# PLOS GENETICS

RESEARCH ARTICLE

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

**Nicolas Arning**[1]*, **Samuel K. Sheppard**[2], **Sion Bayliss**[2], **David A. Clifton**[3], **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, United Kingdom, **2** The Milner Centre of Evolution, Department of Biology & Biochemistry, University of Bath, Claverton Down, Bath, United Kingdom, **3** Department of Engineering Science, University of Oxford, Oxford, UK; Oxford-Suzhou Centre for Advanced Research, Suzhou, China

\* Nicolas.arning@bdi.ox.ac.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 the classifier we named aiSource. 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 is 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

*Campylobacter jejuni* and *Campylobacter coli* are among the most common causes of gastroenteritis globally and are responsible for approximately nine million annual cases in the European Union [1,2]. These zoonotic bacteria are a common commensal constituent of the gut microbiota of bird and animal species [3,4] but cause serious infections in humans. Symptoms include nausea, fever, abdominal pain, and severe diarrhoea, with potential for the development of debilitating, and sometimes fatal, sequelae [5,6]. Various infection sources have been identified including animal faeces, contaminated drinking water and especially raw or undercooked poultry and other meats [7]. However, effectively combating disease requires a detailed understanding of the relative contribution of different sources to human infection.

As in many other bacterial species, *Campylobacter* populations represent diverse assemblages of strains [3,8–10]. Within this structured population, some lineages are more commonly observed in particular host species [3,4,11]. Because of this host association, DNA sequence comparisons of bacteria from human gastroenteritis and potential reservoir populations have potential to reveal the infection source. This has identified contaminated poultry as a major source of human infection [12,13]. Based on the body of evidence including DNA sequence analysis [14], targeted interventions have been implemented, including improved biosecurity measures on poultry farms, which have halved recorded campylobacteriosis cases in New Zealand [15,16].

Extending the principal of linking source-sink populations using genotype data, methods have been developed to attribute *C. jejuni* to the likely source based on bacterial gene frequencies in potential reservoir populations [17,18]. Among the most common genotyping approaches for *C. jejuni* has been multi-locus sequence typing (MLST) that catalogues DNA sequence variation across seven housekeeping genes that are common to all strains [19,20]. Isolates with identical alleles at all loci are assigned to the same sequence type (ST) and those with identical sequences at most or all loci are grouped within the same clonal complex (CC). Using these data, and allele frequencies, it has been possible to probabilistically assign clinical isolates (STs and CCs) to host source using source attribution models such as the asymmetric island model implemented in *iSource* [17] and the Bayesian population assignment model STRUCTURE [18,21]. Both methods have been instructive in estimating the relative contribution of a range of domestic and wild animal hosts to human infection, with poultry often identified as the principal source of human campylobacteriosis across different regions and countries [17,18,22–25].

There are two main limitations when using genotype data to for bacterial source attribution. The first is that the ability to attribute is only as good as the degree of genotype segregation. For example, in *C. jejuni* there are host restricted genotypes [3,26] that can be readily attributed to a given host source when observed in human infections, as well as ecological generalists [27,28] that have relatively recently transitioned between hosts and cannot therefore be attributed with confidence [29]. While host switching potentially imposes a biological

constraint on quantitative attribution models, the second limitation is far more tractable. Specifically, most current source attribution methods are subject to limitations imposed by the underlying data. Reflecting the technology of the time, MLST-based source attribution is based only on a small fraction of the genome (approximately 0.2% for *C. jejuni* [25]) and there is considerable potential for better strain differentiation using current techniques.

The increasing availability of large whole genome sequence (WGS) datasets has greatly enhanced analyses of bacterial population structure and diversity [30]. However, exploiting the full information can be challenging due to variable gene content and the complexity of interpreting the short reads produced by next generation sequencing. Notwithstanding this, some studies have attempted to overcome the limited discriminatory power of MLST in attribution studies by screening WGS data to identify elements (SNPs and genes) that segregate by host [31–34]. Using these host segregating markers as input data has improved the resolution of existing attribution models, including STRUCTURE, and provided information about potential infection reservoirs and the UK and France. However, using bespoke marker selection approaches with software designed for MLST data does not maximize the potential of WGS data for source attribution.

Here, we present a machine learning approach using WGS data to predict the source of human *C. jejuni* infection. This has two principal advantages over existing techniques. First, building on WGS-based machine learning source attribution approaches applied to *Salmonella enterica* and *Escherichia coli* [35,36], we take an agnostic approach to identify which machine learning tool performs best from a broad range of available algorithms. Second, we use a WGS input capture approach using data types conveniently available in public databases such as PubMLST [37]) allowing the analysis of existing MLST, core-genome MLST and WGS datasets and the reuse of data for continuous updatable monitoring in a generalizable framework. Thus, we aimed to overcome limitations of the currently available methods and use the output to investigate the infective potential of *C. jejuni* strains.

## Methods

### Dataset acquisition

A total of 5,799 *C. jejuni* and *C. coli* genomes isolated from various sources and host species were available on the public database for molecular typing and microbial genome diversity: PubMLST (https://pubmlst.org/) with the following source distribution: (chicken: 4147, cattle: 716, sheep: 584, bird: 212, environment: 140). WGS data corresponded to MLST ST and CC designations as well as core genome (cg) MLST classes. The dataset was divided into training (75%) and testing (25%) sets, but we diverged from the more common independent random drawing of individual samples. Instead we used phylogeny-aware sorting, wherein all members of one ST were sorted entirely into either training or testing sets (S1 Table). The ST based sorting accounts for the phylogenetic non-independence of samples [38]. To allow for sufficient sample sizes per reservoir population (hereafter "class"), only the five most prevalent classes for MLST and cgMLST were used (chicken, cattle, sheep, wild bird and environment). For farm animals the classes "chicken" and "chicken offal or meat" were combined to "chicken" (likewise for sheep and cattle), whilst "environment", "sand" and "river water" were combined into "environment", consistent with previous studies [18,39].

