1   # A supervised statistical learning approach for accurate

2   # *Legionella pneumophila* source attribution during outbreaks

3

4

5   Andrew H. Buultjens[1,2], Kyra Y. L. Chua[3], Sarah L. Baines[1], Jason Kwong[1,2,3], Wei Gao[1], Zoe Cutcher[4,5],

6   Stuart Adcock[4], Susan Ballard[3], Mark B. Schultz[3], Takehiro Tomita[3], Nela Subasinghe[3], Glen P. Carter[1,2],

7   Sacha J. Pidot[1], Lucinda Franklin[4], Torsten Seemann[3,6], Anders Gonçalves Da Silva[2,3], Benjamin P.

8   Howden[1,2,3,*], Timothy P. Stinear[1,2,*]

9

10   **Affiliations:**

11   [1]Department of Microbiology and Immunology at the Peter Doherty Institute for Infection and Immunity, The

12   University of Melbourne, Victoria, Australia

13   [2]Doherty Applied Microbial Genomics, The Peter Doherty Institute for Infection and Immunity, Victoria, Australia

14   [3]Microbiological Diagnostic Unit Public Health Laboratory at the Peter Doherty Institute for Infection and Immunity,

15   The University of Melbourne, Victoria, Australia.

16   [4]Health Protection Branch, Department of Health and Human Services, Victoria, Australia

17   [5]National Centre for Epidemiology and Population Health, Australian National University, Canberra, Australia

18   [6]Victorian Life Sciences Computational Initiative, The University of Melbourne, Victoria, Australia

19

20   **\*To whom correspondence may be addressed:** tstinear@unimelb.edu.au, bhowden@unimelb.edu.au

21

22   **Running title:** *L. pneumophila* source tracking using genomics

## 23  Abstract

24    Public health agencies are increasingly relying on genomics during Legionnaires' disease investigations.

25    However, the causative bacterium (*Legionella pneumophila*) has an unusual population structure with

26    extreme temporal and spatial genome sequence conservation. Furthermore, Legionnaires' disease

27    outbreaks can be caused by multiple *L. pneumophila* genotypes in a single source. These factors can

28    confound cluster identification using standard phylogenomic methods. Here, we show that a statistical

29    learning approach based on

30    *L. pneumophila* core genome single nucleotide polymorphism (SNP) comparisons eliminates ambiguity for

31    defining outbreak clusters and accurately predicts exposure sources for clinical cases.  We illustrate the

32    performance of our method by genome comparisons of 234 *L. pneumophila* isolates obtained from patients

33    and cooling towers in Melbourne, Australia between 1994 and 2014. This collection included one of the

34    largest reported Legionnaires' disease outbreaks, involving 125 cases at an aquarium. Using only sequence

35    data from *L. pneumophila* cooling tower isolates and including all core genome variation, we built a

36    multivariate model using discriminant analysis of principal components (DAPC) to find cooling tower-

37    specific genomic signatures, and then used it to predict the origin of clinical isolates. Model assignments

38    were 93% congruent with epidemiological data, including the aquarium Legionnaires' outbreak and three

39    other unrelated outbreak investigations.  We applied the same approach to a recently described

40    investigation of Legionnaires' disease within a UK hospital and observed model predictive ability of 86%.

41    We have developed a promising means to breach *L. pneumophila* genetic diversity extremes and provide

42    objective source attribution data for outbreak investigations.

43

## 44  Importance

45    Microbial outbreak investigations are moving to a paradigm where whole genome sequencing and

46    phylogenetic trees are used to support epidemiological investigations. It's critical that outbreak source

47    predictions are accurate, particularly for pathogens like *Legionella pneumophila*, which can spread widely

48    and rapidly via cooling system aerosols causing Legionnaires' disease. Here, by studying hundreds of

49    *Legionella pneumophila* genomes collected over 21 years around a major Australian city, we uncovered

50    limitations with the phylogenetic approach that could lead to misidentification of outbreak sources. We

51    implement instead a statistical learning technique that eliminates the ambiguity of inferring disease

52    transmission from phylogenies. Our approach takes geolocation information and core genome variation

53    from environmental *L. pneumophila* isolates to build statistical models that predict with high confidence

54    the environmental source of clinical *L. pneumophila* during disease outbreaks. We show the versatility of

55    the technique by applying it to unrelated Legionnaires' disease outbreaks in Australia and the UK.

56

## Introduction

58    Legionellae are Gram-negative bacteria that replicate within free-living aquatic amoebae and are present in

59    aquatic environments worldwide. These bacteria can proliferate in man-made water systems and cause

60    large outbreaks of pneumonia known as Legionnaires' disease when contaminated water is aerosolized and

61    inhaled (1). The majority of human infections are caused by *Legionella pneumophila* serogroup 1 (2). Public

62    health investigations of Legionnaires' disease outbreaks are typically supported by molecular typing

63    methods to establish the likely source of the bacteria and the extent of the outbreak. Investigations usually

64    proceed with the assumption that a single *Legionella* genotype is responsible for an environmental point

65    source reservoir (3). Traditional molecular typing methods described for fingerprinting Legionellae include

66    pulsed-field gel electrophoresis (PFGE) and sequence-based typing (SBT) (4). Increasingly, whole genome

67    sequencing (WGS) is being employed to investigate individual *Legionella* outbreaks and the insights

68    obtained from these high-resolution comparisons are challenging our expectations regarding common-

69    source outbreaks, which usually are characterized by a single strain or genotype (5-9). It is becoming

70    evident that outbreaks can be caused by multiple co-circulating *L. pneumophila* genotypes (5, 10) and that

71    *L. pneumophila* core genomes can be surprisingly conserved across space and time (8, 11-13).

72

73    Melbourne is in the state of Victoria and it is the second largest city in Australia with a population

74    approaching five million inhabitants, and considered the ninth largest city in the Southern Hemisphere.

75    Legionellosis has been a notifiable disease in Victoria since 1979 and there are 50-100 cases reported each

76    year, most occurring in the greater metropolitan region of Melbourne (14). The Microbiological Diagnostic

77    Unit Public Health Laboratory (MDU PHL) is Victoria's State Reference Laboratory for the characterization

78    and typing of *Legionella* spp. The laboratory's collection includes isolates from a particularly noteworthy

79    outbreak at the Melbourne Aquarium in April 2000. This was the largest single episode of Legionellosis

80    reported in Australia (15), approximately three months after the aquarium was opened to visitors, with

81    construction of the site completed in December 1999. It resulted in 125 confirmed cases, with positive

82    cultures obtained from 11 patients. Our isolate collection also spanned 28 other potential legionellosis

83    outbreaks or infection clusters, for which at least one culture isolate had been obtained.

84

85    In this study, we used comparative genomics to explore the population structure of 234 *Legionella*

86    *pneumophila* isolates recovered from human and environmental sources submitted to the MDU PHL in

87    Melbourne over a 21-year period. This collection included 11 clinical and 14 environmental isolates from

88    the Aquarium outbreak and 42 clinical and 50 environmental isolates from 28 other likely point source case

3

89    clusters. We also assessed genomic data from a recently described investigation of Legionnaires' cases at a

90    UK hospital (8). The aim of this project was to develop a robust genomic approach that would surmount the

91    unusual population structure of

92    *L. pneumophila* and assist identification of case clusters and source tracking efforts during Legionnaires'

93    disease outbreak investigations.

