

Methods paper template

1 **SeroBA: rapid high-throughput serotyping**  
2 **of *Streptococcus pneumoniae* from whole**  
3 **genome sequence data**

4

5 Authors: Lennard Epping<sup>1,2</sup>, Andries J. van Tonder<sup>3</sup>, Rebecca A. Gladstone<sup>3</sup>, The Global  
6 Pneumococcal Sequencing consortium, Stephen D. Bentley<sup>3</sup>, Andrew J. Page<sup>\*,1,+</sup>, Jacqueline A.  
7 Keane<sup>\*,1,+</sup>

8

9 <sup>1</sup>Pathogen Informatics, Wellcome Trust Sanger Institute, Hinxton, Cambs, UK, CB10 1SA.

10 <sup>2</sup>Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin, Germany

11 <sup>3</sup>Infection Genomics, Wellcome Trust Sanger Institute, Hinxton, Cambs, UK, CB10 1SA.

12

13 \* Corresponding email: [jm15@sanger.ac.uk](mailto:jm15@sanger.ac.uk), [ap13@sanger.ac.uk](mailto:ap13@sanger.ac.uk)

14 + joint corresponding authors

15

---

16 **ABSTRACT**

17 *Streptococcus pneumoniae* is responsible for 240,000 - 460,000 deaths in children under 5 years of  
18 age each year. Accurate identification of pneumococcal serotypes is important for tracking the  
19 distribution and evolution of serotypes following the introduction of effective vaccines. Recent efforts  
20 have been made to infer serotypes directly from genomic data but current software approaches are  
21 limited and do not scale well. Here, we introduce a novel method, SeroBA, which uses a hybrid  
22 assembly and mapping approach. We compared SeroBA against real and simulated data and present  
23 results on the concordance and computational performance against a validation dataset, the robustness  
24 and scalability when analysing a large dataset, and the impact of varying the depth of coverage in the  
25 *cps* locus region on sequence-based serotyping. SeroBA can predict serotypes, by identifying the *cps*  
26 locus, directly from raw whole genome sequencing read data with 98% concordance using a *k*-mer  
27 based method, can process 10,000 samples in just over 1 day using a standard server and can call  
28 serotypes at a coverage as low as 10x. SeroBA is implemented in Python3 and is freely available  
29 under an open source GPLv3 license from: <https://github.com/sanger-pathogens/seroba>

30

31

---

32 **DATA SUMMARY**

- 33 1. The reference genome *Streptococcus pneumoniae* ATCC 700669 is available from National  
34 Center for Biotechnology Information (NCBI) with the accession number: FM211187  
35 2. Simulated paired end reads for experiment 2 have been deposited in FigShare:  
36 <https://doi.org/10.6084/m9.figshare.5086054.v1>  
37 3. Accession numbers for all other experiments are listed in Supplementary Table S1 and  
38 Supplementary Table S2.

39 **I/We confirm all supporting data, code and protocols have been provided within the article or**  
40 **through supplementary data files. ☒**

41

---

## 42 IMPACT STATEMENT

43 This article describes SeroBA, a *k*-mer based method for predicting the serotypes of *Streptococcus*  
44 *pneumoniae* from Whole Genome Sequencing (WGS) data. SeroBA can identify 92 serotypes and 2  
45 subtypes with constant memory usage and low computational costs. We showed that SeroBA is able  
46 to reliably predict serotypes at a depth of coverage as low as 10x and is scalable to large datasets.