### Feature engineering

The allelic profiles of MLST and cgMLST were used directly. MLST samples that had missing alleles on any loci and cgMLST samples with more than 10% missing loci were discarded, with the missing alleles in cgMLST encoded as -1. To potentially exploit the gradient of separation

encoded in the sequences underlying the MLST allelic profiles, we downloaded the underlying allele sequences for every loci of the MLST scheme and encoded the nucleotides as dummy variables and k-mers (k = 21) using DSK [40]. DSK was also used for encoding the WGS as k-mers, as they have previously been successfully used on *C. jejuni* WGS analysis, namely for determining the genetic basis of C. jejuni host affinity [41] and survival [42]. Using k = 21 led to a prohibitively large input vector due to the number of unique k-mers found in all genomes (109,675,176). We reduced the number of k-mers by applying a variance threshold where k-mers which were present or absent in more than 99% of the samples were discarded, reducing the numbers of unique k-mers to 7,285,583. Furthermore, we performed feature selection by testing the dependence of the source labels on every individual k-mer using the Chi-Square statistic. To avoid data-leakage we only performed the feature selection using the training data and labels to select the 100,000 k-mers with the highest score.

## Algorithm training

All machine learning and deep learning was performed in Python (for a list of all algorithms see Fig 1). The xgboost library [43] was used for the gradient boosting classifiers with all other machine learners implemented in scikit-learn [44]. The hyper-parameters for each classifier were chosen using Cartesian grid search on five-fold cross-validation of the training set. The Keras library (https://keras.io/) was used to construct deep learning algorithms aimed at supplying a wide range of commonly used architectures. We found this to work best, empirically,



**Fig 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 one standard error 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.

given that there is no principled means of architecture selection for such models. Specifically: (i) A recurrent neural network consisting (RNN) of a layer with 64 gated recurrent units, a 50% dropout layer and Rectified Linear Unit (ReLU) activation layer; (ii) A 1-dimensional convolutional network with two convolutional layers of kernel size 3 and 5 respectively and 30 filters, both followed by 50% dropout layers and a ReLU layer; (iii) A Long short-term memory network (LSTM) consisting of one LSTM layer with 64 units and a 50% dropout layer; (iv) A Shallow dense network with one dense layer with 64 units followed by a 50% dropout layer and a ReLU activation layer; (v) A Deep dense network with 6 dense layers starting with 128 units and halving units with each successive layer. All individual dense layers are followed by a 50% dropout layer and a ReLU layer.

To all deep learning architectures, we added an output layer comprising a dense layer with soft-max activation with one unit for every class. We encoded the labels as dummy variables and used categorical cross-entropy as a loss function together with the Adam optimiser [45]. Cyclical learning rates were used with a maximum learning rate of 0.1 and a minimum learning rate of 0.0001 to overcome local minima. The accuracy on the test set was measured at every epoch and the overall best performing weights were stored as a checkpoint. The data was deployed in batches of 128 samples with every batch randomly undersampled so that each class was represented in equal proportions. The training was run for 500 generations with early stopping after 50 generations.

## Algorithm testing

Both machine learning and deep learning were tested on the same 25% test set. The original data were skewed in source composition by ratios which did not necessarily reflect source origin of infection. We therefore used two methods to rebalance the classes in testing. The first test set featured an even distribution of classes, whereas the second undersampled the over-abundant chicken-origin genomes to emulate relative contribution to human disease. We used the ratios predicted by Wilson et al. (12), where *Campylobacter* genomes from chickens were 1.61 times more common than those from cattle. In both methods, rebalancing the classes was achieved by undersampling, which we repeated 200 times with replacement and averaged the accuracy over all iterations whilst also recording one standard error. As our balanced test set is limited by the number of available samples from the minority source (35 environment samples), the repeated undersampling allows us to use all available samples of the residual classes in testing. For performance metrics we registered accuracy, precision (positive predictive value), recall (sensitivity), F1, negative predictive value, specificity and speed. Speed was measured relative to other classifiers where a scale was defined with 0 being the slowest classifier and 1 being the quickest and all intermediate values being normalised within these confines. For comparison to previous methods, iSource was applied to the test dataset [17]. Having established that XGBoost on cgMLST was the best performing source attribution method, we retrained the classifier with both training and testing data and applied it to all 15,988 human cgMLST samples available on the PubMLST database. The prediction took 892 milliseconds on a Dell OptiPlex 7060 desktop using ten threads on an Intel Core i7-8700 CPU and 16 GB RAM. Our algorithm named aiSource can be found and applied from: https://github.com/narning1992/aiSource

## Phylogenetic analysis

We defined the generalist index as the number of sources the ST was found in across all isolates in the dataset, which included additional samples for which only MLST data was available (S1 Table). We built a phylogeny of CC21 genomes from both source-associated and human

isolates using Neighbour Joining, based on pairwise hamming distances of k-mer presence/ absence in the WGS dataset, as described by Hedge and Wilson [46]. We used TreeBreaker [47] to infer the evolution of phenotypes across the phylogenetic tree of ST-21 and the most closely related sequence types. The known labels of the source-associated samples were used as phenotypic information for input into TreeBreaker together with the phylogeny of CC21. TreeBreaker was run for 5,500,000 iterations with 500,000 iterations as burn-in and 1000 iterations between sampling. The phylogenetic trees were visualised with Microreact [48] and arranged alongside the results of TreeBreaker in Inkscape.

