# High performance *Legionella pneumophila* source attribution using genomics-based machine learning classification

*Andrew H. Buultjens[1,2,#], Koen Vandelannoote[3], Karolina Mercoulia[4], Susan Ballard[4], Clare Sloggett[4], Benjamin P. Howden[2,4,5], Torsten Seemann[4] and Timothy P. Stinear[1,2]*

1. Department of Microbiology and Immunology, Doherty Institute for Infection and Immunity, University of Melbourne, Melbourne, Victoria, Australia.
2. Centre for Pathogen Genomics, University of Melbourne, Melbourne, Victoria, Australia
3. Bacterial Phylogenomics Group, Institut Pasteur du Cambodge, Phnom Penh, Cambodia
4. Microbiology Diagnostic Unit, Department of Microbiology and Immunology, Doherty Institute for Infection and Immunity, University of Melbourne, Melbourne, Victoria, Australia
5. Department of Infectious Diseases, Austin Health, Heidelberg, Victoria, Australia

# Corresponding author: buultjensa@unimelb.edu.au

21    **ABSTRACT:**

22    Fundamental to effective Legionnaires' disease outbreak control is the ability to

23    rapidly identify the environmental source(s) of the causative agent, *Legionella*

24    *pneumophila*. Genomics has revolutionised pathogen surveillance but *L.*

25    *pneumophila* has a complex ecology and population structure that can limit source

26    inference based on standard core genome phylogenetics. Here we present a

27    powerful machine learning approach that assigns the geographical source of

28    Legionnaires' disease outbreaks more accurately than current core genome

29    comparisons. Models were developed upon 534 *L. pneumophila* genome sequences,

30    including 149 genomes linked to 20 previously reported Legionnaires' disease

31    outbreaks through detailed case investigations. Our classification models were

32    developed in a cross-validation framework using only environmental *L. pneumophila*

33    genomes. Assignments of clinical isolate geographic origins demonstrated high

34    predictive sensitivity and specificity of the models, with no false positives or false

35    negatives for 13 out of 20 outbreak groups, despite the presence of within-outbreak

36    polyclonal population structure. Analysis of the same 534-genome panel with a

37    conventional phylogenomic tree and a core genome multi-locus sequence type

38    allelic distance-based classification approach revealed that our machine learning

39    method had the highest overall classification performance – agreement with

40    epidemiological information. Our multivariate statistical learning approach

41    maximises use of genomic variation data and is thus well-suited for supporting

42    Legionnaires' disease outbreak investigations.

43

44

**INTRODUCTION:**

*Legionella pneumophila* is a gram negative bacterium that can thrive in warm, moist built environments and then cause Legionnaires' disease (LD) in humans when contaminated water is aerosolised and inhaled (David et al., 2016; Fields, Benson, & Besser, 2002; Mercante & Winchell, 2015; Schwake, Garner, Strom, Pruden, & Edwards, 2016). The vast majority of clinical infections are caused by *L. pneumophila* serogroup 1 (Yu et al., 2002). To combat LD outbreaks, public health authorities must rapidly investigate and determine the environmental sources to then intervene to prevent further disease transmission. A major difficulty in pin-pointing source(s) is the fact that there often exist a multitude of possible origins, particularly in densely populated urban settings.

The advent of bacterial genotyping has been advantageous for LD outbreak investigations, helping to 'rule in' or 'rule out' suspected environmental sources by attempting to match the genotypes of *L. pneumophila* recovered from patients to those derived from a given environmental source. In particular, Sequence Based Typing (SBT) compares DNA sequence variations across seven core genes to generate a sequence type (ST) that is standardised and internationally recognised (Lück, Fry, Helbig, Jarraud, & Harrison, 2013). An ST can be used to assign isolates from clinical specimens to specific environmental sources. Despite its popularity and simple interpretation, the SBT scheme lacks discriminatory power. The scheme captures only a tiny fraction of bacterial genomic variation and this is problematic when the majority of LD cases are caused by just a handful of STs (Borchardt, Helbig,

3

68      & Lück, 2008; David et al., 2016; Harrison, Afshar, Doshi, Fry, & Lee, 2009). SBT is

69      thus largely inadequate for LD source investigations.

70

71      Whole genome sequencing is used increasingly routinely for public health

72      surveillance and infectious disease outbreak investigations and recent efforts have

73      utilised the power of genomics to confirm suspected bacterial pathogen

74      environmental sources (Abrams & Trees, 2017; Goldberg, Sichtig, Geyer, Ledeboer,

75      & Weinstock, 2015; Krøvel et al., 2022; Petzold, Prior, Moran-Gilad, Harmsen, &

76      Lück, 2017; Ricci et al., 2022; Rousseau et al., 2022; Schoonmaker-Bopp et al., 2021;

77      Wüthrich et al., 2019). In particular, genomic analyses that assess core-genome

78      variation (sites present in all isolate genomes) such as phylogenomic trees and

79      pairwise SNP distances, have been useful to investigate disease transmission (Gorrie

80      et al., 2021; Ingle, Howden, & Duchene, 2021; Kwong et al., 2016; Sintchenko &

81      Holmes, 2015).

82

83      Another genomics-based approach for *L. pneumophila* source tracking is the core

84      genome multi locus sequence typing (cgMLST) scheme that builds upon the SBT

85      concept but greatly expands the genomic variation that is considered (Moran-Gilad

86      et al., 2015). In cgMLST, the allele scheme is enlarged from seven core genes to a

87      panel of 1,521 genes to produces an allele-type integer for each novel variant

88      combination (Moran-Gilad et al., 2015). This systematised and expanded approach

89      provides greater discrimination compared to conventional SBT, however like

90      phylogenomic approaches, it is still limited to only core-genome variation. Despite

91      the increased utility of such core-genome based approaches compared with SBT,

92    they still lack adequate discriminatory power for investigation of some *L.*

93    *pneumophila* outbreaks where isolate genomes are often near identical at the core-

94    genome level (Buultjens et al., 2017; McAdam et al., 2014; Sánchez-Busó et al.,

95    2016).

96

97    An alternative to core-genome analyses is to incorporate variation in accessory

98    genome sites; that is, to use DNA sequences present in some but not all isolates.

99    Here, to make better use of all the available genomic variation, we have developed a

100    machine learning statistical modelling method that utilises SNP variation in both the

101    accessory and core genome (pan-genome SNP variation) to classify genomes by

102    likely environmental source. Our approach integrates pan-genome SNP variation

103    using multivariate algorithms that model interrelationships among multiple variables

104    to assign source with greater accuracy than a standard core-genome SNP

105    comparison approach. This advance builds on our previously reported *L.*

106    *pneumophila* source tracking modelling approach that had high positive classification

107    capacity (rule-in) but had no negative classification ability (rule-out) (Buultjens et al.,

108    2017).

109

110    In this study, we have implemented 'one-versus-rest' machine learning classifier

111    algorithms with the ability to reject *L. pneumophila* clinical isolate genomes that

112    don't belong to classes used to train models, achieving both high classification

113    sensitivity and specificity. We have benchmarked the classification performance of

114    our machine learning method against phylogenomic and cgMLST allele distance

115    approaches using epidemiological assignments. Our machine learning algorithms

5

116     built with pan-genome SNP variants allowed us to assign the environmental sources

117     of LD outbreaks and make objective assignments of clinical isolate genome origins. It

118     is envisioned that future LD public health investigations may make use of such

119     sensitive and specific multivariate modelling advancements to rapidly identify the

120     environmental source of *L. pneumophila* and reduce the spread of this preventable

121     disease.

122

123

124     **METHODS:**

125     **Bacterial genomes used in this study:**

126     The isolate genomes originating from this study were cultured and sequenced as per

127     previously described (Buultjens et al., 2017). WGS data for an international collection

128     of diverse *L. pneumophila* (spanning 23 STs) was included in this study

129     (Supplementary Table. S1). A total of 246 isolates in this study were newly

130     sequenced while 288 were publicly available as either draft genome assemblies or

131     raw reads.