94

## Results

96    **Isolates and epidemiology.** There were 234 *Legionella pneumophila* serogroup 1 (Lpn-SG 1) isolates

97    obtained across a 21-year period between 1994 and 2014. Initial MLST analysis indicated that 180 isolates

98    (77%) belonged to ST30. The collection comprised 180 clinical isolates of respiratory origin (sputum or

99    bronchoscopy specimens) and 64 environmental isolates recovered from cooling tower water samples. All

100   isolates were collected in the state of Victoria with the exception of six isolates from patients who were

101   exposed elsewhere. Further information for each isolate is available in Table S1, including NCBI SRA

102   accession numbers. One hundred and ten of the 234 isolates were epidemiologically associated with 29

103   formally investigated case clusters or outbreaks, designated as outbreaks A-AC (Table S1). The majority of

104   these cases occurred within a 42 km radius of Melbourne city center and over a 16-year period. Outbreak

105   A, the Melbourne Aquarium outbreak, was the largest (15).

106

107   **Complete genome sequence of *Legionella pneumophila* serogroup 1 isolate Lpm7613.** Before this study,

108   there were no closed, fully assembled ST30 *L. pneumophila* genomes. Thus, to ensure identification of

109   maximum genetic variation among this dominant ST in our collection, we first established a ST30 reference

110   genome sequence, selecting a clinical isolate from the Melbourne Aquarium outbreak (Lpm7613). The

111   finished genome consisted of a single circular 3,261,562 bp chromosome (38.3% GC) and a 129,875 bp

112   circular plasmid (pLpm7613) (Fig. S1). Although the chromosome indicated this genome belonged to the

113   same lineage as *L. pneumophila* Philadelphia (Fig. 1A), the plasmid shared 100% nucleotide identity with

114   pLPP reported in *L. pneumophila* Paris, but 2kb shorter in length (16). A total of 2,891 chromosomal

115   protein-coding sequences (CDS), 43 tRNA genes and nine rRNA loci were predicted using Prokka (17).

116   CRISPR-Cas regions were not detected (18).

117

118   **Assessment of *L. pneumophila* population structure.** Sequence reads from the other 233 genomes and the

119   ten selected publicly available completed genomes were mapped to the chromosome of reference strain

120   Lpm7613. Approximately 90% of the Lpm7613 genome was present in all genomes (*i.e.*, core), with 188,049

4

121    variable core nucleotide positions identified. Population structure analyses using an unsupervised Bayesian

122    clustering approach revealed six distinct groups (BAPS groups) (Fig. 1A). Comparison of intra- and inter-

123    BAPS group pairwise SNP distances confirmed the validity of these clusters and highlighted the extensive

124    genetic variation among this Lpn-SG1 population (Fig. 1C). The exceptions were BAPS groups 3 and 4, which

125    classified isolates across two clades, and is likely explained by recombination. Most striking, however, was

126    the lack of diversity within the 186 genomes comprising BAPS group 5 (hereafter referred to as BAPS-5),

127    with a median core SNP distance of only 5 SNPs (IQR 3 – 7). Isolates dispersed across time and space

128    (including isolates from England, New South Wales, South Australia and Tasmania) were scattered

129    throughout the entire phylogeny. All 180 ST30 isolates were encompassed by BAPS-5, as was ST37 *L.*

130    *pneumophila* Philadelphia (Philadelphia, USA), ST211 *L. pneumophila* ATCC 43290 (Denver, USA) and ST733

131    *L. pneumophila* Thunder Bay (Ontario, Canada) (Fig. 1A,B). The median inter-BAPS group distances ranged

132    between 27,506 to 63,136 SNPs (Fig. 1C), highlighting that there is also substantial genetic diversity among

133    Lpn-SG1 isolates circulating in Melbourne.

134

135    A rooted maximum likelihood phylogeny of the population was then inferred using the 181,633 non-

136    recombining core SNP loci. The phylogenomic tree reflected the BAPS clusters with BAPS-5 forming a

137    distinct, well-supported lineage (Fig. 1A). The separation of the three North American reference isolates

138    from the Melbourne ST30 isolates is suggestive of contemporaneous global dispersal of this BAPS-5 lineage

139    (Fig. 1A,B). All BAPS groups displayed monophyletic origins with the exception of BAPS-3 and BAPS-4. BAPS-

140    3 had a single isolate of paraphyletic origin that shared a most recent common ancestor (MRCA) with BAPS-

141    2 while BAPS-4 contained two paraphyletic sub-clades, one of which shared a MRCA with the majority of

142    BAPS-3 isolates.

143

144    **Impact of recombination.** Recombination is a driving force in the evolution of the Legionellae (5, 7, 19-22).

145    Therefore, to further understand the structure and evolution of this Lpn-SG1 population we assessed the

146    impact of DNA exchange. There was evidence of extensive recombination among isolates across BAPS

147    groups 1-4, and 6 with approximately 3% of variable nucleotide sites impacted relative to the Lpm7613

148    reference chromosome. The detection of two paraphyletic groups (BAPS-3 and BAPS-4) is likely explained

149    by ancestral recombination among the component sub-clades. In comparison, there was little

150    recombination evident among BAPS-5 isolates (Fig. S2), in accord with the core SNP phylogeny described

151    above, and suggesting the relatively recent emergence of this *L. pneumophila* lineage. After removal of

152    putative sequences affected by recombination, tree branch lengths showed no correlation with isolation

153    dates ($r^2$=0.116). This observation indicates that nucleotide substitutions in the population have not been

154　evolving under a molecular clock model, thus limiting estimates for dates of emergence for particular

155　lineages.

156

157　**Genomic molecular epidemiology of local outbreaks.** We next compared only the 180 ST30 genomes to

158　our Lpm7613 reference genome, and again confirmed the very restricted genomic diversity within this

159　lineage (median core SNP distance was 6 SNPs (IQR 4 – 9), with five outlier genomes, impacted by

160　recombination. (Fig. 1B,C Fig. S2). Within this reconstructed core-genome ST30-specific phylogeny, many

161　but not all epidemiologically-related isolates formed distinct, well-supported, monophyletic clades. In some

162　instances, epidemiologically-associated isolates spanned multiple clades (outbreaks A, B, C, D and K) (Fig.

163　1B). In addition, Outbreak A (the Melbourne aquarium outbreak), which was previously considered to

164　represent infections caused by a single clone (Table S1) (15), actually contained five distinct genotypes (A1-

165　A5) (Fig 1B,C).

166

167　The analysis of environmental surveillance isolates provided an ideal means to gain insights into the

168　diversity within potential reservoirs of Lpn-SG1 - diversity that might enable prospective source tracking. A

169　phylogeographic analysis was therefore undertaken to assess the relationship between 64 environmental

170　Melbourne metropolitan isolates against their 11 cooling tower sampling locations. Based on variation in

171　core SNPs, striking geographical structure was observed, with the majority of isolates from common cooling

172　towers tightly clustering in the phylogeny (Fig. 2A). Comparisons of pairwise core SNPs depicted smaller

173　within group diversity and larger between location group diversity, further indicating the existence of

174　geographical population structure (Fig. 2B). This structure among the environmental Lpn-SG1 isolates

175　suggested it might be possible to use the genome data to build models predictive of environmental source

176　to assist epidemiological efforts during outbreak investigations.

177

178　**A multivariate statistical model for source attribution.** To enhance resolution and try to detect outbreak-

179　specific genomic signals, a supervised statistical learning approach called Discriminant Analysis of Principal

180　Components (DAPC) (23) was employed. DAPC is a linear discriminant analysis (LDA) that accommodates

181　discrete genetic-based predictors by first transforming the genetic data into continuous Principal

182　Components (PC) and building predictive classification models. The PCs are used to build discriminant

183　functions (DF) under the constraint that they must minimize within group variance, and maximize variance

184　between groups. Infection clusters were defined *a priori* from the epidemiological findings, and training

185　(environmental) isolates were used to establish the discriminant functions. The model was then be used to

186　estimate the posterior probability of membership for an unknown (*e.g.* clinical) isolate for each pre-

6

187    specified infection cluster given the training data. Here, we used 43 of the 64 environmental isolates in the

188    training set (cooling tower isolates originating from epidemiologically defined infection clusters that

189    possessed at least one environmental and clinical representative), under the assumption that each

190    outbreak was caused by exposure to a point source of Lpn-SG1. We used core genome SNPs from only

191    environmental Lpn-SG1 genomes to build the classifier (24).