## Results and discussion

### Machine learning outperforms popular attribution models for MLST data

In order to anchor our source attribution performance to previous efforts, we compared results using the machine learning classifiers to source probabilities estimated using the asymmetric island model implemented in iSource, which is based on MLST and the most commonly used source attribution method to date [49]. The best performing machine learner on the MLST allelic profile was a random forest (61.9%/68.5% balanced/unbalanced) which performed slightly better than iSource (61%/64%) (Fig 1). Since loci within allelic profiles are deemed either to match or not, and underlying nucleotides sequences are ignored, we investigated whether exploiting the gradient of nucleotide differentiation would lead to better attribution. We used dummy variables and generated k-mers from the sequences underlying the MLST allele labels. The additional feature encodings boosted the top achieving accuracies on MLST to 67.9%/70.7% from dummy variables and 63%/67.5% from k-mers, showing the value of the additional nucleotide-level information.

### 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* [50] and other species are a logical extension of seven-locus MLST in response to the increasing availability of large WGS datasets. Formalizing this approach to derive an approximation of the core genome for *C. jejuni* allowed the implementation of a cgMLST scheme containing 1,343 genes, that are present in the majority (>95%) of *C. jejuni* genomes [51]. This has potential to increase the power of attribution models to discriminate the source of *Campylobacter* isolates based on host segregating genetic variation within the genome [39]. The strong performance of tree-based ensemble classifiers continued when using cgMLST data where the XGBoost classifier achieved 81.3 ±2%/84.6±0% accuracy, the highest accuracy over all data types and classifiers.

Next, we assessed the relative performance of machine learners when applied to k-mers produced from WGS, where the average attribution performance was the highest among all datasets. The best-performing algorithm was a 1-D convolutional neural net (75.0/78.3%), performing better than the top-achieving classifier on MLST but worse than the best classifier on cgMLST despite WGS encoding more genomic information. This may be explained by the feature selection used to limit the input vector to 100,000 k-mers. Beyond comparing classifier performance on different data types, we also wanted to investigate what led to the difference in performance.

The comparison of average accuracy across all data types reveals that with an increase in genomic content being encoded the average performance across all algorithms improves. This

is especially apparent in MLST where, although capturing the same 0.3% of the genome in all isolates, the additional variation in the underlying sequences can be leveraged for better performance. When comparing the average accuracy between classifiers we observed that decision-tree based ensemble learners performed well across all datasets, with random forests performing best on average. The excellent performance of ensemble tree learners on genomic data has been reported on genomic data [52–54] and is linked to their ability to handle correlation as well as interaction of features which is an inherent feature of genomic data [54].
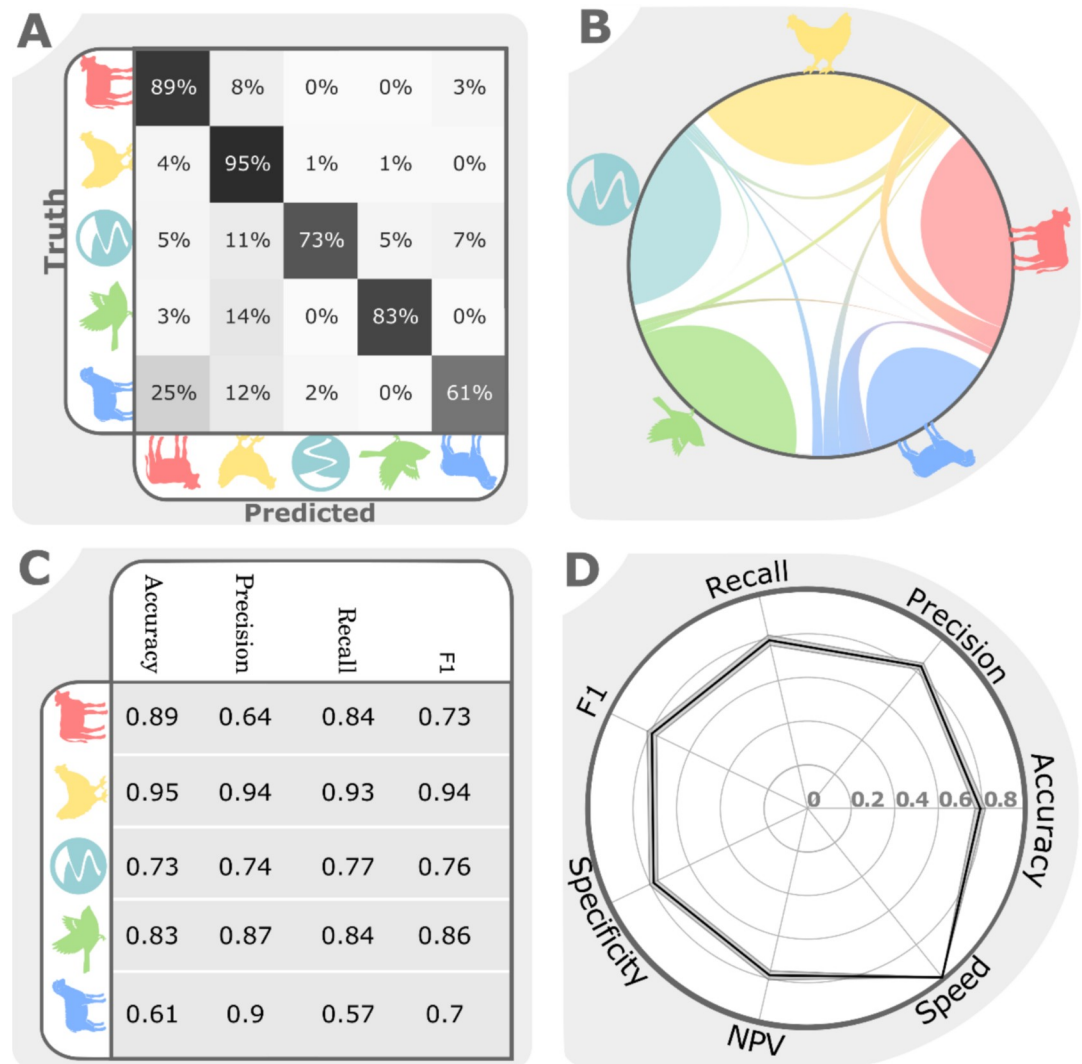
Amongst simple learners the K-nearest neighbour algorithm (KNN) performed best, probably owing to the hereditary nature of the phenotypic trait used as classes here. Host association is inherited both genetically, in the ability to colonise different hosts, and environmentally, in the colocation of parent and offspring cells. These patterns of inheritance result in more closely related sequences being more likely to be associated with the same phenotype. Heritability could explain the success of the KNN algorithm which is based on proximity in hyperdimensional feature space [55], which in our case is genetic similarity which is a proxy for relatedness.