132

133     **Reference based core genome SNP calling:**

134     Snippy v4.4.5 was used to map reads and contigs to a previously described fully

135     assembled *L. pneumophila* clinical isolate genome Lpm7613 originating from

136     Melbourne, Australia (GenBank assembly accession: GCA_900092465.1) using a

137     'minfrac' setting of 0.8 (https://github.com/tseemann/snippy). The snippy-core

138     subcommand was used to generate a core genome SNP alignment - SNP variation in

139     the fraction of the genome shared by all isolates. Pairwise SNP differences were

6

140    assessed using a custom R script (https://github.com/MDU-

141    PHL/pairwise_snp_differences).

142

143    **Phylogenomic tree analysis:**

144    Clonal Frame ML was used to infer sites impacted by recombination (Didelot &

145    Wilson, 2015). The regions predicted to have been affected by recombination were

146    used to generate a bed file that was subsequently used for masking of the core

147    genome alignment with Snippy (see above). A maximum likelihood phylogenomic

148    tree was built from the alignment of non-recombining core SNPs using FastTree

149    v2.1.10 (Price, Dehal, & Arkin, 2009). Trees were displayed using FigTree v1.4.4

150    (http://tree.bio.ed.ac.uk/software/figtree). The cophenetic function of the ape R

151    package v5.6-2 (Paradis, Claude, & Strimmer, 2004) was used to compute a 534 ×

152    534 patristic distance matrix from the tree newick file.

153

154    **Core genome Multi Locus Sequence Typing:**

155    Core genome MLST analysis was undertaken using Coreugate v2.0.5

156    (https://github.com/kristyhoran/Coreugate). Draft genome assemblies were

157    generated by shovill v0.9.0 (https://github.com/tseemann/shovill) using the SPAdes

158    genome assembler v3.15.2 (Bankevich et al., 2012) and provided as input for

159    Coreugate (filter_samples_threshold=0.85). The *L. pneumophila* allele scheme used

160    with Coreugate was described by (Moran-Gilad et al., 2015).

161

162    **Reference independent pan genome SNP calling:**

7

163    Split Kmer Analysis (SKA) v1.0 was used to detect pan-genome SNPs (SNPs in core

164    and accessory sites) from reads and assembly contigs (Harris, 2018). Raw reads were

165    trimmed of adapter sequences using Trimmomatic v0.39 using the '-phred33' option

166    (Bolger, Lohse, & Usadel, 2014). Here, the fastq and fasta subcommands were used

167    to generate split kmer files (kmer size of 15) from isolates with reads in the fastq file

168    format and assembly contigs in fasta format, respectively. The split kmer files were

169    combined using the align subcommand (p=0.1) to produce a reference-independent

170    pan-genome SNP alignment and the humanise subcommand was used to generate a

171    SNP matrix from the skf alignment file.

172

173    **Distance-based classification:**

174    Matrices of pairwise distances were generated from both the phylogenomic tree and

175    the cgMLST alleles and used to devise distance-based classifiers. The average

176    distance among the environmental isolates for each outbreak group was calculated

177    and used as the outbreak group specific cut-off threshold to then classify the 113

178    clinical isolate genomes as either being outbreak related or not. This analysis was

179    conducted only for outbreak groups that had at least two environmental isolate

180    genomes available (14 of the 20 groups) (Table. 1).

181

182    **Classifier evaluation:**

183    Performance of all classifiers was assessed using the F1 metric. The F1-score is the

184    harmonic mean of the recall and precision, conveying the balance between these

185    two metrics. Here, a F1-score of 1 indicates that the classifier performs perfectly (no

186    false positives or false negatives). The F1-score is particularly useful to appraise

187    classification models when there is class imbalance.

188

189    **Machine learning classification framework:**

190    **Preparation of test and train datasets:**

191    The SKA pan-genome SNP matrix was one-hot encoded using the scikit-learn library

192    pre-processing module (Géron, 2019). The encoded matrix was divided into separate

193    training and testing datasets upon whether the isolate genomes were sourced from

194    either environmental samples, for training (n=421), or clinical samples, for testing
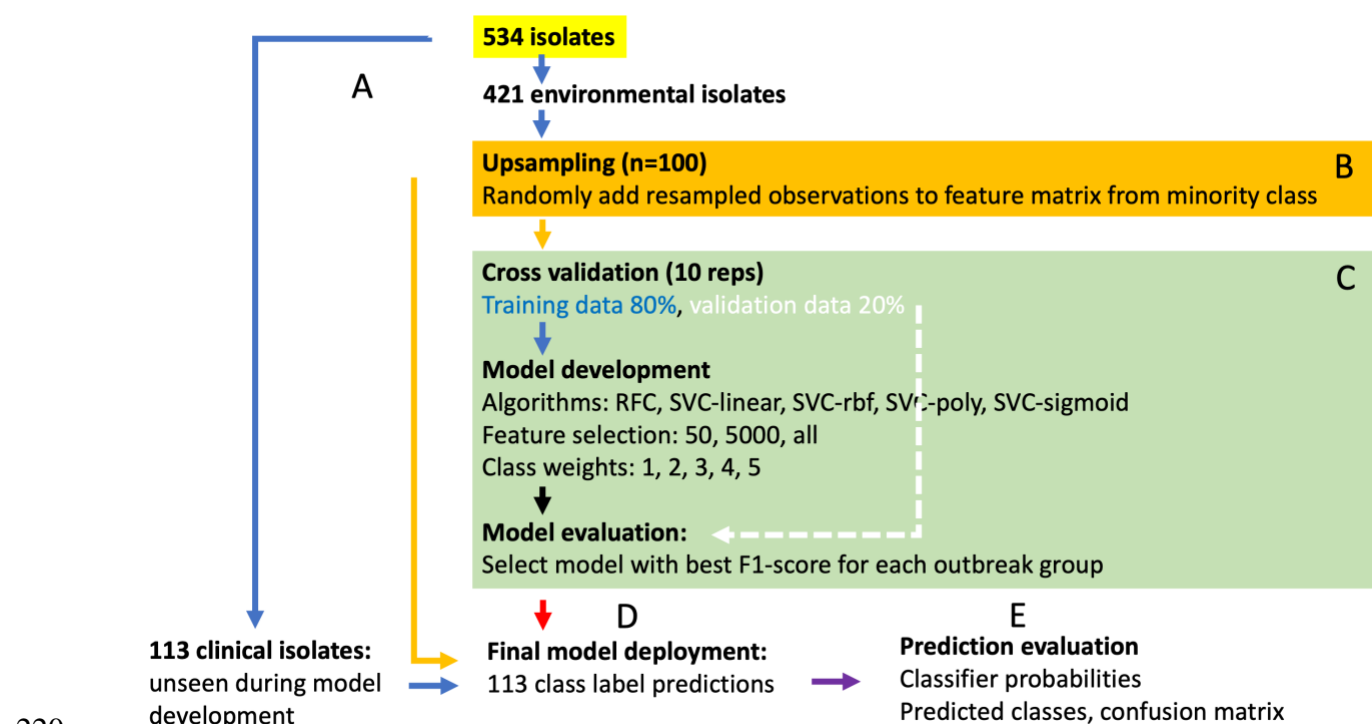
195    (n=113) (Fig. 1A).

196

197    **Model development:**

198    As the available epidemiological information was discrete geographical locations, a

199    supervised classification approach was used. Here the class labels were formatted to

200    represent a binary array of '1' (linked to outbreak) and '0' (not linked to outbreak).

201    The use of separate label files for each outbreak cluster allowed for the

202    implementation of a 'one-vs-rest' classification framework, in which each outbreak

203    group had its own model built, with the learning objective to include isolates of class

204    '1' and reject those of class '0'.

205

206    **Upsampling to redress class imbalance:**

207    Due to the availability of few outbreak-associated environmental isolate genomes

208    compared to that of clinical isolate genomes, there existed a substantial imbalance

209    between the '0' and '1' classes in the training set. To reduce this class imbalance,

210   *upsampling* was implemented in which observations from the minority class were

211   randomly selected (with replacement) and appended to the feature matrix (Fig. 1B).

212   As each outbreak group had a different set of labels, this was undertaken for all 20

213   outbreak groups. Given that there were approximately 20 minority class

214   environmental isolate genomes to 400 majority class observations, an upsampling

215   amount of 100 was chosen as this was approximately 1/4 of the majority class in

216   each situation - a conservative upsampling portion given the severe class imbalance.

217   The remaining class imbalance was addressed through specifying class weights to the

218   classification algorithm (see below).

