer presence/absence in the WGS dataset, as

144  described by Hedge and Wilson (42). We used TreeBreaker to infer the evolution of phenotypes across

145  the phylogenetic tree of ST-21 and the most closely related sequence types. The known labels of the

146  source-associated samples were used as phenotypic information for input into TreeBreaker (43)

147  together with the phylogeny of CC21. TreeBreaker was run for 5,500,000 iterations with 500,000

148  iterations as burn-in and 1000 iterations between sampling. The phylogenetic trees were visualised

149  with Microreact (44) and arranged alongside the results of TreeBreaker in Inkscape.

150

# Results and Discussion

## Machine learning outperforms popular attribution models for MLST data

153  In order to anchor our source attribution performance to previous efforts, we compared results using

154  the machine learning classifiers to source probabilities estimated using the asymmetric island model

155  implemented in iSource, which is based on MLST and the most commonly used source attribution

156  method to date (45). The best performing machine learner on the MLST allelic profile was a random

157  forest (61.9%/68.5% balanced/unbalanced) which performed slightly better than iSource (61%/64%)

158  (Figure 1). Since loci within allelic profiles are deemed either to match or not, and underlying

159  nucleotides sequences are ignored, we investigated whether exploiting the gradient of nucleotide

160  differentiation would lead to better attribution. We used dummy variables and generated k-mers from

161  the sequences underlying the MLST allele labels. The additional feature encodings boosted the top

162  achieving accuracies on MLST to 67.9%/70.7% from dummy variables and 63%/67.5% from k-mers,

163  showing the value of the additional nucleotide-level information.

**Figure 1:** A heatmap showing classifier performance on the class balanced (A) and imbalanced (B) test set. The individual cells are coloured according to the average accuracy on 200 rounds of resampling with replacement with the variance noted next to the average accuracy. The averages of accuracy per classifiers are shown in the rightmost column, whereas the bottom column shows the averages per data type.

## Core genome and WGS datasets increase the power of source attribution models

Having established the competitiveness of machine learning approaches for source attribution using MLST data, we turned our attention to whole genome datasets. Gene-by-gene approaches to cataloguing genomic variation in *Campylobacter* (46) and other species are a logical extension of seven-locus MLST in response to the increasing availability of large WGS datasets. Formalizing this

9

177  approach to derive an approximation of the core genome for *C. jejuni* allowed the implementation of

178  a cgMLST scheme containing 1,343 genes, that are present in the majority (>95%) of *C. jejuni* genomes

179  (47). This has potential to increase the power of attribution models to discriminate the source of

180  *Campylobacter* isolates based on host segregating genetic variation within the genome (37). The

181  strong performance of tree-based ensemble classifiers continued when using cgMLST data where the

182  XGBoost classifier achieved 81.3%/84.6% accuracy, the highest accuracy over all data types and

183  classifiers.

184  q

185  Next, we assessed the relative performance of machine learners when applied to k-mers produced

186  from WGS, where the average attribution performance was the highest among all datasets. The best-

187  performing algorithm was a 1-D convolutional neural net (75.0/78.3%), performing better than the

188  top-achieving classifier on MLST but worse than the best classifier on cgMLST despite WGS encoding

189  more genomic information. This may be explained by the feature selection used to limit the input

190  vector to 100,000 k-mers. Beyond comparing classifier performance on different data types, we also

191  wanted to investigate what led to the difference in performance.

192

193  The comparison of average accuracy across all data types reveals that with an increase in encoded

194  variation the average performance across all algorithms improves. This is especially apparent in MLST

195  where, although capturing the same 0.3% of the genome in all isolates, the additional variation in the

196  underlying sequences can be leveraged for better performance. When comparing the average

197  accuracy between classifiers we observed that decision-tree based ensemble learners performed well

198  across all datasets, with random forests performing best on average. The excellent performance of

199  ensemble tree learners on genomic data has been reported on genomic data (48–50) and is linked to

200  their ability to handle correlation as well as interaction of features which is an inherent feature of

10

201    genomic data (50).

202

203    Amongst simple learners the K-nearest neighbour algorithm (KNN) performed best, probably owing to

204    the hereditary nature of the phenotypic trait used as classes here. Host association is inherited both

205    genetically, in the ability to colonise different hosts, and environmentally, in the colocation of parent

206    and offspring cells. These patterns of inheritance result in more closely related sequences being more

207    likely to be associated with the same phenotype. Heritability could explain the success of the KNN

208    algorithm which is based on proximity in hyperdimensional feature space (51), which in our case is

209    genetic similarity which is a proxy for relatedness.

210

211    The deep learners generally improved in performance with higher dimensionality of the input data -

212    from MLST to WGS data. Among all deep learning architectures, the RNN and LSTM performed best,

213    which was to be expected as DNA is transcribed, and mRNA translated, sequentially 5' to 3'. Both RNNs

214    and LSTMs process input data sequentially and input weights are also adjusted sequentially in back-

215    propagation as opposed to the dense or convolutional architectures where input weights are tweaked

216    concurrently. Having investigated trends across all datasets and algorithms we focused on the best-

217    achieving classifier for a more thorough analysis of how classification performance was driven by

218    different factors within the underlying data.