47

48

---

## 49 INTRODUCTION

50 *Streptococcus pneumoniae* (the pneumococcus) is a clinically important bacterium estimated to cause  
51 700,000 to 1 million deaths in children under 5 years of age annually prior to the introduction of  
52 polysaccharide conjugate vaccines (O'Brien et al. 2009). The capsular polysaccharide biosynthesis  
53 (*cps*) locus, which encodes the serotype, is a major virulence factor in *S. pneumoniae*. The  
54 introduction of multi-valent pneumococcal conjugate vaccines has led to a substantial change in the  
55 circulating serotypes (Menezes et al. 2011) and decreased the number of deaths in children under 5  
56 years of age to 240,000 - 460,000 annually (Wahl et al. 2016). By surveilling the circulating  
57 serotypes, the epidemiological trends of *S. pneumoniae* can be observed, pre- and post-vaccination.  
58 The rapid reduction in the cost of whole genome sequencing (WGS) has led to its extensive use in  
59 the monitoring of pneumococcal serotypes (Lang et al. 2015)

60

61 To date there are nearly 100 known serotypes described for *S. pneumoniae* based on differing  
62 biochemical and antigenic properties of the capsule (Van Tonder et al. 2017). The *cps* locus, which  
63 encodes the serotype, can be very similar between serotypes from the same serogroup (such as  
64 serogroup 6) with some of them distinguished by a single nucleotide polymorphism (SNP), rendering  
65 a gene non-functional or altering the sugar linkage (Bentley et al. 2006). However, dissimilar loci may  
66 be grouped in the serogroup as they elicit a similar antibody response (e.g. serogroup 35). The large  
67 number of identified serotypes, and the high similarity between them, makes it challenging to  
68 computationally predict the serotype based on WGS data. Another challenge is recombination with  
69 other serotypes resulting in a mosaic *cps* locus (Salter et al. 2017) which may or may not affect the  
70 polysaccharide being produced. It is possible to have significant variation across the *cps* locus which  
71 does not lead to a different polysaccharide capsule being produced (Ko et al. 2013). Conversely novel

72 serotypes can be generated through these processes and can go unnoticed by antibody-based  
73 serotyping (Geno et al.; Park et al. 2007). Finally, mixed populations in a single sample and  
74 contamination can lead to ambiguity.

75

76 There are a number of methods available to predict serotypes in *S. pneumoniae*. Besides the gold  
77 standard method, Quellung, which can be subjective in certain cases, there are five additional methods  
78 based on serological tests, at least eight semi-automated molecular tests based on PCR and one  
79 method that uses microarray data for serotyping (Jauneikaite et al. 2015). There are a number of in-  
80 silico methods to detect the *cps* locus, which can then be used to predict serotypes from WGS data  
81 (Croucher et al. 2009; Leung et al. 2012; Kapatai et al. 2016; Metcalf et al. 2016). However, the tool  
82 described by Metcalf *et al.* is an in-house tool, and the tool described by Leung *et al.* only covers half  
83 of the known serotypes.

84

85 The only fully-featured automated pipeline for serotyping *S. pneumoniae* WGS data is PneumoCat,  
86 which was developed by Public Health England (PHE) (Kapatai et al. 2016). PneumoCat provides a  
87 capsular type variant (CTV) database including FASTA sequences for 92 serotypes and 2 subtypes as  
88 well as additional information about alleles, genes and SNPs for serotypes within specific serogroups.  
89 To predict a serotype, PneumoCat uses bowtie2 (Langmead and Salzberg 2012) to align reads to all  
90 serotype sequences. If the serotype belongs to a predefined serogroup or the serotype sequence could  
91 not be unambiguously identified, PneumoCat maps the reads to serogroup specific genes to identify  
92 the genetic variants. It is however computationally and memory intensive, and does not work with  
93 samples where there is a low depth of coverage in the *cps* locus region, as shown below.

94

95 To address these problems, we developed SeroBA, which makes efficient use of computational  
96 resources and can accurately detect the *cps* locus even at low coverage, and thus predict serotypes  
97 from WGS data using a database adapted from PneumoCAT (Kapatai et al. 2016). This accuracy was  
98 evaluated by comparing the results to a standard, validated dataset from PHE (Kapatai et al. 2016).  
99 We showed that it is scalable and robust by calculating the serotypes of 9,477 samples from the GPS  
100 project, an ongoing global pneumococcal sequencing project, on commodity hardware. Simulated  
101 read data, with varying coverage over a known reference sequence, was used to show the minimum  
102 depth of coverage required to call a serotype.