219



221   **Fig. 1.** Flow diagram of the machine learning model development framework. A)

222   Isolate genomes were separated from the input one-hot encoded matrix (n=534)

223   according to being either environmentally (n=421) or clinically derived (n=113). B)

224   Upsampling was performed on the environmental training dataset, where individuals

10

225   from the minority class were randomly upsampled (with replacement). C) Cross

226   validation loop. In each iteration, the training data was randomly split into a training

227   and testing partition of 80% and 20%, respectively, for ten repetitions. Various

228   combinations of model parameters were used, and the classifier was evaluated upon

229   ability to correctly assign the test set component of the data using the F1-score. D)

230   The models for each outbreak group with the greatest F1-score in the cross-

231   validation loop were selected to form a set of final models. Final models were

232   trained with all available upsampled environmental isolate genome data to then

233   assign the classes of the previously unseen 113 clinical isolate genomes. E)

234   Classification outputs were in the form of probabilities that were binarised as either

235   belonging to or not belonging to each specific outbreak group class. Information of

236   clinical isolate known origins was used to establish a confusion matrix and calculate

237   the F1-score.

238

239   **Model development cross-validation:**

240   During model training, supervised classification algorithms learn specific patterns

241   associated with each of the classes with the goal to develop models that are

242   generalisable, in that they can make accurate assignments upon previously unseen

243   observations. To promote optimal model development on the environmental isolate

244   training dataset, an iterative cross-validation procedure was undertaken to

245   determine the best model for each outbreak group (Fig. 1C). Here, the training data

246   was randomly split into training and validation partitions (80% train and 20%

247   validation) 10 times, with models built upon the training portion and used to classify

248   the classes of the validation portion. For each iteration in the cross-validation loop,

11

249    the F1-score was recorded and used to evaluate each model. A different set of

250    model parameter combinations was evaluated with each cross-validation iteration

251    (model parameters: classifier algorithm, class weights, and number of selected

252    features) (Fig. 1C). A total of 1,500 model combinations were evaluated in the cross-

253    validation phase.

254

255    **Multivariate classification algorithms:**

256    Two supervised classifier algorithms were implemented: Random Forest Classifiers

257    (RFC) and Support Vector Classifiers (SVC) (Fig. 1C). RFC indiscriminately select a

258    subset from the training data to create a collection of decision tree predictors to

259    sum the predictions, in effect lowering the variance (Breiman, 1996). Here, each

260    decision tree takes a set of features and provides an individual output, all of which

261    are subsequently summarised to produce a final probabilistic output (Breiman,

262    2001). The scikit-learn RFC module was implemented with default parameters

263    (Géron, 2019). SVC optimise for non-linear combinations of features that best divide

264    the classes across a multi-dimensional hyperplane (Boser, Guyon, & Vapnik, 1992).

265    The scikit-learn SVC module was implemented with default parameters apart from

266    using kernels: 'linear', 'rbf', 'poly' and 'sigmoid' (Géron, 2019).

267

268    **Class weights:**

269    A further approach to combat the occurrence of class imbalance was to specify class

270    weights to the classification algorithms. The reasoning here was that classifiers have

271    default assumptions of class balance and, when faced with class imbalance, a bias

272    exists that favours towards the dominant class. In this case the '0' or 'not outbreak

12

273    related' isolates are likely to cause bias, as they strongly outnumber the amount of

274    '1' or 'outbreak related' isolates. By specifying class weights the classification

275    algorithm is modified to account for the skewed class distribution, enabling

276    improved training and higher performance assignments by penalising

277    misclassification of the minority class. Specifically, the class weights were passed to

278    the scikit-learn classifiers as a dictionary that stipulated class '0' as 0.5 and class '1'

279    as an integer in the range of 1 to 5 (Fig. 1C).

280

281    **Univariate feature selection:**

282    Features that did not vary in proportion between the classes for a particular

283    outbreak group are unlikely to have any classification value for model training and

284    therefore only add noise. To reduce the number of uninformative features and focus

285    on those that are associated with the class labels, feature selection was performed.

286    The SelectKBest univariate module of scikit-learn was employed to assess the

287    independence of individual features against the target variable using a chi-square

288    test, selecting the top 50, 5,000 or all features (Fig. 1C) (Géron, 2019). To avoid any

289    data-leakage, the univariate feature selection was only performed on the training set

290    either during the cross-validation procedure or on all available environmental

291    isolates for the building of the final models (see below).

292

293    **Final model classifications:**

294    Following selection of the top performing model combinations for each outbreak

295    group, final models were built using the model parameters identified and trained

296    with all available environmental isolates (n=421). Here, the final model for each

13

297    outbreak group learned as much as possible about the genomic variability in the

298    data when all available environmental isolates were used (Fig. 1D-E). Thus, this was

299    the optimal way to train the final models to make generalisable source attribution

300    assignments upon the clinical isolate genomes.

301

302    The code used to conduct the abovementioned analyses is detailed in the following

303    github repository:

304    https://github.com/abuultjens/Assign_Legionella_pneumophila_origins

305

306    **RESULTS:**

307    **Selection of *L. pneumophila* genome sequences for classification model**

308    **development:**

309    The overall objective of this research was to attempt to use multivariate statistical

310    learning methods to assign the environmental sources of LD outbreaks. However, to

311    benchmark the performance of such methods it was first necessary to select a set of

312    *L. pneumophila* genomes representing different LD outbreak investigations. Our

313    principles for genome selection were to maximise both genomic and spatial diversity

314    to achieve a collection that spanned many Sequence Types (STs) and originated from

315    various locations worldwide. A review of the literature and publicly available *L.*

316    *pneumophila* genome sequences revealed studies from three different jurisdictions

317    (see details below) spanning 20 distinct LD outbreaks that were suitable to include

318    because they had sufficiently rich epidemiological information and associated *L.*

319    *pneumophila* genomic data from both clinical and environmental sources. In all, 534

320    *L. pneumophila* genomes were identified for use in this study, of which 421 and 113

14

321     represented bacterial isolates from environmental and clinical sources, respectively

322     (Table. S1).

323

324     The outbreak associated group consisted of 149 isolates that were epidemiologically

325     linked to a total of 20 outbreaks across three major geographical regions: 1)

326     Melbourne, Victoria, Australia, 2) Essex, England, and 3) New York State, United

327     States (Table. S1). The Melbourne *L. pneumophila* genomes, hereon referred to with

328     prefix "MELB", represented five different LD outbreaks spread across the Melbourne

329     metropolitan area, occurring between 1998-2018 (Buultjens et al., 2017). The Essex

330     *L. pneumophila* genomes, hereon referred to with prefix "ESSEX", consisted of

331     genomes obtained from *L. pneumophila* isolates linked to LD disease occurring in five

332     distinct wards within a single hospital campus (isolated between 2007-2011) (David

333     et al., 2017). The New York State *L. pneumophila* genomes, hereon referred to with

334     prefix "NY", consisted of 10 separate LD outbreaks across the New York State area (*L.*

335     *pneumophila* isolated between 2004-2012) (Raphael et al., 2016).

336

337     To assist in developing a classification framework with negative classification

338     capacity, *i.e.* the ability of the model to call true negatives, we included genome

339     sequences from 74 *L. pneumophila* clinical isolates not associated with any of the

340     abovementioned outbreaks, hereon referred to as clinical non-outbreak associated

341     (CNOA) (Table. 1). These isolate genomes were isolated between 1986-2014 and

342     originated from across Europe, the United Kingdom and Australia. In a similar way,

343     to challenge the model building process, we included 311 environmental isolates

344     (isolated between 1995-2018) that were not associated with any of the outbreaks

345  (MELB, ESSEX or NY), hereon referred to as the environmental non-outbreak

346  associated (ENOA) (Table. 1).