192

193    Outbreak-associated environmental Lpn-SG1 were grouped *a priori* into training set groups based on the

194    origin of the cooling towers from which they were isolated (see model building details in methods). The DFs

195    were then used to classify 15 clinical isolates that had been independently assigned based on

196    epidemiological data to the training set groups, hereon referred to as the validation genomes (Table S1).

197    The input matrix for DAPC was an alignment of 714 non-recombinogenic SNPs variable among the 43

198    environmental genomes. Plots depicting the separation of isolates according to the first two discriminant

199    functions and the amount of variation explained is shown (Fig. 3). A model was trained using the first four

200    principal components (PC), as this was found to be optimal (see methods). We next classified our clinical

201    validation genomes using the model and found a 93% match between our model's assignment and that

202    proposed by the epidemiological data (Fig. 4A,B). These data show that despite the high level of genome

203    conservation and the presence of multiple genotypes within a single environmental source, it is possible to

204    utilize signature differences in core genome SNPs to build predictive probabilistic classification models. The

205    single discrepancy between model predictions and epidemiological groupings was an infection cluster C

206    genome that was predicted as originating from the Melbourne Aquarium. Interestingly, cluster C was

207    located closest to the Melbourne Aquarium at a distance of approximately 500 metres. Given the proximity

208    of clusters A and C, these data may indicate cooling towers were seeded from a common *L. pneumophila*

209    source. In order to appraise the utility of this method beyond a large urban setting and the ST30 genotype,

210    we built a sister model using 31 ST1 environmental *L. pneumophila* genomes from a previously published

211    hospital investigation in Essex, UK, and used it to predict the origins of seven nosocomial clinical isolates

212    (Fig. 3A, Table S1) (8). Here, the model was trained using an alignment of 59 non-recombinogenic SNPs

213    among the 31 environmental genomes and retaining the first 15 PCs, as this was found to be optimal. As

214    with the Melbourne disease clusters, the model performed very well. For 86% of the clinical isolates there

215    was a match between the model's ward assignment and the origin suggested by epidemiology (Fig. 4A,B).

216    Again, a single discrepancy occurred with a ward G genome predicted to originate from ward A. Wards A

217    and G were co-located on the same corner and level of a common building, again suggesting a common *L.*

218    *pneumophila* source (8). As before, isolates from a common source would be miss-assigned by the model,

219    owing to the lack of location-specific genomic variants.

220

221 **Core genome multilocus sequence typing (cgMLST) has reduced discrimination.** In order to evaluate the
222 utility of the recently described cgMLST scheme for source tracking (25, 26), we trained new DAPC models
223 for both the Melbourne and Essex hospital datasets using a matrix of allelic integers derived from SNP
224 profiles of the 1,529 cgMLST loci (Fig. 3B). When using the first one and seven PCs, we observed only 60%
225 and 71% concordance between our model's assignment and that predicted by the epidemiological data for
226 the Melbourne and Essex hospital datasets, respectively (Fig. 4A,C).

227

228 # Discussion

229 In this study, we have retrospectively examined a large collection of 234 clinical and environmental isolates
230 Lpn-SG1 isolates spanning 29 defined outbreaks. Isolates were collected over wide temporal and spatial
231 scales and detailed genomic comparisons revealed wide extremes Lpn-SG1 genetic diversity among distinct
232 genomic populations; a phenomenon not fully appreciated from previous genomic investigations that have
233 sampled less extensively and focused on single outbreaks (5, 6, 27). Most striking in our collection was the
234 high sequence conservation and dominance of a single genotype (BAPS-5, ST30), shared by 77% of isolates
235 with a median core SNP distance of only 5 SNPs across 21 years. In agreement with our findings, two recent
236 population genomic investigations of Lpn-SG1 also describe the unusual restriction in core genome diversity
237 (8, 12).

238

239 Based on our previous experience with other bacterial pathogens (28) and reports in the literature of
240 Legionnaires' disease outbreak investigations using genomics (27, 29) - we expected to be able to use Lpn-
241 SG1 genomic comparisons and develop genetic rule-in or rule-out criteria to guide outbreak assessment
242 and source attribution. For example, we recently proposed a 'traffic-light' system for *Listeria*
243 *monocytogenes* based on SNP difference cutoffs of 'likely related', 'possibly-related' and 'not-related' (28).
244 This approach has also been proposed for *L. pneumophila* (25). A comparison of genotyping approaches
245 using 335 Lpn isolates, including 106 from the European Society for Clinical Microbiology Study Group's
246 *Legionella Typing Panel*, proposed an escalating, hierarchical approach to genotyping, beginning with an
247 extended 50-gene MLST scheme up to a 1529-gene cgMLST (25, 26).

248

249 The analysis of the population structure of Lpn-SG1 presented here indicates that SNP-based typing with
250 threshold cut-offs, whether they are based on seven genes, 50 genes, 1500 genes or whole genomes will
251 not necessarily provide sufficient discriminative power. These genotyping approaches will be confounded
252 by the presence of (i) indistinguishable Lpn-SG1 genotypes present in unrelated cases and (ii) polyclonal

253     outbreaks. Our retrospective analysis of the Melbourne Aquarium outbreak illustrates clearly both these

254     issues, where five distinct subtypes were recovered from 25 clinical and environmental isolates (Fig. 1B).

255     There is a growing awareness of single source, polyphyletic Lpn-SG1 outbreaks (8, 10, 13, 30, 31). These

256     data all point to the need for a different approach in order to use molecular epidemiology and genomics in

257     support of *Legionella* outbreak investigations.

258

259     We address this issue by exploiting all core genome information to train probabilistic classification models.

260     Our DAPC analysis demonstrates that it is possible to build predictive models based on Lpn-SG1

261     environmentally derived genomes that help in identifying the source of clinical isolates during complex

262     outbreak investigations in both the community and hospital environments (Fig. 4). By including all core SNP

263     variation, DAPC was able to identify outbreak-specific genotypes, even when the source of the outbreak

264     was polyclonal. This enabled us to build robust models that assigned validation set genomes, with known

265     provenance, back to their original groupings with high concordance. The fact that this model was built

266     purely from environmental surveillance isolates demonstrates that such approaches can be developed

267     prospectively and be preexisting, ready to deploy at the onset of outbreaks.

268

269     In contrast to the high performance of the DAPC model developed from core genome SNPs, the model built

270     using variants identified by cgMLST scheme had a lower matching rate when assigning validation genomes

271     back to their putative epidemiological groupings (Fig. 4C). Despite cgMLST being a useful tool for broad

272     Lpn-SG1 population structure assessment, our analysis suggests it may have insufficient resolution and thus

273     predictive capacity for outbreak investigations.

274

275     The DAPC approach however, while promising, does not permit discrimination among isolates that do not

276     belong to defined clusters. This is because the model assumes that the world is composed of only the *k*

277     groups used to train it, and therefore assigns unknown isolates to one of these groups, even if the isolate is

278     known not to be part of any of the groups. One way to address this issue would be to create a single group

279     classifier, which is trained with environmental samples. Isolates with low probability of membership to this

280     single large group would then be excluded before being analyzed with the multi-group model.  Future

281     models could be further improved by adding epidemiological evidence (*e.g.* patient zip codes), and assess

282     how that improves our assignment of a clinical isolate to a particular location. An advantage of a

283     classification-based model is that its output could be distilled down to a zip code (or group of zip codes) and

284     a probability that a clinical isolate is associated with the zip code (indicating uncertainty about the

285     classification). This would obviate the need to interpret, and explain phylogenetic trees. Interpreting trees

9

286  is often not intuitive and trees may fail to communicate what action is required from a public health

287  perspective. Crucial for such a classification approach to work however, is an extensive and temporally

288  dynamic database of environmental Lpn-SG1 genotypes. That is, there would need to be ongoing

289  surveillance and isolation of Lpn-SG1 from environmental sources. We are currently investigating how to

290  implement such models.