103

104

---

## 105 THEORY AND IMPLEMENTATION

106

107 SeroBA takes Illumina paired-end reads in FASTQ format as input as shown in Figure 1.  
108 Precomputed databases are bundled with the application that describe the serotypes. The first of these  
109 is a *k*-mer counts database for every serotype sequence produced by KMC (v3.0.0) (Kokot et al.  
110 2017), the second is an ARIBA (v 2.9.3) (Hunt et al. 2017) compatible database for every serotype,  
111 and the third is a capsular type variant (CTV) database, including FASTA sequences for 92 serotypes

112 and 2 subtypes, as well as additional information about alleles, genes and SNPs for serotypes in  
113 specific serogroups. These databases were adapted from PneumoCAT (Kapatai et al. 2016). A *k*-mer  
114 analysis is performed on the input reads, and the intersection is found between these *k*-mers and the  
115 precomputed *k*-mer database of serotypes. The *k*-mer coverage of the input reads over the serotype  
116 sequences is normalised to the sequence length of the serotype sequence and the serotype with the  
117 highest normalised sequence coverage is selected. This step identifies the possible serotype or  
118 serogroup and ARIBA is used to confirm the presence of the selected serotype from the raw reads. If a  
119 serogroup is selected, the *cps* sequence produced by ARIBA and serotype specific genes are aligned  
120 with nucmer (Kurtz et al. 2004) to find specific variants, such as presence/absence of genes, SNPs, or  
121 gene truncations as defined in the CTV. The output of SeroBA includes the predicted serotype with  
122 detailed information that led to the prediction, as well as an assembly of the *cps* locus sequences.

123

## 124 VALIDATION DATASET

125 A validation dataset consisting of 2,065 UK isolates (Supplementary Table S1) retrieved from the  
126 PHE archive was originally used to evaluate PneumoCat. It consists of 72 out of 92 known serotypes,  
127 including all serotypes contained in commercial vaccines, and 19 non-typeable samples. The serotype  
128 of each sample was confirmed by latex agglutination with Statens Serum Institut typing sera (Kapatai  
129 et al. 2016). PneumoCat v1.1 (Kapatai et al. 2016) and SeroBA v0.1 were evaluated on an AMD  
130 Opteron 6272 server running Ubuntu 12.04.2 LTS, with 32 cores and 256GB of RAM. A single CPU  
131 (Central Processing Unit) was used for each experiment, repeated 10 times, with the mean memory  
132 usage and wall clock times noted.

133 Figure 2 summarises the serotypes called for each sample by each method. As serotyping with latex  
134 agglutination and Quellung can be subjective (Selva et al. 2012) and potentially imprecise, a serotype  
135 was said to be concordant if two or more methods agreed on the same serotype. This gave a  
136 concordance of 98.4% for SeroBA and 98.5% for PneumoCat with latex agglutination method. The  
137 reference sequences in the CTV for the serotypes 24A, 24B, 24F may not be representative for the  
138 circulating strains (Kapatai et al. 2016), so SeroBA will report serogroup 24 instead of reporting the  
139 serotype. As discussed in (Kapatai et al. 2016) serological prediction in serogroup 12 were error-prone,  
140 so a prediction of either 12B or 12F were counted as concordant.

141

142 The overall computational resources required to call the serotypes differed substantially between  
143 PneumoCat and SeroBA (Table 1): SeroBA was fifteen times faster and required five times less  
144 memory than PneumoCat.

## 145 EVALUATION USING A LARGE DATASET

146 To show the scalability of SeroBA to large datasets, we took 9,477 *S. pneumoniae* samples from the  
147 GPS project (Supplementary Table S2) and calculated the serotypes using the hardware setup  
148 previously described. A comparison with serotypes determined using experimental methods gave an  
149 accuracy of 98.2% for SeroBA. The serotypes were determined by different experimental methods as  
150 listed in Supplementary Table S2. Using all 32 cores resulted in a total wall-clock time of 823.78  
151 hours. This showed that SeroBA can robustly scale to large datasets.