347

348  **Table 1.** Attributes of the 534 *L. pneumophila* isolates included in this study.

| Group | Number of environmental isolates | Number of clinical isolates | STs | Reference |
|---|---|---|---|---|
| ENOA | 311 | NA | 15 SBTs | This study; Bartley, PB., et. al., 2016; Buultjens, AH., et. al., 2017; David, S., Rusniok, C., et. al., 2016; David, S., et. al., 2017; Moran-Gilad, J., et. al., 2015; Qin, T., et.al., 2016 |
| CNOA | NA | 74 | 9 SBTs | Bartley, PB., et. al., 2016 ; Buultjens, AH., et. al., 2017 ; David, S., Rusniok, C., et. al., 2016; David, S., et. al., 2017; Moran-Gilad, J., et. al., 2013 |
| MELB-2018 | 20 | 3 | SBT30 | This study |
| MELB-A | 14 | 11 | SBT30 | Buultjens, AH., et. al., 2017 |
| MELB-C | 3 | 1 | SBT30 | Buultjens, AH., et. al., 2017 |
| MELB-G | 18 | 2 | SBT30 | Buultjens, AH., et. al., 2017 |
| MELB-M | 8 | 1 | SBT30 | Buultjens, AH., et. al., 2017 |
| ESSEX-A | 7 | 2 | SBT1 | David et al., 2017 |
| ESSEX-B | 3 | 1 | SBT1 | David et al., 2017 |
| ESSEX-E | 2 | 1 | SBT1 | David et al., 2017 |
| ESSEX-G | 14 | 1 | SBT1 | David et al., 2017 |
| ESSEX-H | 5 | 2 | SBT1 | David et al., 2017 |

| NY-1 | 3 | 1 | SBT1 | Raphael, B. Baker, D., et. al., 2016 |
|------|---|---|------|--------------------------------------|
| NY-2 | 2 | 2 | ND | Raphael, B. Baker, D., et. al., 2016 |
| NY-3 | 1 | 3 | SBT1 | Raphael, B. Baker, D., et. al., 2016 |
| NY-4 | 1 | 1 | SBT1 | Raphael, B. Baker, D., et. al., 2016 |
| NY-5 | 1 | 1 | SBT62 | Raphael, B. Baker, D., et. al., 2016 |
| NY-6 | 3 | 1 | SBT36 | Raphael, B. Baker, D., et. al., 2016 |
| NY-7 | 1 | 1 | SBT36 | Raphael, B. Baker, D., et. al., 2016 |
| NY-8 | 1 | 1 | SBT1204 | Raphael, B. Baker, D., et. al., 2016 |
| NY-9 | 2 | 1 | SBT94 | Raphael, B. Baker, D., et. al., 2016 |
| NY-10 | 1 | 2 | SBT731 | Raphael, B. Baker, D., et. al., 2016 |
| **Total** | 421 | 113 | | |

349

**Population structure of *L. pneumophila* isolates used in this study:**

351 We examined the genomic context of 421 environmental *L. pneumophila* isolate

352 genomes alongside 113 clinical isolate genomes to investigate the ability to make

353 inferences of source attribution. The 20 outbreak groups were from three distinct

354 geographical regions, Melbourne (Australia), Essex (UK) and New York (US).

355 Sequence read alignment against a SBT30 reference genome revealed 221,214 core

356 genome SNPs. There were 144,829 SNP sites inferred to have arisen by

357 recombination, leaving 76,385 SNPs that were derived through vertical transmission.

358 Pairwise SNP comparisons were performed to depict the amount of diversity within

359 each outbreak group (Fig. 2A). Most of the groups had mean intra-group distances

360 between 0-4 SNPs, while MELB-A, NY-3 and NY-8 had elevated within group

361 variations of 41, 61 and 24 SNPs, respectively.

362

363    Phylogenomic analysis has become an important approach to examine pathogen

364    population structure and to investigate the likely origins of *L. pneumophila* clinical

365    isolates using WGS data (David et al., 2016; Gorzynski et al., 2022; Graham, Doyle, &

366    Jennison, 2014; Qin et al., 2016; Reuter et al., 2013; Wüthrich et al., 2019). A

367    phylogenomic tree was estimated from the non-recombining core-genome SNP

368    alignment to depict the clonal ancestry (Fig. 2B). The tree illustrated the same

369    grouping of outbreak related isolates that was observed with the pairwise SNP

370    distance analysis. In particular, the groups with high internal SNP diversity displayed

371    the existence of within-outbreak polyclonal population structure (Fig. 2B). The

372    MELB-A isolate genomes were found to harbor several distinct genotypes, one of

373    which was exclusively represented by clinical isolates (Fig. 1B-C). Outbreak group NY-

374    3 isolate genomes were located across several distinctive subtrees in the phylogeny,

375    indicating a within group polyclonal population structure (Fig. 1B). NY-8 isolate

376    genomes had an elevated within group diversity while also being substantially

377    distinct to all other isolates included in the study (Fig. 1B).
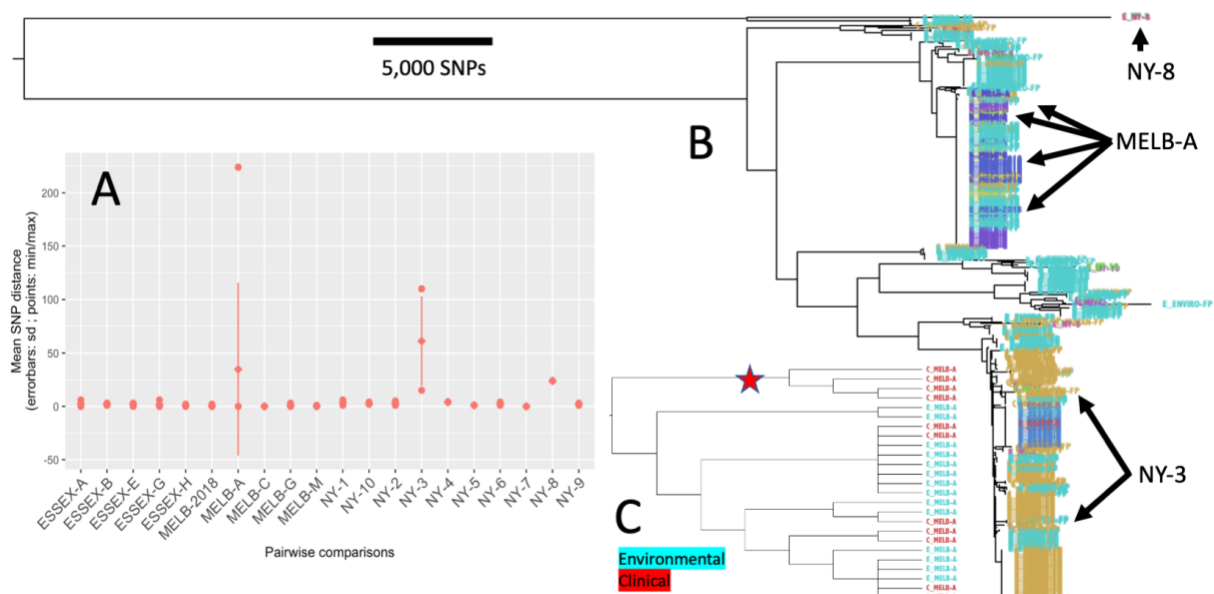
378

**Fig 2.** Assessment of genomic population structure of 534 *L. pneumophila* clinical and environmental isolate genomes. A) Pairwise SNP comparisons of within outbreak group diversity. Three groups had elevated levels of within group diversity: MELB-A, NY-3 and NY-8. B) Phylogenomic tree generated from non-recombining core genome SNPs. Outbreak groups MELB-A, NY-3 and NY-8 are indicated C) Subtree containing isolate genomes associated with the MELB-A outbreak. The subtree is displayed as a cladogram with branch lengths transformed to illustrate the tree topology. Red star indicates a distinct genotype containing only clinical isolate genomes without any environmental representatives.

**Phylogenomic tree distance-based classification:**

To objectively assess the ability to infer clinical isolate origins from the phylogenomic tree, patristic distances were extracted and used to build outbreak group specific classifiers. Here, the patristic distances represent the individual total branch length

19

394     distances between all possible isolate pairs in the tree, represented as a 534 x 534

395     distance matrix. The average distance between the environmental isolate genomes

396     of each outbreak group were calculated for groups that had at least two or more

397     environmental representatives (14 of the 20 outbreak groups). The average distance

398     among environmental isolate genomes was used as a threshold to assign each query

399     clinical isolate as either related or unrelated to the outbreak groups, with each group

400     having a specific threshold distance (14 different thresholds and classifiers). The

401     assumption underlying the use of distance thresholds was that a clinical isolate

402     genome with equal or less patristic distance from the mean distance observed

403     among environmental isolate genomes from a specific outbreak group is likely

404     related to that outbreak while those with greater distances are more divergent and

405     thus likely originated elsewhere.

406

407     The cut-off distance threshold for each outbreak group was determined through

408     analysis of only the environmental isolate genomes for each specific outbreak group.

409     This is an ideal approach, as the thresholds are not biased by the addition of any

410     clinical isolate genomes, therefore building a classification tool that is prospective, in

411     that the system would be ready for deployment before the first clinical isolate

412     genome is reported in an outbreak investigation. The performance of the classifiers

413     was assessed using the F1-score which is the harmonic mean of method recall and

414     precision, conveying the balance between these two metrics. Here, a F1-score of 1

415     indicates that the classifier performs perfectly (no false positives or false negatives).