219

220    **Host transition imposes a biological limit on source attribution models**

221    To better understand the limitations of attribution algorithms we investigated the factors driving

222    misclassification in the different models with different datasets. The XGBoost implementation of

223    gradient boosted decision trees, using the cgMLST dataset, was the overall best-performing classifier

224    in our analyses. Consequently, this was used to investigate attribution performance further (Figure 2).

11

225 Among all source populations the most frequent misclassification was found between sheep and

226 cattle, which is a common source of errors in source attribution (17) owing to strongly overlapping

227 gene pools stemming from frequent cross-species transmission that may reflect commonalities in

228 physiological features of the ruminant gastrointestinal tracts (52). We also looked at factors besides

229 source reservoir of the sample, as circumstances like geographical origin of the isolate (56) and the

230 season in which they were sampled (57) have been shown to influence source attribution. We

231 therefore stratified classification accuracy by continent, year, generalist index and *Campylobacter*

232 species using the full non-undersampled Test dataset (Figure 3, S1 Table).



233

234 **Figure 2:** XGBoost Classifier performance on cgMLST: A) Misclassification matrix per source. The

235    diagonal represents correct classification and off-diagonal fields are misclassifications. The

236    percentages are calculated per row. B) Misclassification matrix as depicted in a flow diagram. C)

237    Classifier performance on the unbalanced test set according to four different metrics per source

238    population. D) Radar plot showing the classifier performance on the unbalanced test by seven metrics

239    averaged over 200 rounds of resampling with replacement. The variation is depicted as a shaded

240    surface underneath the black line representing the average.

241

242    Investigating the accuracy of the XGBoost classifier per sample size revealed that the low number of

243    wild bird samples (212 samples; 84% accuracy) did not impede classification performance when

244    compared to more abundant source samples like cattle (716 samples; 84% accuracy) and sheep (584

245    samples; 57% accuracy), presumably because wild bird STs tend to be atypical compared to the other

246    reservoirs (46). To investigate how the ability to colonise multiple hosts affected performance, we

247    defined a 'generalist index' as the number of hosts in which an ST was found across all PubMLST

248    samples (S1 Table). The performance across generalist indices showed that strains restricted to fewer

249    hosts were predicted with higher accuracy. This is likely due to host switching blurring the source-

250    specific genetic signal, as previously reported (29). Consistent with this, 58% of all wild bird samples

251    belonged to STs only found in this niche, compared to 41% in environment, 9% in cattle, 3% in sheep

252    and 32% in chicken. Besides *C. jejuni*, an estimated 10% of campylobacteriosis cases are caused by

253    *Campylobacter coli* (53). Consistent with previous studies, we found improved accuracy over

254    attribution of *C. jejuni*, potentially reflecting more pronounced strain segregation by host (29), as well

255    as a higher proportion of environmental and sheep associated strains in human infection (11,54,55)

256    (Figure 3).

257

| **A** Sample count host animal isolates | **B** Accuracy of prediction on Test Set | **C** Human sample prediction | Source |
|---|---|---|---|
| 212 | 84% | 1.5% | Bird |
| 716 | 84% | 13.6% | Cattle |
| 4147 | 93% | 74.1% | Chicken |
| 140 | 77% | 0.4% | Environment |
| 584 | 57% | 10.5% | Sheep |

| | | | **Continent** |
|---|---|---|---|
| 47 | 92% | | South America |
| 465 | 82% | | North America |
| 180 | 100% | | Oceania |
| 5082 | 89% | | Europe |
| 25 | 17% | | Asia |

| | | | **Year** |
|---|---|---|---|
| 29 | 100% | | 2000 |
| 15 | 75% | | 2001 |
| 49 | 91% | | 2003 |
| 90 | 85% | | 2004 |
| 84 | 76% | | 2005 |
| 110 | 87% | | 2006 |
| 44 | 94% | | 2007 |
| 136 | 92% | | 2008 |
| 102 | 82% | | 2009 |
| 65 | 94% | | 2010 |
| 344 | 86% | | 2011 |
| 320 | 78% | | 2012 |
| 429 | 88% | | 2013 |
| 743 | 87% | | 2014 |
| 415 | 91% | | 2015 |
| 1576 | 96% | | 2016 |
| 454 | 84% | | 2017 |
| 300 | 70% | | 2018 |

**Legend**

Predicted Sources
- CHICKEN
- CATTLE
- SHEEP
- ENVIRONMENT
- BIRD

Accuracy of prediction
- CORRECT PREDICTION
- INCORRECT PREDICTION

| | | | **Generalist index** |
|---|---|---|---|
| 1593 | 96% | | 1 |
| 446 | 91% | | 2 |
| 643 | 77% | | 3 |
| 1523 | 82% | | 4 |
| 1335 | 75% | | 5 |

| | | | **Campylobacter species** |
|---|---|---|---|
| 4748 | 88% | | *C. jejuni* |
| 1043 | 90% | | *C. coli* |

14

258

259 **Figure 3:** Source attribution per source, continent, year generalist index and *Campylobacter* species.

260 A) Sample sizes across different factors in the imbalanced training set. B) Prediction accuracy on the

261 full test dataset divided by different factors. C) Source attribution stratified into varying factors

262

263 Having analysed the classification accuracy within the dataset, the machine learning method was

264 compared to previous source attribution studies (Figure 4). Attribution of cases to chicken was

265 consistent with higher estimates from previous studies, resulting in less attribution to all other

266 sources, with environment identified as the source of just 0.4% of human infections. This differences

267 in our prediction to previous studies could reflect the greater discriminatory power of cgMLST data

268 over MLST.