The deep learners generally improved in performance with higher dimensionality of the input data—from MLST to WGS data. Among all deep learning architectures, the RNN and LSTM performed best, which was to be expected as DNA is transcribed, and mRNA translated, sequentially 5′ to 3′. Both RNNs and LSTMs process input data sequentially and input weights are also adjusted sequentially in back-propagation as opposed to the dense or convolutional architectures where input weights are tweaked concurrently. Having investigated trends across all datasets and algorithms we focused on the best-achieving classifier for a more thorough analysis of how classification performance was driven by different factors within the underlying data.

## Host transition imposes a biological limit on source attribution models

To better understand the limitations of attribution algorithms we investigated the factors driving misclassification in the different models with different datasets. The XGBoost implementation of gradient boosted decision trees, using the cgMLST dataset, was the overall best-performing classifier in our analyses. We coined the resulting algorithm aiSource to highlight our pedigree to the previously most commonly used iSource[17] and investigated attribution performance further (Fig 2). Among all source populations the most frequent misclassification was found between sheep and cattle, which is a common source of errors in source attribution [17] owing to strongly overlapping gene pools stemming from frequent cross-species transmission that may reflect commonalities in physiological features of the ruminant gastrointestinal tracts [56]. We also looked at factors besides source reservoir of the sample, as circumstances like geographical origin of the isolate (56) and the season in which they were sampled (57) have been shown to influence source attribution. We therefore stratified classification accuracy by continent, year, generalist index and *Campylobacter* species using the full non-under-sampled Test dataset (Fig 3 and S1 Table).

Investigating the accuracy of aiSource per sample size revealed that the low number of wild bird samples (212 samples; 84% accuracy) did not impede classification performance when compared to more abundant source samples like cattle (716 samples; 84% accuracy) and sheep (584 samples; 57% accuracy), presumably because wild bird STs tend to be atypical compared to the other reservoirs [34,50]. Investigating other stratifications reveals that sparsely sampled strata seem to be outliers in classification performance. For example, the relatively few (25) Asian examples seen in training lead to only 17% accuracy on the 12 Asian examples in testing. The 8 samples from 2000 seem to exhibit perfect attribution accuracy whilst showing comparably much higher allocation of cases to cattle (Fig 3). Generally, performance within strata showing few samples should be considered with considerable scepticism. To investigate how

**Fig 2. aiSource (based on XGBoost) performance on cgMLST.** A) Misclassification matrix per source. The diagonal represents correct classification and off-diagonal fields are misclassifications. The percentages are calculated per row. B) Misclassification matrix as depicted in a flow diagram. C) Classifier performance on the unbalanced test set according to four different metrics per source population. D) Radar plot showing the classifier performance on the unbalanced test by seven metrics averaged over 200 rounds of resampling with replacement. The variation is depicted as a shaded surface underneath the black line representing the average.

https://doi.org/10.1371/journal.pgen.1009436.g002

the ability to colonise multiple hosts affected performance, we defined a 'generalist index' as the number of hosts in which an ST was found across all PubMLST samples (S1 Table). The performance across generalist indices showed that strains restricted to fewer hosts were predicted with higher accuracy. This is likely due to host switching blurring the source-specific genetic signal, as previously reported [29]. Consistent with this, 58% of all wild bird samples belonged to STs only found in this niche, compared to 41% in environment, 9% in cattle, 3% in sheep and 32% in chicken.

Having analysed the classification accuracy within the dataset, the aiSource was compared to previous source attribution studies (Fig 4). Attribution of cases to chicken was consistent with higher estimates from previous studies, resulting in less attribution to all other sources, with environment identified as the source of just 0.05% of human infections. These differences

**Fig 3. Source attribution per source, continent, year generalist index and *Campylobacter* species.** A) Sample sizes across different factors in the imbalanced training set. B) Prediction accuracy on the full test dataset divided by different factors. C) Source

attribution of the human samples, as predicted by the XGBoost model trained on the full source associated cgMLST dataset stratified into varying factors.

https://doi.org/10.1371/journal.pgen.1009436.g003

in our prediction to previous studies could reflect the greater discriminatory power of cgMLST data over MLST.

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

Attribution predictions are inferred from the observed frequencies of genotypes in host reservoirs assessed through sampling. However, the relative source composition observed in

| Comparison source attribution to previous studies | | | | | | |
|---|---|---|---|---|---|---|
| % (First Author and Year of Publication) | 🐔 | 🐄 | 🐦 | 🐑 | 〰 | 🐄+🐑 |
| 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 |
| Boysen 2013 | 67 | 17 | | | | |
| Our Study | 75 | 15 | 1 | 9 | 0 | 24 |

**Fig 4. Comparison of source attribution using aiSource to previously published studies.**

https://doi.org/10.1371/journal.pgen.1009436.g004

sampling does not necessarily correspond to host contributions to human infection as some strains that are found at low frequency in the host could be more infectious to humans. For example, some *C. jejuni* strains increase in relative frequency through different stages of the poultry slaughter and production chain because they have genes that promote survival outside of the host [42]. There is also evidence that there is a genetic bottleneck at the point of human infection that promotes colonization by strains that have specific genes conferring human niche tropism [57]. Analysis of WGS or cgMLST data can potentially allow for changes in relative frequency and provide finer-grained source attribution, potentially at the level of the individual genome.