416     The patristic distance-based classifiers demonstrated the ability to correctly assign

417     most clinical isolate genomes to their known origins (0.43 mean false negatives),

20

418    however this approach had a high false positive rate (3.93 mean false positives) with

419    an overall mean F1 score of 0.50 (Table. 2).

420

421    **Table 2.** Distance-based and machine learning classifications for 113 test set clinical

422    isolate genomes when trained on a set of 421 training environmental isolate

423    genomes.

| Outbreak group | Patristic distance-based classifiers | | | cgMLST distance-based classifiers | | | Machine learning classifiers | | |
|---|---|---|---|---|---|---|---|---|---|
| | False positive | False negative | F1-score | False positive | False negative | F1-score | False positive | False negative | F1-score |
| MELB-2018 | 4 | 0 | 0.60 | 2 | 0 | 0.75 | 0 | 0 | 1 |
| MELB-A | 6 | 4 | 0.58 | 18 | 0 | 0.55 | 1 | 5 | 0.67 |
| MELB-C | 11 | 0 | 0.15 | 6 | 0 | 0.25 | 0 | 1 | 0 |
| MELB-G | 7 | 0 | 0.36 | 19 | 0 | 0.17 | 0 | 0 | 1 |
| MELB-M | 11 | 0 | 0.15 | 22 | 0 | 0.08 | 0 | 0 | 1 |
| ESSEX-A | 5 | 0 | 0.44 | 2 | 1 | 0.40 | 1 | 1 | 0.5 |
| ESSEX-B | 4 | 0 | 0.33 | 3 | 0 | 0.4 | 1 | 1 | 0 |
| ESSEX-E | 4 | 0 | 0.33 | 6 | 0 | 0.25 | 0 | 1 | 0 |
| ESSEX-G | 3 | 0 | 0.40 | 4 | 0 | 0.33 | 1 | 0 | 0.67 |
| ESSEX-H | 0 | 1 | 0.67 | 0 | 1 | 0.67 | 0 | 0 | 1 |
| NY-1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| NY-2 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| NY-3 | NA | NA | NA | NA | NA | NA | 0 | 3 | 0 |
| NY-4 | NA | NA | NA | NA | NA | NA | 0 | 0 | 1 |
| NY-5 | NA | NA | NA | NA | NA | NA | 0 | 0 | 1 |
| NY-6 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| NY-7 | NA | NA | NA | NA | NA | NA | 0 | 0 | 1 |
| NY-8 | NA | NA | NA | NA | NA | NA | 0 | 0 | 1 |
| NY-9 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |
| NY-10 | NA | NA | NA | NA | NA | NA | 0 | 0 | 1 |
| AVERAGE | | | | | | | 0.20 | 0.60 | 0.74 |
| | *3.93 | *0.43 | *0.50 | *5.86 | *0.14 | *0.56 | (*0.29) | (*0.64) | (*0.70) |

424    * When considering groups with two or more environmental isolate genomes

425

426 **cgMLST distance-based classification:**

427 In addition to phylogenomics, cgMLST is another genomic comparison approach

428 used to infer the source attribution of *L. pneumophila* clinical isolate genomes which

429 builds upon the established SBT genotyping method by greatly expanding the

430 number of core-genome loci (Moran-Gilad et al., 2015; Qin et al., 2016). The

431 advantage of cgMLST over analyses that consider all core genome SNPs is the

432 standardised framework in which the alleles are called, in that cgMLST is not

433 susceptible to fluctuations in core genome size caused by the addition or removal of

434 isolates from the analysis. We next investigated if the allelic distance derived from

435 the cgMLST scheme, when applied to the 534 isolates, could be used to provide

436 improved source attribution inference. Here, the same threshold derivation and

437 classification approach that was employed for the patristic distances was applied,

438 however using a distance matrix generated from cgMLST allelic variation.

439

440 The cgMLST based classifiers had fewer false negatives than the patristic distance-

441 based classifiers (0.14 mean false negatives) while having a higher false positive rate

442 (5.86 mean false positives) and a marginally higher overall mean F1-score of 0.56

443 (Table. 2). The classifiers performed well for NY outbreak groups that had more than

444 one environmental isolate genome, all achieving F1-scores of 1. While the

445 implementation of phylogenomic tree and cgMLST distance-based classifiers

446 introduced an objective framework to make source inferences, these approaches

447 were based solely on core-genome variation, raising the question of whether

448    approaches built using SNP variation from across the pan-genome may achieve

449    greater assignment capacity.

450

451    **Machine learning classification:**

452    To enhance the classification capacity of the framework, we applied a machine

453    learning approach that utilised an alignment containing 479,480 SNPs detected in

454    both core and non-core sites. The advantage of using pan-genome SNPs for this type

455    of analysis was that additional variation in accessory genome sites is thus

456    considered, improving the discriminatory potential for downstream analyses. In

457    addition to greater SNP variation, the use of a multivariate classification algorithm

458    provides the advantage in that the concerted effects of all input genomic variants

459    are modelled to learn about informative structures in the data.

460

461    To reduce the likelihood of overfitting, a cross-validation framework was established

462    that iteratively split the environmental isolate data into train and validation

463    partitions. A total of 1,500 model combinations consisting of different model

464    parameters using both Random Forest Classifiers (RFC) and Support Vector

465    Classifiers (SVC) (see methods) were evaluated. In this way, the best model

466    combination for each outbreak group was determined using only environmental

467    isolate genomic variation prior to the analysis ever encountering any clinical isolate

468    genomes, thus eliminating the risk of model overfitting, and providing a prospective

469    approach.

470

471    **Machine learning model results:**

472    Application of the final models for the assignment of the clinical isolate genomes

473    provided the lowest false positive rate of all previous distance-based approaches

474    (0.29 mean false positives), the highest level of false negatives (0.64 mean false

475    negatives) and the highest overall mean F1 score of 0.70 when applied to the 14

476    outbreaks with two or more environmental isolate genomes (Table. 2). Models

477    developed for outbreak groups MELB-2018, MELB-G, MELB-M, ESSEX-H, NY-1

478    through NY-2 and NY-4 through NY-10 (13/20) had F1-scores of 1, indicating the

479    absence of any false positives or false negatives – classifications that perfectly align

480    with the epidemiological labels (Table. 2). As the machine learning method used

481    upsampling to artificially replicate the training observations, it was possible to apply

482    this method to outbreak groups with as few as one environmental isolate genome,

483    having an overall mean F1-score of 0.74 when applied to all 20 outbreak groups.

484    (Table. 2). The parameters of the final models are reported in Supplementary Table

485    2.

486

487    **Examination of machine learning model false positives and false negatives:**

488    False positives occurred with models ESSEX-A, ESSEX-B and ESSEX-G. In these

489    instances, the false positives were from other ESSEX outbreak groups clinical isolate

490    genomes (wards within the same hospital). Six of the models MELB-A, MELB-C,

491    ESSEX-A, ESSEX-B, ESSEX-E and NY-3 had one or more false negative classifications.

492    In the case of NY-3, there was an appreciable amount of within outbreak diversity

493    (Fig. 2A) and just a single environmental isolate used for model training (Table. 1).

494    For MELB-A, the four clinical isolates that were classified as false negatives by the

495    machine learning approach were on a branch in the phylogeny that did not contain

496     any MELB-A environmental isolate genomes and therefore were from a specific

497     genotype that was not represented in the training data (Fig. 2C). Despite this, all 11

498     of the MELB-A clinical isolate genomes were within the top 21% of the 113 test-set

499     clinical isolate genomes when ranked according to decreasing classification

500     probability (Fig. 3A).