291

292  The modeling approach, is not intended to be used in isolation, but rather employed as an adjunct to

293  traditional epidemiological investigations. In this way, insights gained through epidemiological

294  investigations can be informed by microbiological evidence from our predictive models. A limitation of our

295  current models are the relatively small sample sizes. Performance measures for validation sets this small

296  are often sensitive to slight perturbations in the data and may be influenced by small features of the data.

297  However, as a proof-of-concept implementation of our approach, we have built two models from

298  independent datasets, and both demonstrate high predictive capacity. More robust appraisals of model

299  performance will require validation with larger datasets, collected prospectively.

300

301  From a biological perspective, the lack of genetic diversity in Lpn-SG1 over such coarse temporal and spatial

302  scales is potentially explained by a reservoir of latent-state bacteria intermittently seeding warm water

303  sources in the greater Melbourne region and is supported by the frequently-reported and widespread

304  presence of *Legionella* species in drinking water supply systems (DWSS) (32-34). Independent studies

305  propose similar hypotheses to explain the surprisingly high sequence conservation among some *L.*

306  *pneumophila* genomes (8, 12).

307

308  This study is, to our knowledge, the largest genomic investigation of environmental and clinical Legionella

309  reported to date from a single jurisdiction and confirms that Lpn-SG1 is an unusual 'edge case' in the

310  application of genomics in public health microbiology. In the absence of a deep understanding of local *L.*

311  *pneumophila* population structure (both clinical and environmental) the combination of extreme genomic

312  monomorphism combined with outbreaks caused by mixed pathogen populations could easily lead to

313  erroneous conclusions regarding source attribution. Thus, we require new approaches that can better

314  utilize the genomic information available, and harmoniously combine it with epidemiological evidence, in

315  order to provide public health officials with useful and timely information.

316

317  **Materials and Methods**

318     **Bacterial strains, growth conditions, case definitions.** *Legionella pneumophila* serogroup 1 isolates were

319     resuscitated from -80°C storage and assessed. Duplicate isolates from the same patient were excluded

320     from the study. Isolates were cultured for 48-72 h at 37°C on BCYE agar and re-confirmed serogroup 1 by

321     latex agglutination (Oxoid). Metadata collected on all isolates included year of isolation and country or city

322     of isolation. Cases resident in the state of Victoria, Australia, were assessed by the Victorian State

323     Government public health unit in accordance with national guidelines and an outbreak investigation was

324     initiated when common exposures were reported by different cases whose onset dates occurred within a

325     two-week window. (http://www.health.gov.au/internet/main/publishing.nsf/content/cdna-song-

326     legionella.htm, accessed 31 August 2015). In this manner, we were able to determine the human cases

327     epidemiologically linked to each other. Many of the outbreaks/infection clusters contained a greater

328     number of cases than there were isolates as the diagnosis of Legionellosis was made by culture-

329     independent methods. Complete, closed genomes of *L. pneumophila* that were publicly available were

330     obtained from GenBank for inclusion in the analysis (Table S1).

331

332     **Sequence based typing.** This was performed as previously described according to the European

333     Legionnaires' Disease Surveillance Network (ELDSNet) method

334     (http://bioinformatics.phe.org.uk/legionella/legionella_sbt/php/sbt_homepage.php, accessed 31 August

335     2015) (35).

336

337     **DNA sequencing.** DNA libraries were prepared using the NexteraXT DNA preparation kit (Illumina) and

338     whole genome sequencing was performed on the NextSeq platform (Illumina) with 2x150 bp chemistry. For

339     single molecule real-time (SMRT) sequencing (Pacific Biosciences), genomic DNA was extracted from

340     agarose plugs using the CDC Pulsenet Protocol to allow for recovery of high molecular weight, intact DNA

341     (http://www.cdc.gov/pulsenet/pathogens, accessed 31 August 2015). Size-selected 10kb DNA libraries

342     were prepared according to manufacturers' instructions and sequenced on the RS II platform (Pacific

343     Biosciences) using P6-C4 chemistry. All sequence reads and the completed genome are available (GenBank

344     BioProject ID: PRJEB13594)

345

346     *Legionella pneumophila* **serogroup 1 isolate Lpm7613 assembly and closure.** A high quality finished ST 30

347     reference genome was established for *L. pneumophila* serogroup 1 clinical isolate Lpm7613 using the

348     SMRT® Analysis System v2.3.0.140936 (Pacific Biosciences). Raw sequence data were *de novo* assembled

349     using the HGAP v3 protocol with a genome size of 4 Mb. Polished contigs were error corrected using Quiver

350     v1. The resulting assembly was then checked using BridgeMapper v1 in the SMRT® Analysis System, and the

11

351     consensus sequence corrected with short-read Illumina data, using the program Snippy

352     (https://github.com/tseemann/snippy). Whole genome annotation was performed using Prokka (17),

353     preferentially using the *L. pneumophila* Paris strain annotation (16). BRIG was used to visualize BLASTn

354     DNA:DNA comparisons of *L. pneumophila* Lpm7613 against other *L. pneumophila* genomes (36).

355     Nomenclature of the genomic islands demonstrated in *L. pneumophila* Lpm7613 was based on previously

356     described islands (37). CRISPR databases were used to search for CRISPR sequences

357     (http://crispi.genouest.org and http://crispr.u-psud.fr/Server/, accessed 14 February 2016).

358

359     **Variant detection and phylogenetic analysis.** The genomes of ten publicly available complete *L.*

360     *pneumophila* genomes (Table S1) were shredded to generate short *in silico* sequence reads of 250bp and all

361     244 *L. pneumophila* reads sets were mapped against the Lpm7613 reference genome using Snippy v3.2. An

362     alignment file from pairwise comparisons of core genome SNPs (with inferred recombining sites removed)

363     was used as input to FastTree v2.1.8 with double precision (38) to infer a maximum likelihood phylogenetic

364     tree using the general time reversible model of nucleotide substitution. Branch support was estimated

365     using 1,000 bootstrap replicates. Resulting trees were visualized in FigTree v1.4.2

366     (http://tree.bio.ed.ac.uk/software/figtree/). Single nucleotide polymorphism (SNP) differences between

367     isolates were tabulated and visualized using a custom R-script (https://github.com/MDU-

368     PHL/pairwise_snp_differences). The core genome SNPs were also used as the input into a Bayesian analysis

369     of population structure (BAPS) using iterative clustering to a depth of 10 levels and a pre-specified

370     maximum of 20 clusters (39).

371

372     **Recombination and molecular clock analysis.** Recombination detection was performed using

373     ClonalFrameML (40), taking as input a full genome alignment (included invariant sites) prepared using

374     Snippy as above and the ML phylogeny as a guide tree with polytomies removed from the FastTree tree

375     using a custom python script (https://github.com/kwongj/nw_multi2bifurcation). Results were visualized

376     using a custom Python script to render separate and superposable images of extant and ancestral inferred

377     recombination regions (https://github.com/kwongj/cfml-maskrc). Molecular clock-likeness of the ML tree

378     with ClonalFrameML-adjusted branch lengths was assessed using TempEst v1.5

379     (http://tree.bio.ed.ac.uk/software/tempest/).

380

381     **Phylogeographic analysis.** Variant detection for the 64 environmental genomes was undertaken by running

382     snippy-core. Core SNPs were used to reconstruct a phylogenomic tree with FastTree that was overlaid upon

383     a base map in GenGIS (41). Victorian population mesh data was downloaded from the Australia Bureau of

384    Statistics webpage

385    (http://www.abs.gov.au/AUSSTATS/abs@.nsf/DetailsPage/1270.0.55.001July%202016?OpenDocument)

386    and Local Government Area data was downloaded from the Victorian Government Data Directory webpage

387    (https://www.data.vic.gov.au/data/dataset/lga-geographical-profiles-2014-beta/resource/f6c49074-0679-

388    4c79-a0db-04dac8eda364).