To identify evidence of differential host affinities, we applied treeBreaker [47] to trace the evolution of a host association along the phylogeny of CC-21, the most commonly found clonal complex in human infection (27). CC-21 frequently colonizes all host sources analysed in this study and is therefore considered a generalist strain, potentially complicating accurate attribution. TreeBreaker detected a change in host association on a branch that groups together a cattle-associated ST-21 subgroup with the cattle-associated lineages ST-982 and ST-806 (Fig 5A). The source composition in this clade (asterisked in Fig 5A) differed from the rest of CC-21, which were predominantly composed of chicken and sheep isolates. Moreover, the asterisked clade differed in its propensity for transmission to humans. Overall, CC-21 was over-represented among human infections, perhaps reflecting its generalist affinities. Yet the asterisked clade was over-represented only 1.7 to 3.6-fold, compared to 5.5 to 6.2-fold for the rest of CC-21 (Fig 5B).

As the host association changed within CC-21, the ability to transmit to humans appears to have changed as well. This in turn induced a change in the source composition of CC-21 sampled from human infections compared to CC-21 sampled from animals. Previous studies analysing source attribution based on MLST would have overlooked these shifts.

## Outlook and conclusions

The increasing availability of large pathogen genome datasets, algorithms and resources for analysing them, has created possibilities for investigating the transmission of zoonotic diseases that are incompletely understood. It is clear from the data presented here that tree-based ensemble methods for machine learning classification using bacterial genomic data provide considerable utility for improving the accuracy host source attribution for human campylobacteriosis. Key to the effectiveness of this approach is leveraging the full gradient of genomic differentiation afforded by WGS or cgMLST analysis. Host associated genetic variation can be observed in both core and accessory genes [41] but using these data presents practical considerations. With more computational resources available, it may be possible to analyse all k-mers present in the WGS samples (here 109,675,176 unique kmers) with multiple algorithms accompanied by cross-validation and bootstrap replication.

Beyond simple attribution to host source, resolving the fine-grained structure of genomic signatures of association has considerable potential to account for differences in the relative frequency of sub-lineages in samples taken from reservoir hosts and human disease. This can provide important clues about the propensity of strains to survive outside of the host for long enough to transmit to humans as well as the capacity to colonize the human gut given the opportunity [42,57]. This of course leads to questions about the genomic basis of bacterial adaptation, specifically the extent to which 'associated' genetic elements represent adaptations and whether the same genes and alleles enable colonisation of different host animals.

Improving on the approaches described here, better sampling and incremental training of aiSource, which is available under https://github.com/narning1992/aiSource, has considerable
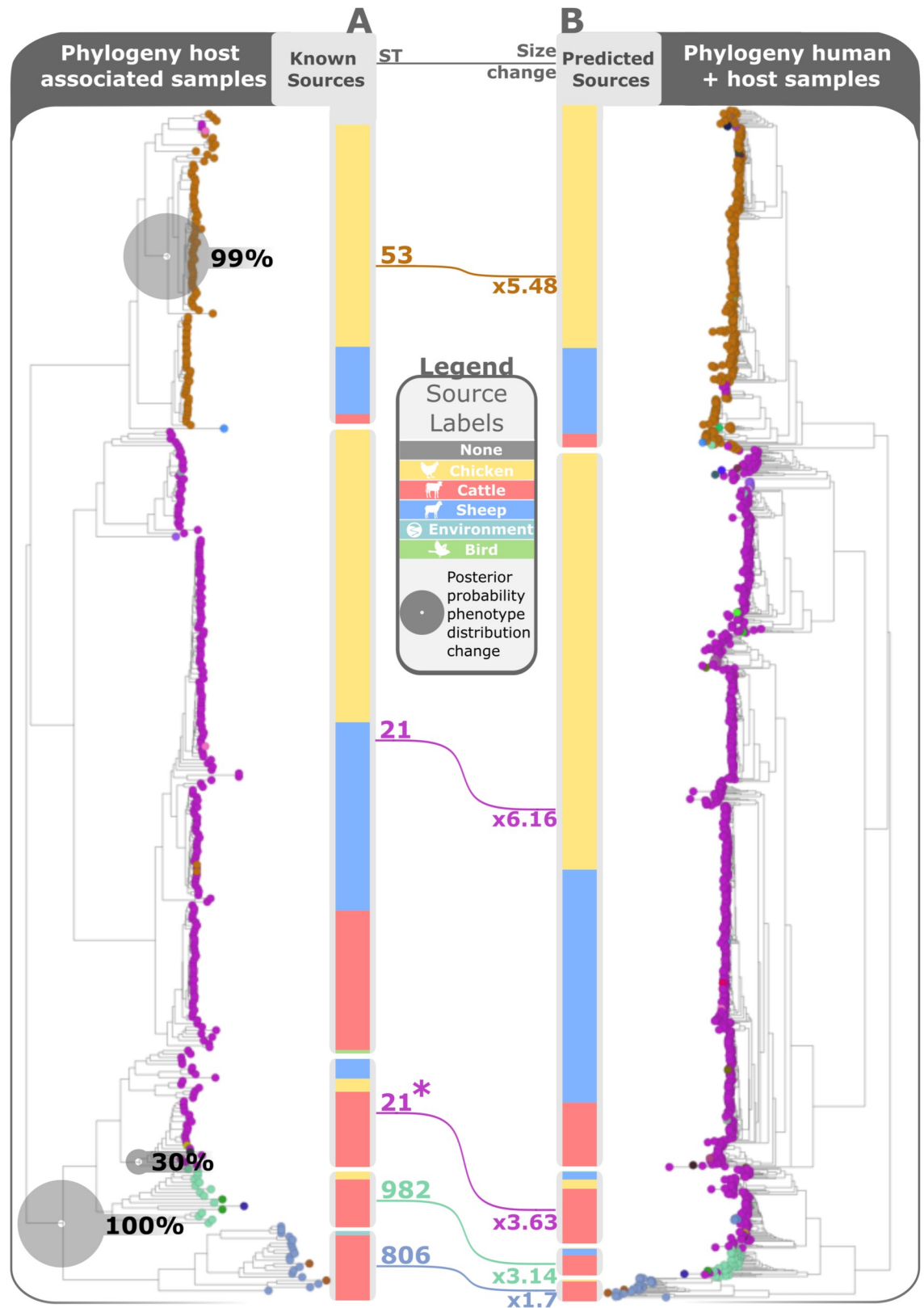
**Fig 5.** Phylogeny of clonal complex 21 of host animal associated samples (A) and bar charts showing the known source distribution and human samples (B) alongside the source distribution predicted by aiSource. The phylogeny is based on Neighbour joining using