501

502     Investigation of the machine learning classifier probabilities for outbreak groups

503     ESSEX-A, ESSEX-B and ESSEX-E also revealed that despite having false negatives at

504     the default classification threshold of 0.5, the classification probabilities were

505     nonetheless informative to rank the clinical isolate genomes (Fig. 3B-D). In this way,

506     when ranked according to decreasing probabilities, the clinical isolate genomes from

507     ESSEX-A, ESSEX-B and ESSEX-E were contained withing the top 25%, 2% and 3% of all

508     clinical isolate genomes, respectively (Fig. 3B-D). In these instances, if the

509     classification threshold were lower than 0.5, these models would have provided

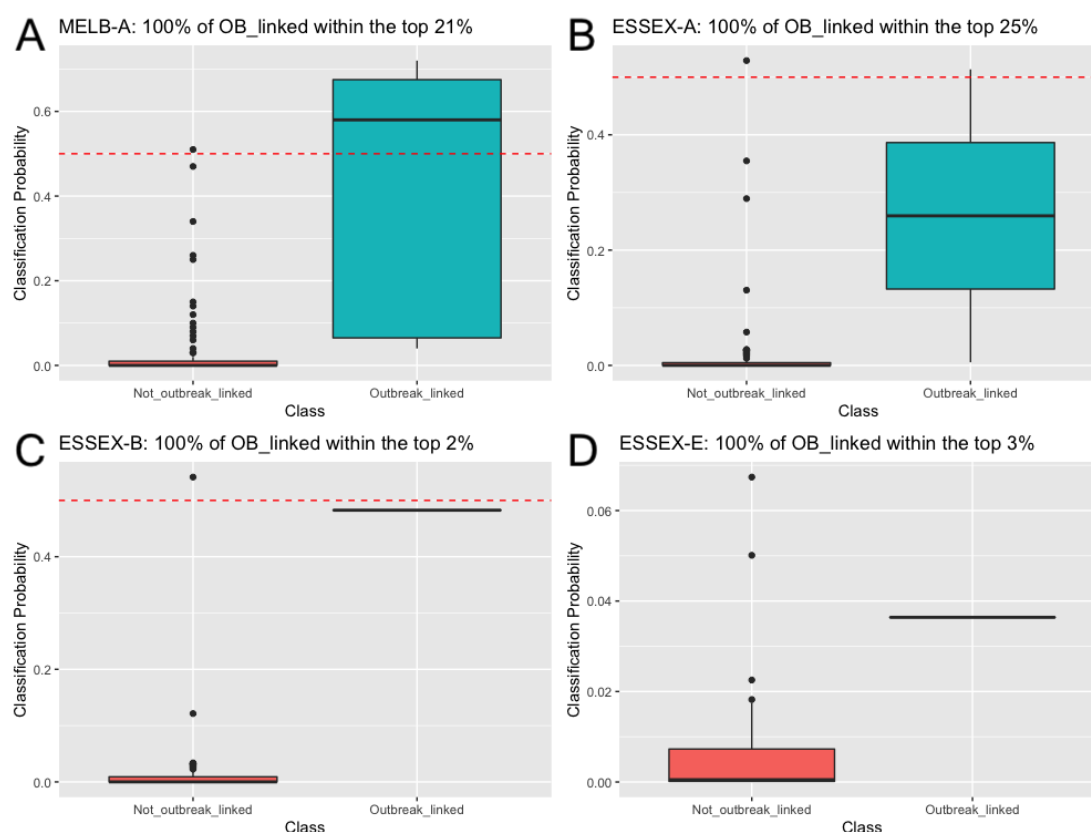510    perfect and near perfect classifications.



511

512    **Fig 3.** Boxplots of classification probabilities for the outbreak linked and non-

513    outbreak linked 113 test set clinical isolate genomes for four outbreak group models

514    that had false negative classifications. Red horizontal dotted lines indicate the

515    classification threshold of 0.5. A: classification probabilities for outbreak group

516    MELB-A. B: classification probabilities for outbreak group ESSEX-A. C: classification

517    probabilities for outbreak group ESSEX-B. D: classification probabilities for outbreak

518    group ESSEX-E.

519

520

521    **DISCUSSION:**

522    Timely and accurate identification of environmental sources of LD is of utmost

523    importance to public health investigations and, in this era of high-resolution genomic

524    technologies, innovative approaches are needed to rapidly distil complex analyses to

525    provide actionable insights. In this study, we have deployed a machine learning

526    classification approach and assessed its' ability alongside alternative approaches to

527    make assignments of clinical isolate origins that align against the known

528    epidemiological information for 20 distinct LD outbreaks.

529

530    This work builds on our previous efforts to build accurate multivariate assignment

531    models, here providing the necessary negative classification capacity that was

532    lacking in our earlier work. To assess the ability of these multivariate approaches to

533    call true negatives, we included 74 clinical isolates that were not associated with any

534    of the 20 outbreak groups that were used to train the models. In a similar way, we

535    also included 311 environmental isolates that were not associated with the outbreak

536    groups to assess how well the model could learn from known outbreaks while faced

537    with a larger than necessary training dataset that contained unrelated

538    environmental isolates. Our improved approach presented in this investigation made

539    use of a set of 'one-vs-rest' classification strategies, in which a separate target

540    variable and model was used for each outbreak group. This had the effect of

541    focusing on genomic variation that was specific to an individual outbreak group,

542    optimising the model to include outbreak linked isolates while rejecting others and

543    therefore affording negative classification capacity.

544

545    The analysis of suspected pathogen transmission with phylogenomic trees built from

546    core genome SNPs has become the *de facto* standard in the field of bacterial

547    genomics. Here we assessed the ability of patristic distances derived from a

27

548    phylogenomic tree to place epidemiologically linked isolate genomes into

549    arrangements that could then permit the inference of clinical isolate source

550    attribution. Classifiers were devised for the 14 of the 20 groups that had at least two

551    environmental isolate genomes, with assignment thresholds derived from the mean

552    distance observed among the environmental representatives of each group. This

553    approach provided an objective and quantitative phylogenomic-based framework

554    for the classification of query clinical isolate genomes with high sensitivity; however,

555    it suffered from low specificity and had an overall mean F1-score of 0.50.

556

557    Another widely employed tool for *L. pneumophila* genomic comparisons is cgMLST,

558    which builds on the established SBT method by greatly expanding the number of

559    core loci. To investigate the utility of this method to infer clinical isolate genome

560    source attribution, a matrix of cgMLST allelic distances was generated in the same

561    way that patristic distances were used to build distance-based classifiers. The results

562    from this approach were a slight improvement over the patristic distance-based

563    classifiers, with a higher overall mean F1 score of 0.56, however there were a higher

564    number of false positives, again offering meagre specificity and poor overall

565    classification capacity.

566

567    A machine learning classification framework was developed using pan-genome SNP

568    variants to make probabilistic assignments by firstly training models upon variation

569    among environmental isolate genomes to then classify the origins of clinical isolate

570    genomes. To achieve this, an extensive cross-validation framework was established

571    that assessed the performance of various model building parameters (see methods)

572    on the ability for an algorithm to learn upon a portion of environmental isolate

573    genomes and then assign the known classes of the remaining environmental

574    representatives (cross-validation), with the best classification models selected to

575    then learn using the entire training set to make assignments upon the previously

576    unseen clinical isolate genomes.

577

578    The application of the machine learning models for the assignment of 113 test set

579    clinical isolate genomes had the greatest classification capacity with 13 out of 20

580    models achieving an F1-score of 1, indicating perfect sensitivity and specificity. The

581    machine learning method also achieved the greatest overall mean F1 score of 0.70

582    when evaluating the 14 groups with two or more environmental representatives and

583    0.74 when applied to all 20 groups. The higher performance of the machine learning

584    modelling approach compared to phylogenomic tree branch length distance and

585    cgMLST allelic distance methods is likely since 1) it considered SNP variation across

586    the pan-genome, 2) it explicitly made use of the underlying sequence composition of

587    the SNP variation and 3) it employed a multivariate approach that modelled the

588    concerted interactions of all input variants. Together, these three aspects of the

589    modelling approach work to make efficient use of the richness of the available SNP

590    allelic variation to achieve greater classification capacity.

591

592    False positives were detected with machine learning models ESSEX-A, ESSEX-B and

593    ESSEX-G. Here, the false positives were from other wards in the same hospital,

594    suggesting a sort of 'cross reactivity' among nearby locations within a common

595    institution. Despite these false positives, the 74 unrelated clinical isolates were

29

596    correctly assigned as true negatives by all final models, indicating overall satisfactory

597    negative classification capacity. False negative assignments occurred with models

598    MELB-A, MELB-C, ESSEX-A, ESSEX-B, ESSEX-E and NY-3. In the case of MELB-C and

599    NY-3, previous analyses have identified that there likely exists an issue with the

600    epidemiological source attribution for these outbreak groups, offering a possible

601    explanation for the inability of the models to accurately assign these isolate

602    genomes to their known origins in previous investigations (Buultjens et al., 2017;

603    Raphael et al., 2016).

604

605    For MELB-A, the four clinical isolates assigned as false negatives by the machine

606    learning approach were on a branch in the phylogeny that did not contain any

607    environmental isolate genomes from the MELB-A outbreak group, meaning this

608    specific genotype was not represented in the training data. Despite this, all MELB-A

609    clinical isolate genomes were within the top 21% of all clinical isolate genomes when

610    ranked according to decreasing classification probability. This suggests that the

611    modelling approach was able to make use of the level of shared ancestry among all

612    MELB-A isolates to nevertheless provide a useful degree of probability ranking even

613    when that specific genotype was not explicitly represented in the training data. Not

614    dissimilar to what was seen with the MELB-A probability ranking, the classification

615    probabilities for the ESSEX-A, ESSEX-B and ESSEX-G clinical isolate genomes revealed

616    that the known positives for each of these groups were ranked highly despite being

617    less than the standard classification threshold of 0.5. This highlights that alternative

618    probability evaluation frameworks besides classification, such as probability ranking,

619    should be considered for these approaches.