389

390    **DAPC model building using core SNPs.** Discriminant analysis of principal components (DAPC) is a

391    multivariate method that tries to reconstruct hypothesized subdivisions in a given population (typically

392    formed from demographic or phenotypic information) using genomic data (42). DAPC was implemented in

393    the R package *adegenet* v2.0.1 (42). For input, we used a matrix of single nucleotide polymorphisms (SNP)

394    for all genomes originating from infection clusters that possessed at least one environmental and clinical

395    representative (Table S1). SNP detection was undertaken by running Snippy and sites that were

396    recombinogenic and or invariant among the environmental genomes were discarded. An input SNP matrix

397    of exclusively environmental isolates (hereon referred to as the training set) was used to develop a DAPC

398    model. The training set subdivisions were based on the geographic origin of the environmental isolates

399    (Table S1) (23). The resultant model was then tested using clinical isolates (hereon referred to as the

400    validation set). The ability of the model to predict the environmental source of the validation set genomes

401    was simulated across the first to the 20th principal components, allowing an optimal number of principal

402    components to be identified. The optimized model was then used to predict the environmental origin of

403    the clinical isolate genomes.

404

405    **DAPC model building using cgMLST variation.** In order to detect variants within the recently described

406    cgMLST regions, reads were mapped to the Lp_Philadelphia chromosome (NC_002942.5) using snippy. SNP

407    profiles from within the cgMLST regions were reduced to allelic integers, with all genes containing zero

408    coverage or uncertain base-calls, excluded. Allelic integers were concatenated into a matrix and, using the

409    same DAPC model-building method as mentioned above, models were established using the training set

410    environmental genomes and used to predict the origin of the validation set clinical isolate genomes.
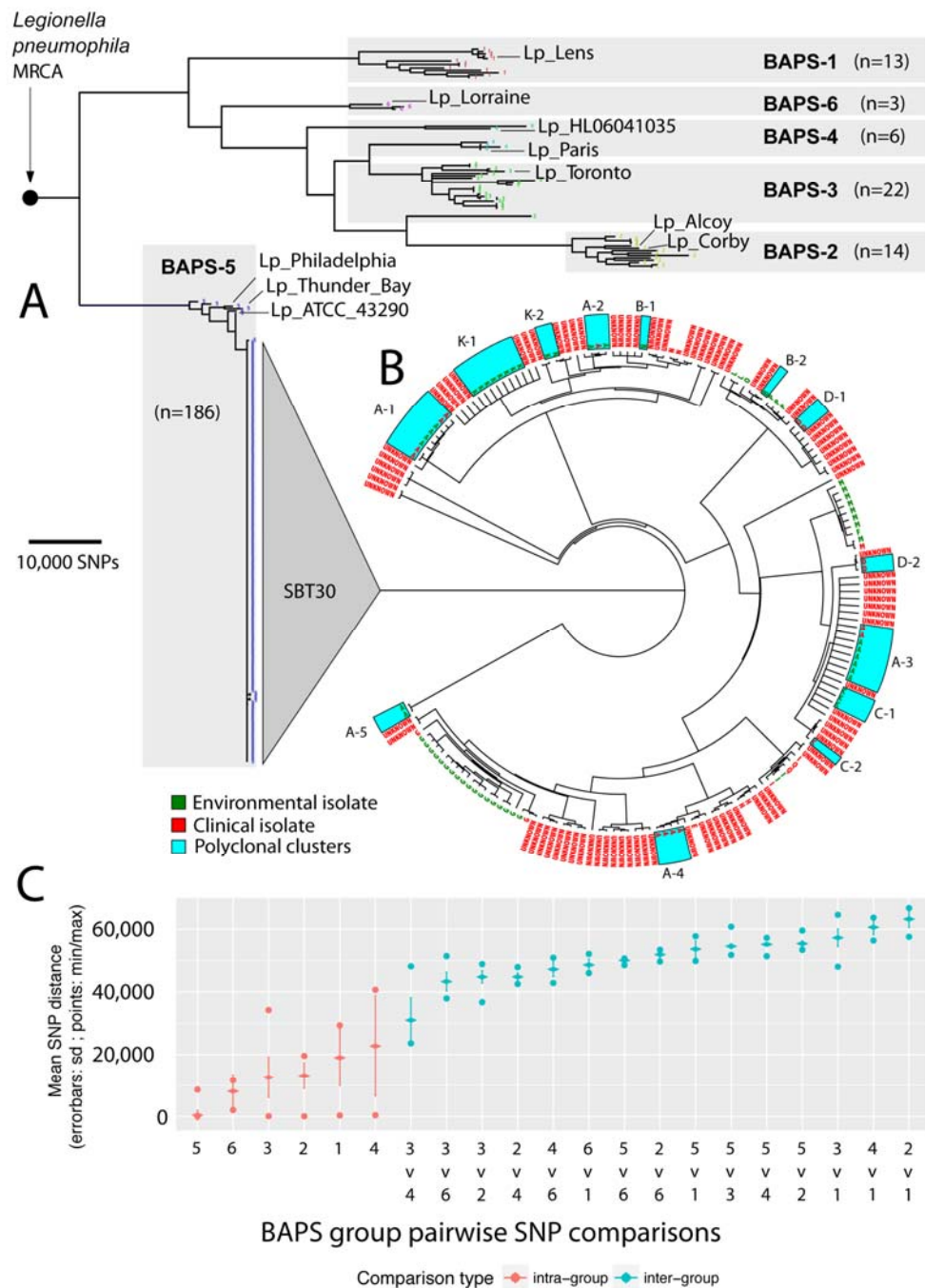
411

412    # References

413    1.      **Fields BS, Benson RF, Besser RE.** 2002. *Legionella* and Legionnaires' disease: 25 years of

414            investigation. Clin Microbiol Rev **15:**506-526.

415  2.  Yu VL, Plouffe JF, Pastoris MC, Stout JE, Schousboe M, Widmer A, Summersgill J, File T, Heath CM,
416      Paterson DL, Chereshsky A. 2002. Distribution of *Legionella* species and serogroups isolated by
417      culture in patients with sporadic community-acquired legionellosis: an international collaborative
418      survey. J Infect Dis **186**:127-128.

419  3.  Luck C, Fry NK, Helbig JH, Jarraud S, Harrison TG. 2013. Typing methods for *Legionella*. Methods
420      Mol Biol **954**:119-148.

421  4.  Mercante JW, Winchell JM. 2015. Current and emerging *Legionella* diagnostics for laboratory and
422      outbreak investigations. Clin Microbiol Rev **28**:95-133.

423  5.  McAdam PR, Vander Broek CW, Lindsay DS, Ward MJ, Hanson MF, Gillies M, Watson M, Stevens
424      JM, Edwards GF, Fitzgerald JR. 2014. Gene flow in environmental *Legionella pneumophila* leads to
425      genetic and pathogenic heterogeneity within a Legionnaires' disease outbreak. Genome Biol
426      **15**:504.

427  6.  Reuter S, Harrison TG, Koser CU, Ellington MJ, Smith GP, Parkhill J, Peacock SJ, Bentley SD, Torok
428      ME. 2013. A pilot study of rapid whole-genome sequencing for the investigation of a *Legionella*
429      outbreak. BMJ Open **3**.

430  7.  Sanchez-Buso L, Comas I, Jorques G, Gonzalez-Candelas F. 2014. Recombination drives genome
431      evolution in outbreak-related *Legionella pneumophila* isolates. Nat Genet **46**:1205-1211.

432  8.  David S, Afshar B, Mentasti M, Ginevra C, Podglajen I, Harris SR, Chalker VJ, Jarraud S, Harrison
433      TG, Parkhill J. 2017. Seeding and Establishment of Legionella pneumophila in Hospitals:
434      Implications for Genomic Investigations of Nosocomial Legionnaires' Disease. Clin Infect Dis
435      **64**:1251-1259.