hamming distance of the k-mers drawn from WGS. The connecting lines show the increase in frequency of the clades in human samples and the size of the grey circles show the posterior probability of a change in phenotypic distribution along the branches of the tree.

https://doi.org/10.1371/journal.pgen.1009436.g005

potential. The low computational requirements of aiSource and its high prediction speed make it an excellent tool for analysing large genome datasets. Furthermore, by using phylogeny-aware train/test splitting for measuring performance, prediction remains accurate when new genetic variants are introduced because the algorithm can be incrementally trained with new data. This has considerable potential for developing automated and continuous disease surveillance systems to reduce campylobacteriosis that remains one of the most common food-borne illness in the world.

## Supporting information

**S1 Table. Table containing all samples used in this study and their corresponding PubMLST accession IDs, sequence types, clonal complexes, source labels, predicted labels, generalist index, country of isolation, year of sampling, Campylobacter species and whether they have been used in either training or testing the machine learner.**
(TSV)

## Acknowledgments

## Author Contributions

**Conceptualization:** Daniel J. Wilson.

**Data curation:** Nicolas Arning, Sion Bayliss.

**Formal analysis:** Nicolas Arning.

**Investigation:** Nicolas Arning.

**Methodology:** Nicolas Arning, Daniel J. Wilson.

**Software:** Nicolas Arning.

**Supervision:** David A. Clifton, Daniel J. Wilson.

**Visualization:** Nicolas Arning.

**Writing – original draft:** Nicolas Arning.

**Writing – review & editing:** Samuel K. Sheppard, David A. Clifton, Daniel J. Wilson.

## References

1.   The European Union One Health 2018 Zoonoses Report. EFSA Journal. 2019; 17(12):e05926. https://doi.org/10.2903/j.efsa.2019.5926 PMID: 32626211

2. Kaakoush NO, Castaño-Rodríguez N, Mitchell HM, Man SM. Global Epidemiology of Campylobacter Infection. Clinical Microbiology Reviews. 2015 Jul; 28(3):687–720. https://doi.org/10.1128/CMR.00006-15 PMID: 26062576

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. Molecular Ecology. 2011; 20(16):3484–90. https://doi.org/10.1111/j.1365-294X.2011.05179.x PMID: 21762392

4. Sheppard SK, Colles F, Richardson J, Cody AJ, Elson R, Lawson A, et al. Host Association of Campylobacter Genotypes Transcends Geographic Variation. Applied and Environmental Microbiology. 2010 Aug; 76(15):5269–77. https://doi.org/10.1128/AEM.00124-10 PMID: 20525862

5. Nachamkin I, Allos BM, Ho T. Campylobacter Species and Guillain-Barré Syndrome. Clinical Microbiology Reviews. 1998 Jul; 11(3):555–67. https://doi.org/10.1128/CMR.11.3.555 PMID: 9665983

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. https://doi.org/10.1111/j.1365-2672.2009.04444.x PMID: 19702866

7. Altekruse SF, Stern NJ, Fields PI, Swerdlow DL. Campylobacter jejuni—An Emerging Foodborne Pathogen. Emerging Infectious Diseases. 1999; 5(1):28–35. https://doi.org/10.3201/eid0501.990104 PMID: 10081669

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. https://doi.org/10.1093/gbe/evw146 PMID: 27333878

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. Scientific Reports. 2018 Jan 30; 8(1):1902. https://doi.org/10.1038/s41598-018-20135-4 PMID: 29382867

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. https://doi.org/10.1371/journal.pone.0015708 PMID: 21179537

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): The Institute. 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. https://doi.org/10.1016/j.ijfoodmicro.2009.02.010 PMID: 19269051

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. Emerging Infectious Diseases. 2011 Jun; 17(6):1007–15. https://doi.org/10.3201/eid/1706.101272 PMID: 21749761

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 Genetics. 2008 Sep; 4(9):e1000203. https://doi.org/10.1371/journal.pgen.1000203 PMID: 18818764

18. Sheppard SK, Dallas JF, Strachan NJC, MacRae M, McCarthy ND, Wilson DJ, et al. Campylobacter Genotyping to Determine the Source of Human Infection. Clinical Infectious Diseases. 2009 Apr; 48(8):1072–8. https://doi.org/10.1086/597402 PMID: 19275496

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.