30

620

621   In addition to the use of a 'one-vs-rest' classification approach, another notable

622   point of difference with this new method was the use of pan-genome SNPs derived

623   from the reference independent kmer-based method, SKA. Our previous work built

624   models using only variation in core-genome SNPs that were called using read

625   alignment to a reference genome (Buultjens et al., 2017). The consequence of using

626   pan-genome variation was particularly important in this application since the core-

627   genome among the diverse group of 534 *L. pneumophila* isolates is abbreviated,

628   therefore reducing the total amount of SNP diversity. Specifically, the pan-genome

629   alignment provided 258,266 more SNPs than when only core genome variants were

630   considered, equating to addition information to be learnt by multivariate

631   approaches.

632

633   The lack of environmental isolates representing a specific MELB-A genotype that was

634   observed exclusively among clinical isolates indicates that the methods used to

635   sample, culture and sequence *L. pneumophila* from environmental sources had failed

636   to adequately capture the true extent of bacterial diversity in that source. Efforts to

637   capture environmental *L. pneumophila* diversity typically involve taking multiple

638   colony picks from environmental samples. While care was taken in this approach to

639   maximise the environmental diversity captured, there evidently was relevant

640   diversity that did not progress to culture isolation and subsequent genome

641   sequencing, presumably due to the limited sensitivity of culture-based methods

642   (Reller, Weinstein, & Murdoch, 2003). Reduced detection of genomic diversity

643   among environmental samples compared to that recovered from clinical specimens

644    has been observed in a previous investigation (Wüthrich et al., 2019). Alternative

645    methods that would likely widen the capture of environmental diversity are shotgun

646    metagenomic or culture independent sequencing approaches that directly sequence

647    all environmental DNA, eliminating the bottleneck of culture (Christiansen et al.,

648    2014; Wéry et al., 2008).

649

650    All outbreak groups, apart from MELB-2018, MELB-G and ESSEX-G, had very few

651    numbers of environmental isolate genomes and in some cases just a single genome.

652    Such limited examples of environmental genomic diversity are not optimal and the

653    inclusion of greater numbers of training genomes for each group would likely

654    improve the ability of the models to learn about outbreak specific signatures and

655    make more accurate classifications.

656

657    While this study focused on SNP variation, there may be further genomic

658    information among additional variant types such as kmers counted directly from raw

659    reads that may further improve model performance. Such kmer variation has the

660    potential to capture additional genomic variations such as structural variations and

661    copy number differences that were not assessed in this study. Further work may also

662    investigate the specific genomic variants that permit the building of accurate

663    classification models. Such outbreak associated variants may be diagnostic of specific

664    point sources and thus may be informative to understand bacterial genomic

665    responses to certain environmental reservoirs or public health control measures

666    (e.g., different decontamination or biocide practices).

667

668     Given the dynamic nature of bacterial populations, routine re-building of the models

669     with newly collected environmental isolates may be required to ensure accuracy as

670     emerging genomic signatures are then learned by the model. Another consideration

671     might be to limit the length of time in which genomes remain in the training

672     database, as older genomic signatures may no longer represent extant *L.*

673     *pneumophila* in environmental sources as time goes by. Here, a temporal sliding

674     window could be used, as has been implemented in other bacterial genomic

675     investigations (Gorrie et al., 2021).

676

677     **CONCLUSION:**

678     The advent of highly accessible bacterial genomics has provided a wealth of *L.*

679     *pneumophila* genomes in publicly assessable databases that are paired with

680     epidemiological information, of which provide the basis to build source attribution

681     classification approaches. Our development of an improved machine learning

682     classification technique now affords models with the ability to call true negatives,

683     offering the previously lacking negative classification capacity. Here we demonstrate

684     that our improved approach provides greater source tracking ability than two widely

685     used methods – phylogenomic trees and cgMLST allelic variation. Given the reported

686     high classification capacity of this improved approach, it is the vision of this work

687     that, soon, future LD public health investigations may make use of such modelling

688     advancements to rapidly pinpoint the correct environmental sources of *L.*

689     *pneumophila* and reduce the incidence of this preventable disease.

690

691     **ACKNOWLEDGEMENTS:**

697    **REFERENCES:**

698    Abrams, A. J., & Trees, D. L. (2017). Genomic sequencing of Neisseria
699        gonorrhoeae to respond to the urgent threat of antimicrobial-resistant
700        gonorrhea. *Pathogens and disease, 75*(4).
701    Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., . . .
702        Prjibelski, A. D. (2012). SPAdes: a new genome assembly algorithm and its
703        applications to single-cell sequencing. *Journal of computational biology,*
704        *19*(5), 455-477.
705    Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: a flexible trimmer for
706        Illumina sequence data. *Bioinformatics, 30*(15), 2114-2120.
707    Borchardt, J., Helbig, J., & Lück, P. (2008). Occurrence and distribution of
708        sequence types among Legionella pneumophila strains isolated from
709        patients in Germany: common features and differences to other regions of
710        the world. *European Journal of Clinical Microbiology & Infectious Diseases,*
711        *27*(1), 29-36.
712    Boser, B. E., Guyon, I. M., & Vapnik, V. N. (1992). *A training algorithm for optimal*
713        *margin classifiers.* Paper presented at the Proceedings of the fifth annual
714        workshop on Computational learning theory.
715    Breiman, L. (1996). Bagging predictors. *Machine learning, 24*(2), 123-140.
716    Breiman, L. (2001). Random forests. *Machine learning, 45*(1), 5-32.
717    Buultjens, A. H., Chua, K. Y., Baines, S. L., Kwong, J., Gao, W., Cutcher, Z., . . . Tomita,
718        T. (2017). A supervised statistical learning approach for accurate
719        Legionella pneumophila source attribution during outbreaks. *Applied and*
720        *environmental microbiology, 83*(21), e01482-01417.
721    Christiansen, M. T., Brown, A. C., Kundu, S., Tutill, H. J., Williams, R., Brown, J. R., . .
722        . Dave, J. (2014). Whole-genome enrichment and sequencing of Chlamydia
723        trachomatisdirectly from clinical samples. *BMC infectious diseases, 14*(1),
724        1-11.
725    David, S., Afshar, B., Mentasti, M., Ginevra, C., Podglajen, I., Harris, S. R., . . .
726        Parkhill, J. (2017). Seeding and establishment of Legionella pneumophila
727        in hospitals: implications for genomic investigations of nosocomial
728        Legionnaires' disease. *Clinical Infectious Diseases, 64*(9), 1251-1259.
729    David, S., Rusniok, C., Mentasti, M., Gomez-Valero, L., Harris, S. R., Lechat, P., . . .
730        Ma, L. (2016). Multiple major disease-associated clones of Legionella
731        pneumophila have emerged recently and independently. *Genome*
732        *research, 26*(11), 1555-1564.