436  9.  Weiss D, Boyd C, Rakeman JL, Greene SK, Fitzhenry R, McProud T, Musser K, Huang L, Kornblum J,
437      Nazarian EJ, Fine AD, Braunstein SL, Kass D, Landman K, Lapierre P, Hughes S, Tran A, Taylor J,
438      Baker D, Jones L, Kornstein L, Liu B, Perez R, Lucero DE, Peterson E, Benowitz I, Lee KF, Ngai S,
439      Stripling M, Varma JK, South Bronx Legionnaires' Disease Investigation T. 2017. A Large
440      Community Outbreak of Legionnaires' Disease Associated With a Cooling Tower in New York City,
441      2015. Public Health Rep **132**:241-250.

442  10. Sanchez-Buso L, Guiral S, Crespi S, Moya V, Camaro ML, Olmos MP, Adrian F, Morera V, Gonzalez-
443      Moran F, Vanaclocha H, Gonzalez-Candelas F. 2015. Genomic Investigation of a Legionellosis
444      Outbreak in a Persistently Colonized Hotel. Front Microbiol **6**:1556.

445 11. Underwood AP, Jones G, Mentasti M, Fry NK, Harrison TG. 2013. Comparison of the *Legionella*
446 *pneumophila* population structure as determined by sequence-based typing and whole genome
447 sequencing. BMC Microbiol **13**:302.

448 12. David S, Rusniok C, Mentasti M, Gomez-Valero L, Harris SR, Lechat P, Lees, J., Ginerva C, Glaser C,
449 Ma L, Bouchier C, Underwood A, Jarraud S, Harrison TG, Parkhill J, Buchrieser C. 2016. Multiple
450 major disease-associated clones of *Legionella pneumophila* have emerged recently and
451 independently. Genome Res **In press**.

452 13. Schjorring S, Stegger M, Kjelso C, Lilje B, Bangsborg JM, Petersen RF, David S, Uldum SA,
453 Infections ESGfL. 2017. Genomic investigation of a suspected outbreak of *Legionella pneumophila*
454 ST82 reveals undetected heterogeneity by the present gold-standard methods, Denmark, July to
455 November 2014. Euro Surveill **22**.

456 14. Anon. 2014. Communicable Disease Surveillance, Victoria, Oct-Dec 2014. Victorian Infectious
457 Diseases Bulletin **17**:22-23.

458 15. Greig JE, Carnie JA, Tallis GF, Ryan NJ, Tan AG, Gordon IR, Zwolak B, Leydon JA, Guest CS, Hart
459 WG. 2004. An outbreak of Legionnaires' disease at the Melbourne Aquarium, April 2000:
460 investigation and case-control studies. Med J Aust **180**:566-572.

461 16. Cazalet C, Rusniok C, Bruggemann H, Zidane N, Magnier A, Ma L, Tichit M, Jarraud S, Bouchier C,
462 Vandenesch F, Kunst F, Etienne J, Glaser P, Buchrieser C. 2004. Evidence in the *Legionella*
463 *pneumophila* genome for exploitation of host cell functions and high genome plasticity. Nat Genet
464 **36**:1165-1173.

465 17. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. Bioinformatics **30**:2068-2069.

466 18. Rao C, Guyard C, Pelaz C, Wasserscheid J, Bondy-Denomy J, Dewar K, Ensminger AW. 2016. Active
467 and Adaptive *Legionella* CRISPR-Cas reveals a recurrent challenge to the pathogen. Cell Microbiol
468 doi:10.1111/cmi.12586.

469 19. Coscolla M, Gonzalez-Candelas F. 2007. Population structure and recombination in environmental
470 isolates of *Legionella pneumophila*. Environ Microbiol **9**:643-656.

471 20. Coscolla M, Gonzalez-Candelas F. 2009. Comparison of clinical and environmental samples of
472 *Legionella pneumophila* at the nucleotide sequence level. Infect Genet Evol **9**:882-888.

473 21. Gomez-Valero L, Rusniok C, Jarraud S, Vacherie B, Rouy Z, Barbe V, Medigue C, Etienne J,
474 Buchrieser C. 2011. Extensive recombination events and horizontal gene transfer shaped the
475 *Legionella pneumophila* genomes. BMC Genomics **12**:536.

15

476 22. **Costa J, Teixeira PG, d'Avo AF, Junior CS, Verissimo A.** 2014. Intragenic recombination has a critical
477       role on the evolution of *Legionella pneumophila* virulence-related effector sidJ. PLoS One
478       **9**:e109840.

479 23. **Jombart T, Devillard S, Balloux F.** 2010. Discriminant analysis of principal components: a new
480       method for the analysis of genetically structured populations. BMC Genet **11**:94.

481 24. **McNally A, Oren Y, Kelly D, Pascoe B, Dunn S, Sreecharan T, Vehkala M, Valimaki N, Prentice MB,**
482       **Ashour A, Avram O, Pupko T, Dobrindt U, Literak I, Guenther S, Schaufler K, Wieler LH, Zhiyong Z,**
483       **Sheppard SK, McInerney JO, Corander J.** 2016. Combined Analysis of Variation in Core, Accessory
484       and Regulatory Genome Regions Provides a Super-Resolution View into the Evolution of Bacterial
485       Populations. PLoS Genet **12**:e1006280.

486 25. **David S, Mentasti M, Tewolde R, Aslett M, Harris SR, Afshar B, Underwood A, Fry NK, Parkhill J,**
487       **Harrison TG.** 2016. Evaluation of an optimal epidemiologic typing scheme for *Legionella*
488       *pneumophila* with whole genome sequence data using validation guidelines. J Clin Microbiol
489       doi:10.1128/JCM.00432-16.

490 26. **Moran-Gilad J, Prior K, Yakunin E, Harrison TG, Underwood A, Lazarovitch T, Valinsky L, Luck C,**
491       **Krux F, Agmon V, Grotto I, Harmsen D.** 2015. Design and application of a core genome multilocus
492       sequence typing scheme for investigation of Legionnaires' disease incidents. Euro Surveill **20**.

493 27. **Bartley PB, Ben Zakour NL, Stanton-Cook M, Muguli R, Prado L, Garnys V, Taylor K, Barnett TC,**
494       **Pinna G, Robson J, Paterson DL, Walker MJ, Schembri MA, Beatson SA.** 2016. Hospital-wide
495       Eradication of a Nosocomial *Legionella pneumophila* Serogroup 1 Outbreak. Clin Infect Dis **62**:273-
496       279.

497 28. **Kwong JC, Mercoulia K, Tomita T, Easton M, Li HY, Bulach DM, Stinear TP, Seemann T, Howden**
498       **BP.** 2016. Prospective Whole-Genome Sequencing Enhances National Surveillance of *Listeria*
499       *monocytogenes*. J Clin Microbiol **54**:333-342.

500 29. **Graham RM, Doyle CJ, Jennison AV.** 2014. Real-time investigation of a *Legionella pneumophila*
501       outbreak using whole genome sequencing. Epidemiol Infect **142**:2347-2351.

502 30. **Sanchez-Buso L, Olmos MP, Camaro ML, Adrian F, Calafat JM, Gonzalez-Candelas F.** 2015.
503       Phylogenetic analysis of environmental *Legionella pneumophila* isolates from an endemic area
504       (Alcoy, Spain). Infect Genet Evol **30**:45-54.

505 31. **Coscolla M, Fernandez C, Colomina J, Sanchez-Buso L, Gonzalez-Candelas F.** 2014. Mixed infection
506       by *Legionella pneumophila* in outbreak patients. Int J Med Microbiol **304**:307-313.

507    32.    **Diederen BM, de Jong CM, Aarts I, Peeters MF, van der Zee A.** 2007. Molecular evidence for the
508            ubiquitous presence of *Legionella* species in Dutch tap water installations. J Water Health **5:**375-
509            383.