Proceedings of the National Academy of Sciences of the United States of America. 1998 Mar; 95 (6):3140–5. https://doi.org/10.1073/pnas.95.6.3140 PMID: 9501229

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. https://doi.org/10.1128/JCM.39.1.14-23.2001 PMID: 11136741

21. Pritchard JK, Stephens M, Donnelly P. Inference of Population Structure Using Multilocus Genotype Data. Genetics. 2000 Jun; 155(2):945–59. https://doi.org/10.1093/genetics/155.2.945 PMID: 10835412

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. Infection, Genetics and Evolution. 2009 Dec; 9(6):1311–9. https://doi.org/10.1016/j.meegid.2009.09.003 PMID: 19778636

23. Boysen L, Rosenquist H, Larsson JT, Nielsen EM, Sørensen G, Nordentoft S, et al. Source attribution of human campylobacteriosis in Denmark. Epidemiology & Infection. 2014 Aug; 142(8):1599–608. https://doi.org/10.1017/S0950268813002719 PMID: 24168860

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/ https://doi.org/10.3389/fmicb.2016.00887 PMID: 27379033

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. 2013 Nov; 8(11):e81796. https://doi.org/10.1371/journal.pone.0081796 PMID: 24244747

26. Mourkas E, Taylor AJ, Méric G, Bayliss SC, Pascoe B, Mageiros L, et al. Agricultural intensification and the evolution of host specialism in the enteric pathogen Campylobacter jejuni. PNAS. 2020 May 19; 117 (20):11018–28. https://doi.org/10.1073/pnas.1917168117 PMID: 32366649

27. Sheppard SK, Cheng L, Méric G, Haan CPA de, Llarena A-K, Marttinen P, et al. Cryptic ecology among host generalist Campylobacter jejuni in domestic animals. Molecular Ecology. 2014; 23(10):2442–51. https://doi.org/10.1111/mec.12742 PMID: 24689900

28. Woodcock DJ, Krusche P, Strachan NJC, Forbes KJ, Cohan FM, Méric G, et al. Genomic plasticity and rapid host switching can promote the evolution of generalism: a case study in the zoonotic pathogen Campylobacter. Scientific Reports. 2017 Aug; 7(1):1–13. https://doi.org/10.1038/s41598-016-0028-x PMID: 28127051

29. Dearlove BL, Cody AJ, Pascoe B, Méric G, Wilson DJ, Sheppard SK. Rapid host switching in generalist Campylobacter strains erodes the signal for tracing human infections. The ISME Journal. 2016 Mar; 10 (3):721–9. https://doi.org/10.1038/ismej.2015.149 PMID: 26305157

30. Sheppard SK, Guttman DS, Fitzgerald JR. Population genomics of bacterial host adaptation. Nature Reviews Genetics. 2018 Sep; 19(9):549–65. https://doi.org/10.1038/s41576-018-0032-z PMID: 29973680

31. Thépault A, Rose V, Quesne S, Poezevara T, Béven V, Hirchaud E, et al. Ruminant and chicken: important sources of campylobacteriosis in France despite a variation of source attribution in 2009 and 2015. Scientific Reports. 2018 Jun; 8(1):9305. https://doi.org/10.1038/s41598-018-27558-z PMID: 29915208

32. Jehanne Q, Pascoe B, Bénéjat L, Ducournau A, Buissonnière A, Mourkas E, et al. Genome-Wide Identification of Host-Segregating Single-Nucleotide Polymorphisms for Source Attribution of Clinical Campylobacter coli Isolates. Appl Environ Microbiol [Internet]. 2020 Nov 24 [cited 2021 Feb 3]; 86(24). Available from: https://aem.asm.org/content/86/24/e01787-20 https://doi.org/10.1128/AEM.01787-20 PMID: 33036986

33. Berthenet E, Thépault A, Chemaly M, Rivoal K, Ducournau A, Buissonnière A, et al. Source attribution of Campylobacter jejuni shows variable importance of chicken and ruminants reservoirs in non-invasive and invasive French clinical isolates. Scientific Reports. 2019 May 30; 9(1):8098. https://doi.org/10.1038/s41598-019-44454-2 PMID: 31147581

34. Weis AM, Storey DB, Taff CC, Townsend AK, Huang BC, Kong NT, et al. Genomic Comparison of Campylobacter spp. and Their Potential for Zoonotic Transmission between Birds, Primates, and Livestock. Appl Environ Microbiol. 2016 Dec 15; 82(24):7165–75. https://doi.org/10.1128/AEM.01746-16 PMID: 27736787

35. Zhang S, Li S, Gu W, den Bakker H, Boxrud D, Taylor A, et al. Zoonotic Source Attribution of Salmonella enterica Serotype Typhimurium Using Genomic Surveillance Data, United States. Emerging Infectious Diseases. 2019; 25(1):82–91. https://doi.org/10.3201/eid2501.180835 PMID: 30561314

36. Lupolova N, Dallman TJ, Holden NJ, Gally DL. Patchy promiscuity: machine learning applied to predict the host specificity of Salmonella enterica and Escherichia coli. Microbial Genomics [Internet]. 2017 Oct [cited 2019 Sep 16]; 3(10). Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5695212/ https://doi.org/10.1099/mgen.0.000135 PMID: 29177093

**37.** Jolley KA, Bray JE, Maiden MCJ. Open-access bacterial population genomics: BIGSdb software, the PubMLST.org website and their applications. Wellcome Open Res. 2018 Sep 24; 3:124. https://doi.org/10.12688/wellcomeopenres.14826.1 PMID: 30345391