733    Didelot, X., & Wilson, D. J. (2015). ClonalFrameML: efficient inference of
734        recombination in whole bacterial genomes. *PLoS computational biology,*
735        *11*(2), e1004041.
736    Fields, B. S., Benson, R. F., & Besser, R. E. (2002). Legionella and Legionnaires'
737        disease: 25 years of investigation. *Clin Microbiol Rev, 15*(3), 506-526.
738    Géron, A. (2019). *Hands-on machine learning with Scikit-Learn, Keras, and*
739        *TensorFlow: Concepts, tools, and techniques to build intelligent systems*: "
740        O'Reilly Media, Inc.".
741    Goldberg, B., Sichtig, H., Geyer, C., Ledeboer, N., & Weinstock, G. M. (2015).
742        Making the leap from research laboratory to clinic: challenges and
743        opportunities for next-generation sequencing in infectious disease
744        diagnostics. *MBio, 6*(6), e01888-01815.
745    Gorrie, C. L., Da Silva, A. G., Ingle, D. J., Higgs, C., Seemann, T., Stinear, T. P., . . .
746        Sherry, N. L. (2021). Key parameters for genomics-based real-time
747        detection and tracking of multidrug-resistant bacteria: a systematic
748        analysis. *The Lancet Microbe, 2*(11), e575-e583.
749    Gorzynski, J., Wee, B., Llano, M., Alves, J., Cameron, R., McMenamin, J., . . .
750        Fitzgerald, J. R. (2022). Epidemiological analysis of Legionnaires' disease
751        in Scotland: a genomic study. *The Lancet Microbe, 3*(11), e835-e845.
752    Graham, R., Doyle, C., & Jennison, A. (2014). Real-time investigation of a
753        Legionella pneumophila outbreak using whole genome sequencing.
754        *Epidemiology & Infection, 142*(11), 2347-2351.
755    Harris, S. R. (2018). SKA: Split kmer analysis toolkit for bacterial genomic
756        epidemiology. *bioRxiv*, 453142.
757    Harrison, T., Afshar, B., Doshi, N., Fry, N., & Lee, J. (2009). Distribution of
758        Legionella pneumophila serogroups, monoclonal antibody subgroups and
759        DNA sequence types in recent clinical and environmental isolates from
760        England and Wales (2000–2008). *European Journal of Clinical*
761        *Microbiology & Infectious Diseases, 28*(7), 781-791.
762    Ingle, D. J., Howden, B. P., & Duchene, S. (2021). Development of phylodynamic
763        methods for bacterial pathogens. *Trends in Microbiology, 29*(9), 788-797.
764    Krøvel, A. V., Bernhoff, E., Austerheim, E., Soma, M. A., Romstad, M. R., & Löhr, I. H.
765        (2022). Legionella pneumophila in Municipal Shower Systems in
766        Stavanger, Norway; A Longitudinal Surveillance Study Using Whole
767        Genome Sequencing in Risk Management. *Microorganisms, 10*(3), 536.
768    Kwong, J. C., Mercoulia, K., Tomita, T., Easton, M., Li, H. Y., Bulach, D. M., . . .
769        Howden, B. P. (2016). Prospective whole-genome sequencing enhances
770        national surveillance of Listeria monocytogenes. *Journal of clinical*
771        *microbiology, 54*(2), 333-342.
772    Lück, C., Fry, N. K., Helbig, J. H., Jarraud, S., & Harrison, T. G. (2013). Typing
773        methods for Legionella. In *Legionella* (pp. 119-148): Springer.
774    McAdam, P. R., Vander Broek, C. W., Lindsay, D. S., Ward, M. J., Hanson, M. F.,
775        Gillies, M., . . . Fitzgerald, J. R. (2014). Gene flow in environmental
776        Legionella pneumophila leads to genetic and pathogenic heterogeneity
777        within a Legionnaires' disease outbreak. *Genome biology, 15*(11), 1-10.
778    Mercante, J. W., & Winchell, J. M. (2015). Current and emerging Legionella
779        diagnostics for laboratory and outbreak investigations. *Clinical*
780        *microbiology reviews, 28*(1), 95-133.

781  Moran-Gilad, J., Prior, K., Yakunin, E., Harrison, T., Underwood, A., Lazarovitch, T.,
782          . . . Agmon, V. (2015). Design and application of a core genome multilocus
783          sequence typing scheme for investigation of Legionnaires' disease
784          incidents. *Eurosurveillance, 20*(28), 21186.
785  Paradis, E., Claude, J., & Strimmer, K. (2004). APE: analyses of phylogenetics and
786          evolution in R language. *Bioinformatics, 20*(2), 289-290.
787  Petzold, M., Prior, K., Moran-Gilad, J., Harmsen, D., & Lück, C. (2017).
788          Epidemiological information is key when interpreting whole genome
789          sequence data–lessons learned from a large Legionella pneumophila
790          outbreak in Warstein, Germany, 2013. *Eurosurveillance, 22*(45), 17-
791          00137.
792  Price, M. N., Dehal, P. S., & Arkin, A. P. (2009). FastTree: computing large
793          minimum evolution trees with profiles instead of a distance matrix.
794          *Molecular biology and evolution, 26*(7), 1641-1650.
795  Qin, T., Zhang, W., Liu, W., Zhou, H., Ren, H., Shao, Z., . . . Xu, J. (2016). Population
796          structure and minimum core genome typing of Legionella pneumophila.
797          *Scientific reports, 6*(1), 1-10.
798  Raphael, B. H., Baker, D. J., Nazarian, E., Lapierre, P., Bopp, D., Kozak-Muiznieks,
799          N. A., . . . Musser, K. A. (2016). Genomic resolution of outbreak-associated
800          Legionella pneumophila serogroup 1 isolates from New York State.
801          *Applied and environmental microbiology, 82*(12), 3582-3590.
802  Reller, L. B., Weinstein, M. P., & Murdoch, D. R. (2003). Diagnosis of Legionella
803          infection. *Clinical Infectious Diseases, 36*(1), 64-69.
804  Reuter, S., Harrison, T. G., Köser, C. U., Ellington, M. J., Smith, G. P., Parkhill, J., . . .
805          Török, M. E. (2013). A pilot study of rapid whole-genome sequencing for
806          the investigation of a Legionella outbreak. *BMJ open, 3*(1), e002175.
807  Ricci, M. L., Fillo, S., Ciammaruconi, A., Lista, F., Ginevra, C., Jarraud, S., . . . Lindsay,
808          D. (2022). Genome analysis of Legionella pneumophila ST23 from various
809          countries reveals highly similar strains. *Life science alliance, 5*(6).
810  Rousseau, C., Ginevra, C., Simac, L., Fiard, N., Vilhes, K., Ranc, A.-G., . . . Campese, C.
811          (2022). A Community Outbreak of Legionnaires' Disease with Two Strains
812          of L. pneumophila Serogroup 1 Linked to an Aquatic Therapy Centre.
813          *International Journal of Environmental Research and Public Health, 19*(3),
814          1119.
815  Sánchez-Busó, L., Guiral, S., Crespi, S., Moya, V., Camaró, M. L., Olmos, M. P., . . .
816          Vanaclocha, H. (2016). Genomic investigation of a legionellosis outbreak
817          in a persistently colonized hotel. *Frontiers in microbiology, 6*, 1556.
818  Schoonmaker-Bopp, D., Nazarian, E., Dziewulski, D., Clement, E., Baker, D. J.,
819          Dickinson, M. C., . . . Lapierre, P. (2021). Improvements to the Success of
820          Outbreak Investigations of Legionnaires' Disease: 40 Years of Testing and
821          Investigation in New York State. *Applied and environmental microbiology,*
822          *87*(16), e00580-00521.
823  Schwake, D. O., Garner, E., Strom, O. R., Pruden, A., & Edwards, M. A. (2016).
824          Legionella DNA markers in tap water coincident with a spike in
825          Legionnaires' disease in Flint, MI. *Environmental Science & Technology*
826          *Letters, 3*(9), 311-315.
827  Sintchenko, V., & Holmes, E. C. (2015). The role of pathogen genomics in
828          assessing disease transmission. *Bmj, 350*.

829    Wéry, N., Bru-Adan, V., Minervini, C., Delgénes, J.-P., Garrelly, L., & Godon, J.-J.
830            (2008). Dynamics of Legionella spp. and bacterial populations during the
831            proliferation of L. pneumophila in a cooling tower facility. *Applied and*
832            *environmental microbiology, 74*(10), 3030-3037.
833    Wüthrich, D., Gautsch, S., Spieler-Denz, R., Dubuis, O., Gaia, V., Moran-Gilad, J., . . .
834            Tschudin-Sutter, S. (2019). Air-conditioner cooling towers as complex
835            reservoirs and continuous source of Legionella pneumophila infection
836            evidenced by a genomic analysis study in 2017, Switzerland.
837            *Eurosurveillance, 24*(4), 1800192.
838    Yu, V. L., Plouffe, J. F., Pastoris, M. C., Stout, J. E., Schousboe, M., Widmer, A., . . .
839            Paterson, D. L. (2002). Distribution of Legionella species and serogroups
840            isolated by culture in patients with sporadic community-acquired
841            legionellosis: an international collaborative survey. *The Journal of*
842            *infectious diseases, 186*(1), 127-128.
843