269

270

271

**Comparison source attribution to previous studies**

| First Author and Year of Publication | % (chicken) | (cow) | (bird) | (sheep) | (environment) | (cow+sheep) |
|---|---|---|---|---|---|---|
| Wilson 2009 | 57 | 36 | 1 | 4 | 2 | |
| Mullner 2009 | 67 | 19 | | 11 | 12 | |
| Sheppard 2009 | 78 | | 4 | | 4 | 18 |
| Kittl 2013 | 69 | 21 | | | | |
| Strachan 2009 | 43 | 35 | 6 | 15 | | |
| Gras 2012 | 66 | 21 | | 3 | 10 | |
| Mossong 2016 | 61 | | | | 5 | 33 |
| Ravel 2017 | 69 | 14 | | | 2 | |
| Rosner 2017 | 74 | 0 | | | | |
| Thepault 2018 | 56 | | | | 6 | 37 |
| Our Study | 74 | 14 | 1 | 11 | 0 | 25 |

272

15

273    **Figure 4:** Comparison of our source attribution to previously published studies

274

## The fine-grained structure of source attribution can be identified with machine learning

277    Attribution predictions are inferred from the observed frequencies of genotypes in host reservoirs

278    assessed through sampling. However, the relative source composition observed in sampling does not

279    necessarily correspond to host contributions to human infection as some strains that are found at low

280    frequency in the host could be more infectious to humans. For example, some *C. jejuni* strains increase

281    in relative frequency through different stages of the poultry slaughter and production chain because

282    they have genes that promote survival outside of the host (58). There is also evidence that there is a

283    genetic bottleneck at the point of human infection that promotes colonization by strains that have

284    specific genes conferring human niche tropism (59). Analysis of WGS or cgMLST data can potentially

285    allow for changes in relative frequency and provide finer-grained source attribution, potentially at the

286    level of the individual genome.

287

288    To identify evidence of differential host affinities, we applied treeBreaker (43) to trace the evolution

289    of a host association along the phylogeny of CC-21, the most commonly found clonal complex in

290    human infection (27). CC-21 frequently colonizes all host sources analysed in this study and is

291    therefore considered a generalist strain, potentially complicating accurate attribution. TreeBreaker

292    detected a change in host association on a branch that groups together a cattle-associated ST-21

293    subgroup with the cattle-associated lineages ST-982 and ST-806 (Figure 5A). The source composition

294    in this clade (asterisked in Figure 5A) differed from the rest of CC-21, which were predominantly

295    composed of chicken and sheep isolates. Moreover, the asterisked clade differed in its propensity for

296    transmission to humans. Overall, CC-21 was over-represented among human infections, perhaps

16

297    reflecting its generalist affinities. Yet the asterisked clade was over-represented only 1.7 to 3.6-fold,

298    compared to 5.5 to 6.2-fold for the rest of CC-21 (Figure 5B).

299

**Figure 5:** Phylogeny of clonal complex 21 of host animal associated samples (A) and bar charts showing the known source distribution and human samples (B) alongside the predicted source distribution. The

303    phylogeny is based on Neighbour joining using hamming distance of the k-mers drawn from WGS. The

304    connecting lines show the increase in frequency of the clades in human samples and the size of the

305    grey circles show the posterior probability of a change in phenotypic distribution along the branches

306    of the tree.

307

308    As the host association changed within CC-21, the ability to transmit to humans appears to have

309    changed as well. This in turn induced a change in the source composition of CC-21 sampled from

310    human infections compared to CC-21 sampled from animals. Previous studies analysing source

311    attribution based on MLST would have overlooked these shifts.

## 312   Outlook and conclusions

313    The increasing availability of large pathogen genome datasets, algorithms and resources for

314    analysing them, has created possibilities for investigating the transmission of zoonotic diseases that

315    are incompletely understood. It is clear from the data presented here that tree-based ensemble

316    methods for machine learning classification using bacterial genomic data provide considerable utility

317    for improving the accuracy host source attribution for human campylobacteriosis. Key to the

318    effectiveness of this approach is leveraging the full gradient of genomic differentiation afforded by

319    WGS or cgMLST analysis. Host associated genetic variation can be observed in both core and

320    accessory genes (60) but using these data presents practical considerations. With more

321    computational resources available, it may be possible to analyse all k-mers present in the WGS

322    samples (here 109,675,176 unique kmers) with multiple algorithms accompanied by cross-validation

323    and bootstrap replication.

324

325    Beyond simple attribution to host source, resolving the fine-grained structure of genomic signatures

19

326     of association has considerable potential to account for differences in the relative frequency of sub-

327     lineages in samples taken from reservoir hosts and human disease. This can provide important clues

328     about the propensity of strains to survive outside of the host for long enough to transmit to humans

329     as well as the capacity to colonize the human gut given the opportunity (58,59). This of course leads

330     to questions about the genomic basis of bacterial adaptation, specifically the extent to which

331     'associated' genetic elements represent adaptations and whether the same genes and alleles enable

332     colonisation of different host animals.