510    33.    **Rivera JM, Aguilar L, Granizo JJ, Vos-Arenilla A, Gimenez MJ, Aguiar JM, Prieto J.** 2007. Isolation of
511            *Legionella* species/serogroups from water cooling systems compared with potable water systems in
512            Spanish healthcare facilities. J Hosp Infect **67:**360-366.

513    34.    **Storey MV, Langmark J, Ashbolt NJ, Stenstrom TA.** 2004. The fate of legionellae within distribution
514            pipe biofilms: measurement of their persistence, inactivation and detachment. Water Sci Technol
515            **49:**269-275.

516    35.    **Gaia V, Fry NK, Afshar B, Luck PC, Meugnier H, Etienne J, Peduzzi R, Harrison TG.** 2005. Consensus
517            sequence-based scheme for epidemiological typing of clinical and environmental isolates of
518            *Legionella pneumophila*. J Clin Microbiol **43:**2047-2052.

519    36.    **Alikhan NF, Petty NK, Ben Zakour NL, Beatson SA.** 2011. BLAST Ring Image Generator (BRIG):
520            simple prokaryote genome comparisons. BMC Genomics **12:**402.

521    37.    **D'Auria G, Jimenez-Hernandez N, Peris-Bondia F, Moya A, Latorre A.** 2010. *Legionella*
522            *pneumophila* pangenome reveals strain-specific virulence factors. BMC Genomics **11:**181.

523    38.    **Price MN, Dehal PS, Arkin AP.** 2010. FastTree 2--approximately maximum-likelihood trees for large
524            alignments. PLoS One **5:**e9490.

525    39.    **Cheng L, Connor TR, Siren J, Aanensen DM, Corander J.** 2013. Hierarchical and spatially explicit
526            clustering of DNA sequences with BAPS software. Mol Biol Evol **30:**1224-1228.

527    40.    **Didelot X, Wilson DJ.** 2015. ClonalFrameML: efficient inference of recombination in whole bacterial
528            genomes. PLoS Comput Biol **11:**e1004041.

529    41.    **Parks DH, Mankowski T, Zangooei S, Porter MS, Armanini DG, Baird DJ, Langille MG, Beiko RG.**
530            2013. GenGIS 2: geospatial analysis of traditional and genetic biodiversity, with new gradient
531            algorithms and an extensible plugin framework. PLoS One **8:**e69885.

532    42.    **Jombart T, Ahmed I.** 2011. adegenet 1.3-1: new tools for the analysis of genome-wide SNP data.
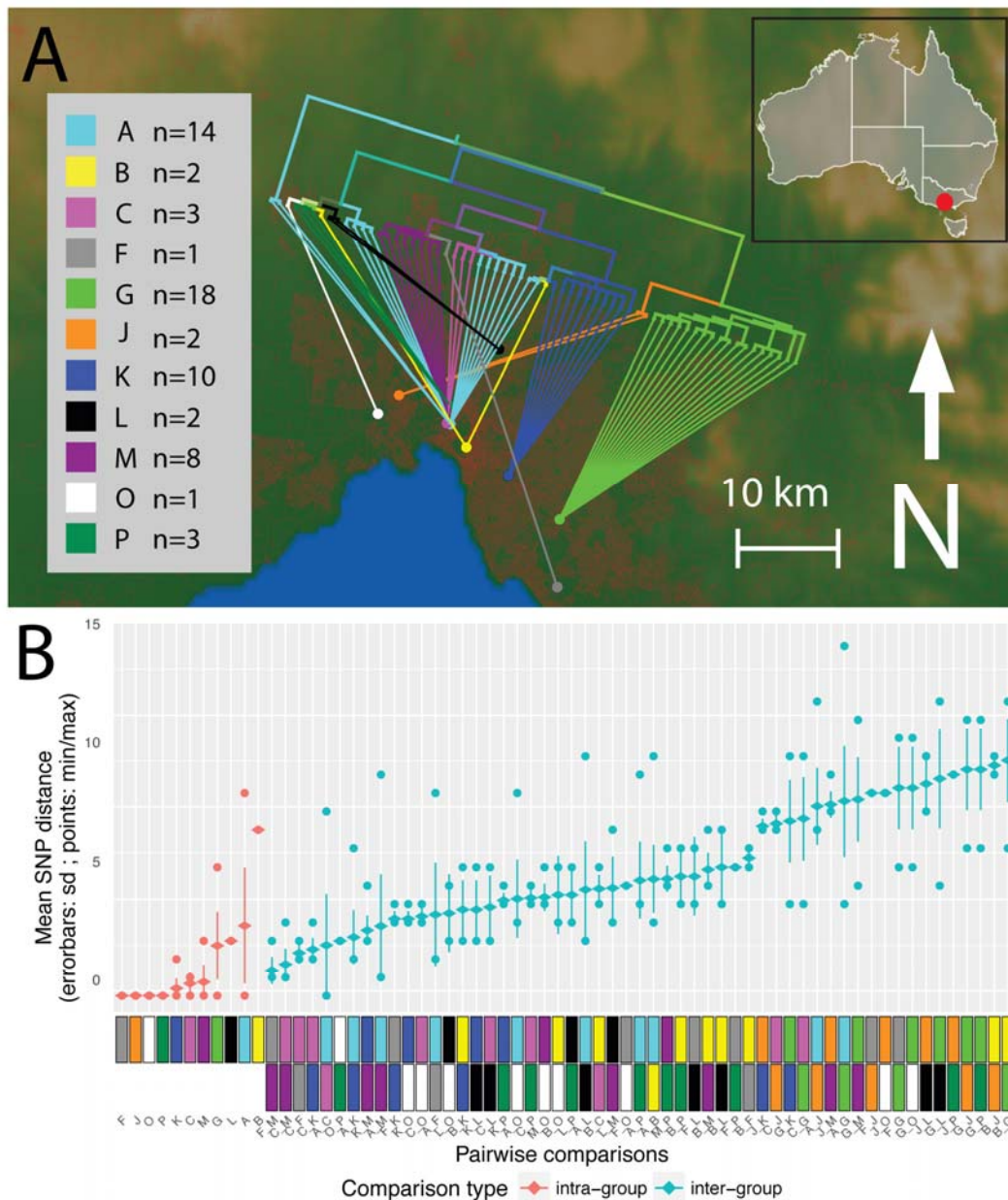533            Bioinformatics **27:**3070-3071.

534

**Figures**

**Fig. 1**: Global *Legionella pneumophila* population clustering, phylogenomics and genomic molecular epidemiology of local outbreaks. (A) Core genome phylogeny estimated using maximum likelihood corresponds with six BAPS groups. Branches with less than 70% bootstrap support were collapsed and scale indicates the number of core SNPs. The locations of the ten international genomes are labeled. (B) ST30

18

541    core genome phylogeny. Tree tips are labeled with outbreak codes. Environmental and clinical isolates are

542    colored according to the key. Polyclonal outbreaks/case clusters are highlighted with blue boxes. Branch

543    lengths have been transformed and are proportional to the number of nodes under each parent node. (C)