**38.** Lees JA, Mai TT, Galardini M, Wheeler NE, Horsfield ST, Parkhill J, et al. Improved Prediction of Bacterial Genotype-Phenotype Associations Using Interpretable Pangenome-Spanning Regressions. mBio [Internet]. 2020 Aug 25 [cited 2021 Feb 3]; 11(4). Available from: https://mbio.asm.org/content/11/4/e01344-20 https://doi.org/10.1128/mBio.01344-20 PMID: 32636251

**39.** Thépault A, Méric G, Rivoal K, Pascoe B, Mageiros L, Touzain F, et al. Genome-Wide Identification of Host-Segregating Epidemiological Markers for Source Attribution in Campylobacter jejuni. Appl Environ Microbiol. 2017 Apr 1; 83(7). https://doi.org/10.1128/AEM.03085-16 PMID: 28115376

**40.** Rizk G, Lavenier D, Chikhi R. DSK: k-mer counting with very low memory usage. Bioinformatics. 2013 Mar; 29(5):652–3. https://doi.org/10.1093/bioinformatics/btt020 PMID: 23325618

**41.** Sheppard SK, Didelot X, Meric G, Torralbo A, Jolley KA, Kelly DJ, et al. Genome-wide association study identifies vitamin B5 biosynthesis as a host specificity factor in Campylobacter. Proceedings of the National Academy of Sciences. 2013 Jul; 110(29):11923–7. https://doi.org/10.1073/pnas.1305559110 PMID: 23818615

**42.** Yahara K, Méric G, Taylor AJ, Vries SPW de, Murray S, Pascoe B, et al. Genome-wide association of functional traits linked with Campylobacter jejuni survival from farm to fork. Environmental Microbiology. 2017; 19(1):361–80. https://doi.org/10.1111/1462-2920.13628 PMID: 27883255

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

**44.** Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine Learning in Python. Journal of Machine Learning Research. 2011; 12:2825–30.

**45.** Kingma DP, Ba J. Adam: A Method for Stochastic Optimization. arXiv:14126980 [cs] [Internet]. 2014 Dec [cited 2019 Sep 17]; Available from: http://arxiv.org/abs/1412.6980

**46.** Hedge J, Wilson DJ. Bacterial Phylogenetic Reconstruction from Whole Genomes Is Robust to Recombination but Demographic Inference Is Not. mBio [Internet]. 2014 Dec 31 [cited 2020 Nov 18]; 5(6). Available from: https://mbio.asm.org/content/5/6/e02158-14 https://doi.org/10.1128/mBio.02158-14 PMID: 25425237

**47.** Ansari MA, Didelot X. Bayesian Inference of the Evolution of a Phenotype Distribution on a Phylogenetic Tree. Genetics. 2016 Sep 1; 204(1):89–98. https://doi.org/10.1534/genetics.116.190496 PMID: 27412711

**48.** Argimón S, Abudahab K, Goater RJE, Fedosejev A, Bhai J, Glasner C, et al. Microreact: visualizing and sharing data for genomic epidemiology and phylogeography. Microbial Genomics,. 2016; 2(11): e000093. https://doi.org/10.1099/mgen.0.000093 PMID: 28348833

**49.** Cody AJ, Maiden MC, Strachan NJ, McCarthy ND. A systematic review of source attribution of human campylobacteriosis using multilocus sequence typing. Eurosurveillance [Internet]. 2019 Oct [cited 2020 Jan 27]; 24(43). Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6820127/ https://doi.org/10.2807/1560-7917.ES.2019.24.43.1800696 PMID: 31662159

**50.** Sheppard SK, Jolley KA, Maiden MCJ. A Gene-By-Gene Approach to Bacterial Population Genomics: Whole Genome MLST of Campylobacter. Genes (Basel). 2012 Apr 12; 3(2):261–77. https://doi.org/10.3390/genes3020261 PMID: 24704917

**51.** Cody AJ, Bray JE, Jolley KA, McCarthy ND, Maiden MCJ. Core Genome Multilocus Sequence Typing Scheme for Stable, Comparative Analyses of Campylobacter jejuni and C. coli Human Disease Isolates. Journal of Clinical Microbiology. 2017 Jul; 55(7):2086–97. https://doi.org/10.1128/JCM.00080-17 PMID: 28446571

**52.** Austerlitz F, David O, Schaeffer B, Bleakley K, Olteanu M, Leblois R, et al. DNA barcode analysis: a comparison of phylogenetic and statistical classification methods. BMC Bioinformatics. 2009 Nov; 10 (14):S10. https://doi.org/10.1186/1471-2105-10-S14-S10 PMID: 19900297

**53.** Deneke C, Rentzsch R, Renard BY. PaPrBaG: A machine learning approach for the detection of novel pathogens from NGS data. Scientific Reports. 2017 Jan; 7:39194. https://doi.org/10.1038/srep39194 PMID: 28051068

**54.** Chen X, Ishwaran H. Random Forests for Genomic Data Analysis. Genomics. 2012 Jun; 99(6):323–9. https://doi.org/10.1016/j.ygeno.2012.04.003 PMID: 22546560

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

**56.** Kwan PSL, Birtles A, Bolton FJ, French NP, Robinson SE, Newbold LS, et al. Longitudinal Study of the Molecular Epidemiology of Campylobacter jejuni in Cattle on Dairy Farms. Applied and Environmental Microbiology. 2008 Jun; 74(12):3626–33. https://doi.org/10.1128/AEM.01669-07 PMID: 18424539

**57.** Méric G, McNally A, Pessia A, Mourkas E, Pascoe B, Mageiros L, et al. Convergent Amino Acid Signatures in Polyphyletic Campylobacter jejuni Subpopulations Suggest Human Niche Tropism. Genome Biology and Evolution. 2018 Mar 1; 10(3):763–74. https://doi.org/10.1093/gbe/evy026 PMID: 29452359