333

334     Improving on the approaches described here, better sampling and incremental training of the XGBoost

335     classifier has considerable potential. The classifier's low computational requirements and high

336     prediction speed make it an excellent tool for analysing large genome datasets. Furthermore, by using

337     phylogeny-aware train/test splitting for measuring performance, prediction remains accurate when

338     new genetic variants are introduced because the algorithm can be incrementally trained with new

339     data. This has considerable potential for developing automated and continuous disease surveillance

340     systems to reduce campylobacteriosis that remains one of the most common food-borne illness in the

341     world.

342

343

# Acknowledgments

## Conflicts of interest

360    DAC declares grants from GlaxoSmithKline and personal fees from Oxford University Innovation,
361    BioBeats, and Sensyne Health, in areas unrelated to this work

362

## Supporting information

364

365    S1 Table. Metadata of all *Campylobacter* isolates used in this study. Contains the accession numbers,
366    year and country of isolation, source label, generalist index, ST, CC ,prediction by our classifier,
367    *Campylobacter species* and whether the samples were used in training or testing.

368

# References

370

1.   The European Union One Health 2018 Zoonoses Report. EFSA J. 2019;17(12):e05926.

2.   Kaakoush NO, Castaño-Rodríguez N, Mitchell HM, Man SM. Global Epidemiology of Campylobacter Infection. Clin Microbiol Rev. 2015 Jul;28(3):687–720.

3.   Sheppard SK, Colles FM, McCARTHY ND, Strachan NJC, Ogden ID, Forbes KJ, et al. Niche segregation and genetic structure of Campylobacter jejuni populations from wild and agricultural host species. Mol Ecol. 2011;20(16):3484–90.

4.   Sheppard SK, Colles F, Richardson J, Cody AJ, Elson R, Lawson A, et al. Host Association of Campylobacter Genotypes Transcends Geographic Variation. Appl Environ Microbiol. 2010 Aug;76(15):5269–77.

5.   Nachamkin I, Allos BM, Ho T. Campylobacter Species and Guillain-Barré Syndrome. Clin Microbiol Rev. 1998 Jul;11(3):555–67.

6.   Nielsen LN, Sheppard SK, McCarthy ND, Maiden MCJ, Ingmer H, Krogfelt KA. MLST clustering of Campylobacter jejuni isolates from patients with gastroenteritis, reactive arthritis and Guillain–Barré syndrome. J Appl Microbiol. 2010 Feb;108(2):591–9.

7.   Altekruse SF, Stern NJ, Fields PI, Swerdlow DL. Campylobacter jejuni—An Emerging Foodborne Pathogen. Emerg Infect Dis. 1999;5(1):28–35.

8.   Gilbert MJ, Miller WG, Yee E, Zomer AL, van der Graaf-van Bloois L, Fitzgerald C, et al. Comparative Genomics of Campylobacter fetus from Reptiles and Mammals Reveals Divergent Evolution in Host-Associated Lineages. Genome Biol Evol. 2016 Jul 2;8(6):2006–19.

9.   Kirk KF, Méric G, Nielsen HL, Pascoe B, Sheppard SK, Thorlacius-Ussing O, et al. Molecular epidemiology and comparative genomics of Campylobacter concisus strains from saliva, faeces and gut mucosal biopsies in inflammatory bowel disease. Sci Rep. 2018 Jan 30;8(1):1902.

10.  Sheppard SK, Dallas JF, Wilson DJ, Strachan NJC, McCarthy ND, Jolley KA, et al. Evolution of an Agriculture-Associated Disease Causing Campylobacter coli Clade: Evidence from National Surveillance Data in Scotland. PLOS ONE. 2010 Dec 15;5(12):e15708.

11.  Ogden ID, Dallas JF, MacRae M, Rotariu O, Reay KW, Leitch M, et al. Campylobacter excreted into the environment by animal sources: prevalence, concentration shed, and host association. Foodborne Pathog Dis. 2009 Dec;6(10):1161–70.

12.  Institute of Environmental Science and Research Ltd. Notifiable and other diseases in New Zealand: Annual Report 2006. Porirua NZ Inst. 2007;

13.  Sheppard SK, Dallas JF, MacRae M, McCarthy ND, Sproston EL, Gormley FJ, et al.

Campylobacter genotypes from food animals, environmental sources and clinical disease in Scotland 2005/6. Int J Food Microbiol. 2009 Aug 31;134(1–2):96–103.

14.  Nichols GL, Richardson JF, Sheppard SK, Lane C, Sarran C. Campylobacter epidemiology: a descriptive study reviewing 1 million cases in England and Wales between 1989 and 2011. BMJ Open. 2012 Jan 1;2(4):e001179.

15.  Sears A, Baker MG, Wilson N, Marshall J, Muellner P, Campbell DM, et al. Marked Campylobacteriosis Decline after Interventions Aimed at Poultry, New Zealand. Emerg Infect Dis. 2011 Jun;17(6):1007–15.

16.  Nohra A, Grinberg A, Marshall JC, Midwinter AC, Collins-Emerson JM, French NP. Shifts in the Molecular Epidemiology of Campylobacter jejuni Infections in a Sentinel Region of New Zealand following Implementation of Food Safety Interventions by the Poultry Industry. Appl Environ Microbiol [Internet]. 2020 Feb 18 [cited 2021 Jan 6];86(5). Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7028974/

17.  Wilson DJ, Gabriel E, Leatherbarrow AJH, Cheesbrough J, Gee S, Bolton E, et al. Tracing the Source of Campylobacteriosis. PLOS Genet. 2008 Sep;4(9):e1000203.