544    Core genome pairwise SNP comparisons of within and between BAPS groups. All groups had smaller within

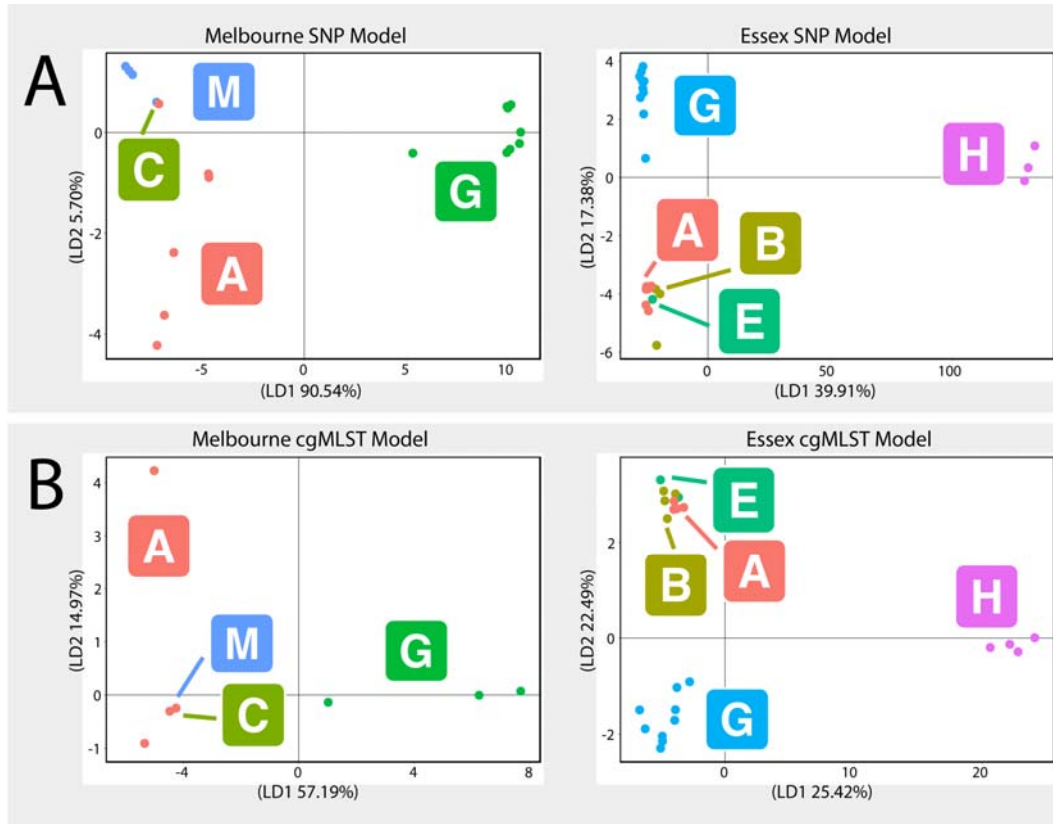545    group diversity compared to comparisons between groups.

546

547

**Fig. 2**: **Phylogeography of 64 Lpn-SG1 environmental isolate genomes.** (A) Map of the greater Melbourne area, showing the location of the 11 cooling towers assessed during Legionellosis outbreaks, designated by colored circles. 'A' (light blue) represents the location Melbourne aquarium outbreak and is close to the centre of Melbourne. Inset shows the location of Melbourne (red circle) within the State of Victoria in south east Australia. Overall the phylogeny aligns closely with the geography of originating cooling towers. For several outbreak codes polyclonality is apparent, as some common origins have connecting lines drawn from different sub-clades of the phylogeny. Red coloration on the base map represents population density within the greater Melbourne region. The branch lengths of the trees have been transformed and are proportional to the number of nodes under each parent node. (B) Core genome pairwise SNP comparisons

557     of within, and between, cooling tower isolate groups. Comparisons of specific epidemiologically defined

558     groups (infection clusters) are indicated with color codes as defined in the key. All groups had smaller

559     within diversity than between group diversity.
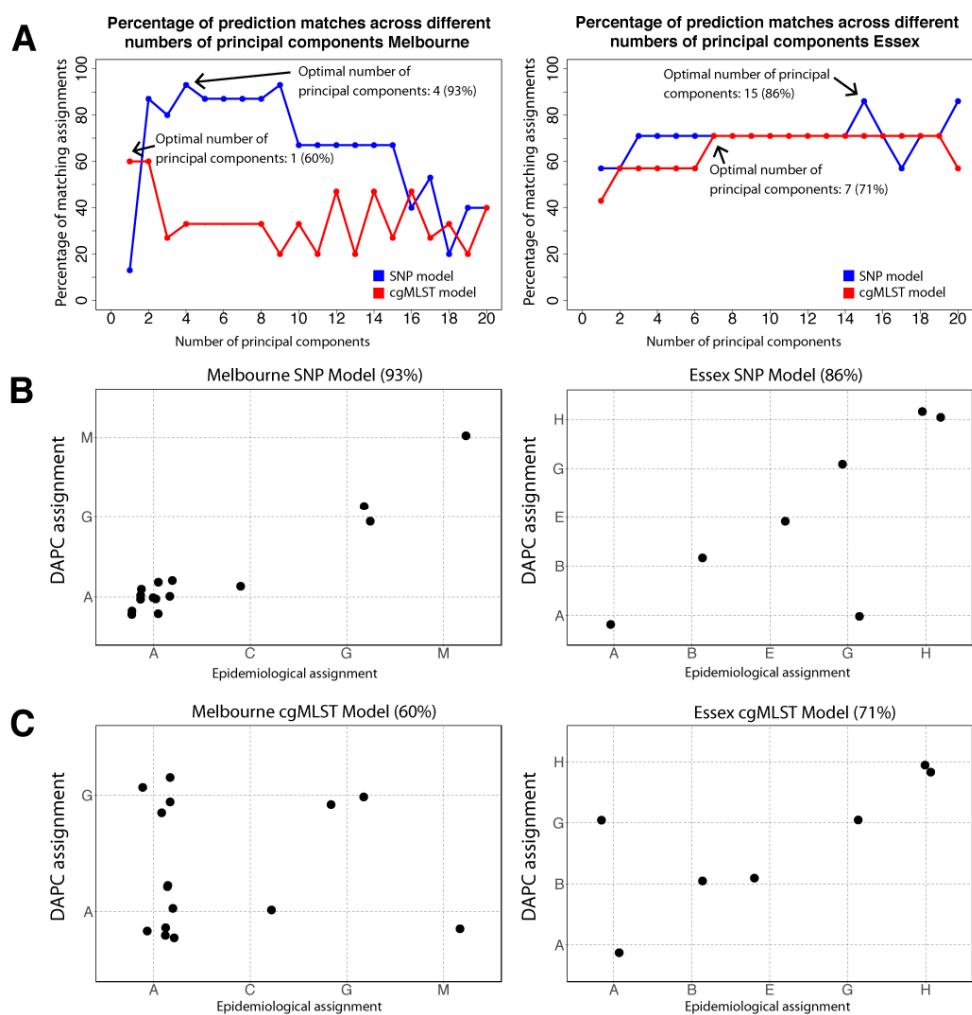
560

561

562

563

**Fig. 3. Scatterplots resulting from discriminant analysis of principal components (DAPC).** (A): Core genome single nucleotide polymorphisms (SNP) based models of the Melbourne (top left) and Essex (top right) datasets and; (B): core genome MLST (cgMLST) based models of the Melbourne (bottom left) and Essex (bottom right) datasets. The membership of each point within an epidemiologically defined cluster (*e.g.* "A" is the Melbourne Aquarium outbreak) is indicated by the colored circles and the corresponding letters labeled within squares. The amount of variation explained by the first and second discriminant functions are specified on the axes of each plot.

571

572

573

574



575

**Fig. 4**: **Discriminatory analysis of principal components (DAPC) modeling of Lpn-SG1 genomic data.** (A) Model comparison plots depicting the percentage of matches between the predicted and epidemiologically determined groupings of the validation set genomes across a range of 1-20 principal components for single nucleotide polymorphisms (SNP) and core genome MLST (cgMLST) DAPC models for the Melbourne and Essex datasets. The retention of four and one principal components was found to be optimal for the SNP (93% match) and cgMLST (60% match) models in Melbourne, respectively, while 15 and seven principal components were found to be optimal for the SNP (86% match) and cgMLST (71% match) models in the Essex hospital, respectively. (B) Assignment plots depicting the ability of the SNP models to predict the source attribution of the validation set clinical isolate genomes for the Melbourne and Essex hospital datasets. (C) Assignment plots depicting the ability of the cgMLST models to predict the source attribution of the validation set clinical isolate genomes for the Melbourne and Essex hospital datasets.