18.  Sheppard SK, Dallas JF, Strachan NJC, MacRae M, McCarthy ND, Wilson DJ, et al. Campylobacter Genotyping to Determine the Source of Human Infection. Clin Infect Dis. 2009 Apr;48(8):1072–8.

19.  Maiden MCJ, Bygraves JA, Feil E, Morelli G, Russell JE, Urwin R, et al. Multilocus sequence typing: A portable approach to the identification of clones within populations of pathogenic microorganisms. Proc Natl Acad Sci U S A. 1998 Mar;95(6):3140–5.

20.  Dingle KE, Colles FM, Wareing DR, Ure R, Fox AJ, Bolton FE, et al. Multilocus sequence typing system for Campylobacter jejuni. J Clin Microbiol. 2001 Jan;39(1):14–23.

21.  Pritchard JK, Stephens M, Donnelly P. Inference of Population Structure Using Multilocus Genotype Data. Genetics. 2000 Jun;155(2):945–59.

22.  Mullner P, Spencer SEF, Wilson DJ, Jones G, Noble AD, Midwinter AC, et al. Assigning the source of human campylobacteriosis in New Zealand: A comparative genetic and epidemiological approach. Infect Genet Evol. 2009 Dec;9(6):1311–9.

23.  Boysen L, Rosenquist H, Larsson JT, Nielsen EM, Sørensen G, Nordentoft S, et al. Source attribution of human campylobacteriosis in Denmark. Epidemiol Infect. 2014 Aug;142(8):1599–608.

24.  Di Giannatale E, Garofolo G, Alessiani A, Di Donato G, Candeloro L, Vencia W, et al. Tracing Back Clinical Campylobacter jejuni in the Northwest of Italy and Assessing Their Potential Source. Front Microbiol [Internet]. 2016 Jun 13 [cited 2021 Feb 3];7. Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4904018/

25.  Kittl S, Heckel G, Korczak BM, Kuhnert P. Source Attribution of Human Campylobacter Isolates by MLST and Fla-Typing and Association of Genotypes with Quinolone Resistance. PLOS ONE.

23

439        2013 Nov;8(11):e81796.

440   26.   Mourkas E, Taylor AJ, Méric G, Bayliss SC, Pascoe B, Mageiros L, et al. Agricultural
441       intensification and the evolution of host specialism in the enteric pathogen Campylobacter
442       jejuni. Proc Natl Acad Sci. 2020 May 19;117(20):11018–28.

443   27.   Sheppard SK, Cheng L, Méric G, Haan CPA de, Llarena A-K, Marttinen P, et al. Cryptic ecology
444       among host generalist Campylobacter jejuni in domestic animals. Mol Ecol. 2014;23(10):2442–
445       51.

446   28.   Woodcock DJ, Krusche P, Strachan NJC, Forbes KJ, Cohan FM, Méric G, et al. Genomic plasticity
447       and rapid host switching can promote the evolution of generalism: a case study in the zoonotic
448       pathogen Campylobacter. Sci Rep. 2017 Aug;7(1):1–13.

449   29.   Dearlove BL, Cody AJ, Pascoe B, Méric G, Wilson DJ, Sheppard SK. Rapid host switching in
450       generalist Campylobacter strains erodes the signal for tracing human infections. ISME J. 2016
451       Mar;10(3):721–9.

452   30.   Sheppard SK, Guttman DS, Fitzgerald JR. Population genomics of bacterial host adaptation. Nat
453       Rev Genet. 2018 Sep;19(9):549–65.

454   31.   Thépault A, Rose V, Quesne S, Poezevara T, Béven V, Hirchaud E, et al. Ruminant and chicken:
455       important sources of campylobacteriosis in France despite a variation of source attribution in
456       2009 and 2015. Sci Rep. 2018 Jun;8(1):9305.

457   32.   Jehanne Q, Pascoe B, Bénéjat L, Ducournau A, Buissonnière A, Mourkas E, et al. Genome-Wide
458       Identification of Host-Segregating Single-Nucleotide Polymorphisms for Source Attribution of
459       Clinical Campylobacter coli Isolates. Appl Environ Microbiol [Internet]. 2020 Nov 24 [cited 2021
460       Feb 3];86(24). Available from: https://aem.asm.org/content/86/24/e01787-20

461   33.   Berthenet E, Thépault A, Chemaly M, Rivoal K, Ducournau A, Buissonnière A, et al. Source
462       attribution of Campylobacter jejuni shows variable importance of chicken and ruminants
463       reservoirs in non-invasive and invasive French clinical isolates. Sci Rep. 2019 May 30;9(1):8098.

464   34.   Zhang S, Li S, Gu W, den Bakker H, Boxrud D, Taylor A, et al. Zoonotic Source Attribution of
465       Salmonella enterica Serotype Typhimurium Using Genomic Surveillance Data, United States.
466       Emerg Infect Dis. 2019;25(1):82–91.

467   35.   Lupolova N, Dallman TJ, Holden NJ, Gally DL. Patchy promiscuity: machine learning applied to
468       predict the host specificity of Salmonella enterica and Escherichia coli. Microb Genomics
469       [Internet]. 2017 Oct [cited 2019 Sep 16];3(10). Available from:
470       https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5695212/

471   36.   Lees JA, Mai TT, Galardini M, Wheeler NE, Horsfield ST, Parkhill J, et al. Improved Prediction of
472       Bacterial Genotype-Phenotype Associations Using Interpretable Pangenome-Spanning
473       Regressions. mBio [Internet]. 2020 Aug 25 [cited 2021 Feb 3];11(4). Available from:
474       https://mbio.asm.org/content/11/4/e01344-20

475   37.   Thépault A, Méric G, Rivoal K, Pascoe B, Mageiros L, Touzain F, et al. Genome-Wide

476          Identification of Host-Segregating Epidemiological Markers for Source Attribution in
477          Campylobacter jejuni. Appl Environ Microbiol. 2017 Apr 1;83(7).

478    38.   Rizk G, Lavenier D, Chikhi R. DSK: k-mer counting with very low memory usage. Bioinformatics.
479          2013 Mar;29(5):652–3.

480    39.   Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: Proceedings of the 22Nd
481          ACM SIGKDD International Conference on Knowledge Discovery and Data Mining [Internet].
482          New York, NY, USA: ACM; 2016 [cited 2019 Sep 17]. p. 785–94. (KDD '16). Available from:
483          http://doi.acm.org/10.1145/2939672.2939785

484    40.   Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn:
485          Machine Learning in Python. J Mach Learn Res. 2011;12:2825–30.

486    41.   Kingma DP, Ba J. Adam: A Method for Stochastic Optimization. ArXiv14126980 Cs [Internet].
487          2014 Dec [cited 2019 Sep 17]; Available from: http://arxiv.org/abs/1412.6980

488    42.   Hedge J, Wilson DJ. Bacterial Phylogenetic Reconstruction from Whole Genomes Is Robust to
489          Recombination but Demographic Inference Is Not. mBio [Internet]. 2014 Dec 31 [cited 2020
490          Nov 18];5(6). Available from: https://mbio.asm.org/content/5/6/e02158-14

491    43.   Ansari MA, Didelot X. Bayesian Inference of the Evolution of a Phenotype Distribution on a
492          Phylogenetic Tree. Genetics. 2016 Sep 1;204(1):89–98.

493    44.   Argimón S, Abudahab K, Goater RJE, Fedosejev A, Bhai J, Glasner C, et al. Microreact: visualizing
494          and sharing data for genomic epidemiology and phylogeography. Microb Genomics.
495          2016;2(11):e000093.

496    45.   Cody AJ, Maiden MC, Strachan NJ, McCarthy ND. A systematic review of source attribution of
497          human campylobacteriosis using multilocus sequence typing. Eurosurveillance [Internet]. 2019
498          Oct [cited 2020 Jan 27];24(43). Available from:
499          https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6820127/

500    46.   Sheppard SK, Jolley KA, Maiden MCJ. A Gene-By-Gene Approach to Bacterial Population
501          Genomics: Whole Genome MLST of Campylobacter. Genes. 2012 Apr 12;3(2):261–77.

502    47.   Cody AJ, Bray JE, Jolley KA, McCarthy ND, Maiden MCJ. Core Genome Multilocus Sequence
503          Typing Scheme for Stable, Comparative Analyses of Campylobacter jejuni and C. coli Human
504          Disease Isolates. J Clin Microbiol. 2017 Jul;55(7):2086–97.

505    48.   Austerlitz F, David O, Schaeffer B, Bleakley K, Olteanu M, Leblois R, et al. DNA barcode analysis:
506          a comparison of phylogenetic and statistical classification methods. BMC Bioinformatics. 2009
507          Nov;10(14):S10.

508    49.   Deneke C, Rentzsch R, Renard BY. PaPrBaG: A machine learning approach for the detection of
509          novel pathogens from NGS data. Sci Rep. 2017 Jan;7:39194.

510    50.   Chen X, Ishwaran H. Random Forests for Genomic Data Analysis. Genomics. 2012
511          Jun;99(6):323–9.

512    51.    Kotsiantis SB, Zaharakis ID, Pintelas PE. Machine learning: a review of classification and
513           combining techniques. Artif Intell Rev. 2006 Nov;26(3):159–90.

514    52.    Kwan PSL, Birtles A, Bolton FJ, French NP, Robinson SE, Newbold LS, et al. Longitudinal Study of
515           the Molecular Epidemiology of Campylobacter jejuni in Cattle on Dairy Farms. Appl Environ
516           Microbiol. 2008 Jun;74(12):3626–33.

517    53.    Sheppard SK, Maiden MCJ. The Evolution of Campylobacter jejuni and Campylobacter coli. Cold
518           Spring Harb Perspect Biol [Internet]. 2015 Aug [cited 2019 Sep 3];7(8). Available from:
519           https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4526750/

520    54.    Roux F, Sproston E, Rotariu O, MacRae M, Sheppard SK, Bessell P, et al. Elucidating the
521           Aetiology of Human Campylobacter coli Infections. PLoS ONE [Internet]. 2013 May [cited 2020
522           Feb 14];8(5). Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3667194/

523    55.    Strachan NJC, Gormley FJ, Rotariu O, Ogden ID, Miller G, Dunn GM, et al. Attribution of
524           Campylobacter Infections in Northeast Scotland to Specific Sources by Use of Multilocus
525           Sequence Typing. J Infect Dis. 2009 Apr;199(8):1205–8.

526    56.    Pérez-Reche FJ, Rotariu O, Lopes BS, Forbes KJ, Strachan NJC. Mining whole genome sequence
527           data to efficiently attribute individuals to source populations. Sci Rep. 2020 Jul 22;10(1):12124.

528    57.    STRACHAN NJC, ROTARIU O, SMITH-PALMER A, COWDEN J, SHEPPARD SK, O'BRIEN SJ, et al.
529           Identifying the seasonal origins of human campylobacteriosis. Epidemiol Infect. 2013
530           Jun;141(6):1267–75.

531    58.    Yahara K, Méric G, Taylor AJ, Vries SPW de, Murray S, Pascoe B, et al. Genome-wide association
532           of functional traits linked with Campylobacter jejuni survival from farm to fork. Environ
533           Microbiol. 2017;19(1):361–80.

534    59.    Méric G, McNally A, Pessia A, Mourkas E, Pascoe B, Mageiros L, et al. Convergent Amino Acid
535           Signatures in Polyphyletic Campylobacter jejuni Subpopulations Suggest Human Niche Tropism.
536           Genome Biol Evol. 2018 Mar 1;10(3):763–74.

537    60.    Sheppard SK, Didelot X, Meric G, Torralbo A, Jolley KA, Kelly DJ, et al. Genome-wide association
538           study identifies vitamin B5 biosynthesis as a host specificity factor in Campylobacter. Proc Natl
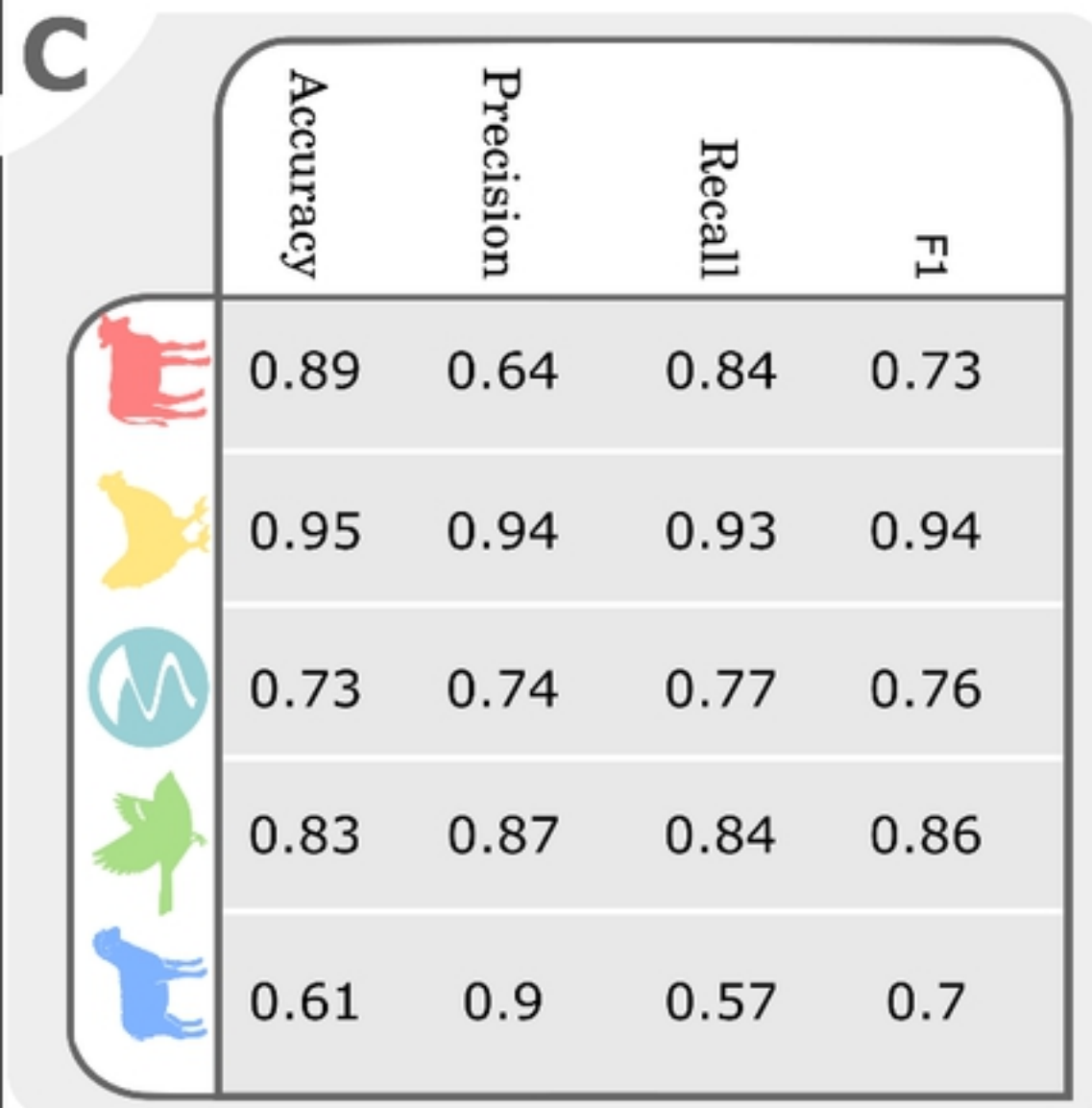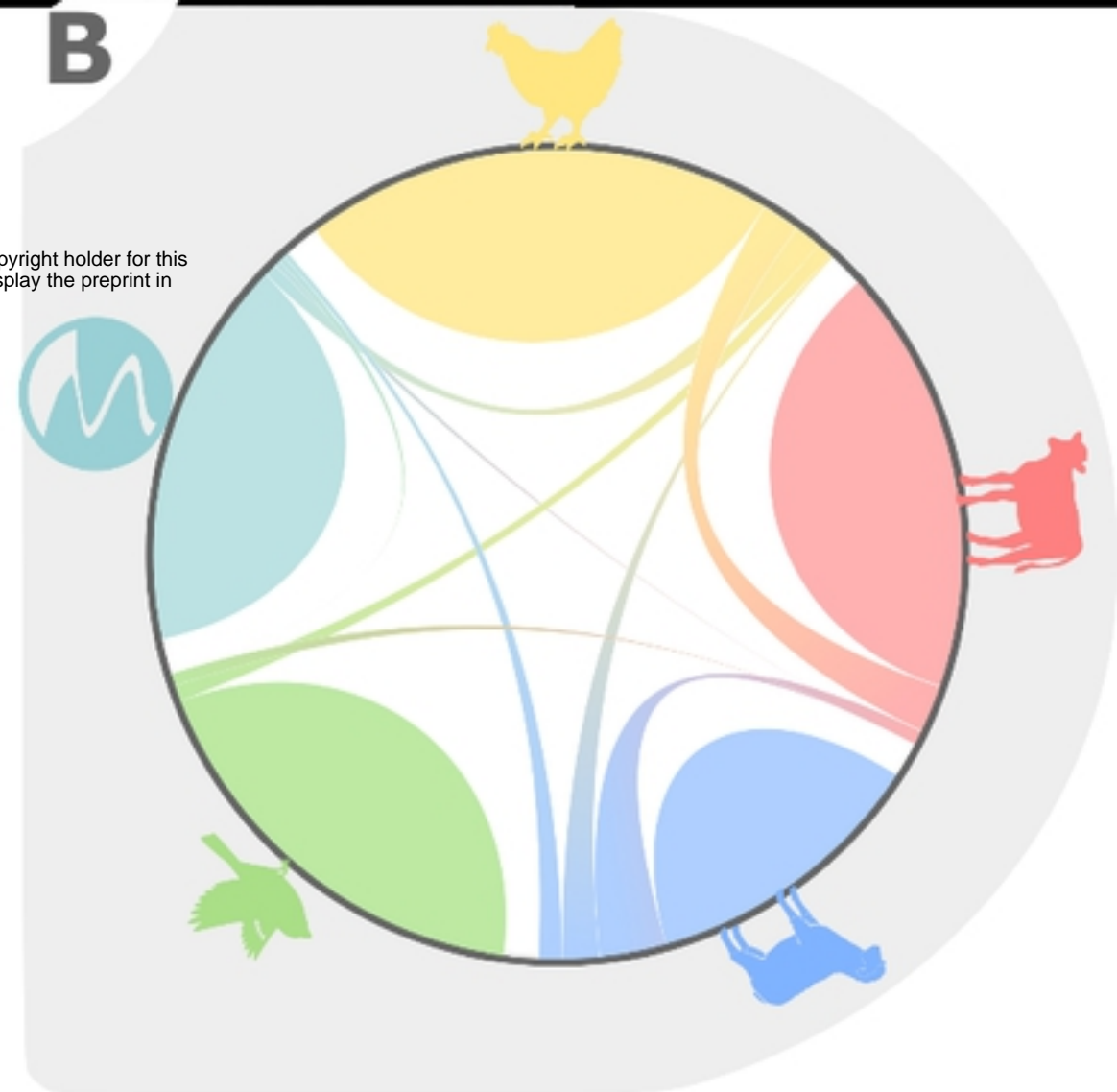539           Acad Sci. 2013 Jul;110(29):11923–7.

540

Figure 1

Figure 2

Figure 3

**Comparison source attribution to previous studies**

| % | 🐔 | 🐄 | 🐦 | 🐑 | 🧬 | 🐄+🐑 |
|---|---|---|---|---|---|---|
| Wilson 2009 | 57 | 36 | 1 | 4 | 2 | |
| Mullner 2009 | 67 | 19 | | 11 | 12 | |
| Sheppard 2009 | 78 | | 4 | | 4 | 18 |
| Kittl 2013 | 69 | 21 | | | | |
| Strachan 2009 | 43 | 35 | 6 | 15 | | |
| Gras 2012 | 66 | 21 | | 3 | 10 | |
| Mossong 2016 | 61 | | | | 5 | 33 |
| Ravel 2017 | 69 | 14 | | | 2 | |
| Rosner 2017 | 74 | 0 | | | | |
| Thepault 2018 | 56 | | | | 6 | 37 |
| Our Study | 74 | 14 | 1 | 11 | 0 | 25 |

Figure 4

Figure 5