## 152 **IMPACT OF DEPTH OF COVERAGE**

153 The effect of depth of coverage on the serotyping results produced by SeroBA and PneumoCat was  
154 evaluated by simulating perfect paired end reads over the serotype 23F *cps* locus from the  
155 *Streptococcus pneumoniae* ATCC 700669 (accession code: FM211187) reference genome (Croucher  
156 et al. 2009). Flanking regions of 1,000 bases were included on either side of the *cps* locus to eliminate  
157 confounding effects of low coverage at the locus boundaries. The reads with a length of 125 base  
158 pairs were generated by FASTAQ (v3.15.0) (<https://github.com/sanger-pathogens/Fastaq>) with an  
159 insert size of 500 bases and standard deviation of 50 with varying depth of coverage from 1x to 50x  
160 and from 100x to 350x in steps of 50. SeroBA started to predict serotype 23F at a depth of coverage  
161 of 10x while PneumoCat required nearly twice as much, needing at least 19x coverage. The  
162 computational resources required by SeroBA remained constant with increasing depth of coverage;  
163 however, the computational resource requirements of PneumoCat continue to grow linearly (Figure  
164 3). At 350x coverage, PneumoCat took 3 times longer than SeroBA. Similarly, the amount of memory  
165 required by SeroBA stabilised at 150MB, regardless of coverage, whereas PneumoCat's memory  
166 requirement grew with the depth of coverage, requiring 3 times more than SeroBA at 350x coverage.  
167 Each experiment was repeated 10 times and the mean was calculated.

168

169

---

## 170 **CONCLUSION**

171 In this paper, we described SeroBA a method for predicting serotypes from *S. pneumoniae* Illumina  
172 NGS reads. We compared SeroBA and PneumoCat to a gold standard experimental serotyping  
173 method and showed that they had approximately the same level of concordance. However, SeroBA  
174 was fifteen times faster and required five times less memory than PneumoCat. The assembly of the  
175 *cps* locus sequence provides by SeroBA is another key feature that is very useful for further analyses  
176 and reference free comparisons. SeroBA was able to predict the serotype from only 10x read depth  
177 and scaled well on a large dataset of nearly 10,000 samples with a prediction accuracy of over 98%.

178

---

## 179 **AUTHOR STATEMENTS**

180

181 [REMOVED FOR BLIND REVIEW]

182

183

---

## 184 **ABBREVIATIONS**

185 SNP: Single nucleotide polymorphism

186 WGS: Whole genome sequencing

187 CTV: Capsular Type Variant database

188 CPS: Capsular polysaccharide biosynthesis

189 GPS: The Global Pneumococcal Sequencing

190

191

---

## 192 REFERENCES

193 Bentley SD, Aanensen DM, Mavroidi A, Saunders D, Rabinowitsch E, Collins M, et al. Genetic  
194 Analysis of the Capsular Biosynthetic Locus from All 90 Pneumococcal Serotypes. *PLoS Genet*  
195 [Internet]. 2006;2(3):e31–e31. Available from:  
196 <http://dx.plos.org/10.1371/journal.pgen.0020031>

197 Croucher NJ, Walker D, Romero P, Lennard N, Paterson GK, Bason NC, et al. Role of Conjugative  
198 Elements in the Evolution of the Multidrug-Resistant Pandemic Clone *Streptococcus*  
199 *pneumoniae* Spain23F ST81. *J Bacteriol* [Internet]. 2009a Mar;191(5):1480–9. Available from:  
200 <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2648205/>

201 Geno KA, Saad JS, Nahm MH. Discovery of Novel Pneumococcal Serotype 35D, a Natural WciG-  
202 Deficient Variant of Serotype 35B.

203 Hunt M, Mather AE, Sánchez-Busó L, Page AJ, Parkhill J, Keane JA, et al. ARIBA: rapid antimicrobial  
204 resistance genotyping directly from sequencing reads. *bioRxiv*. 2017;118000.

205 Jauneikaite E, Tocheva AS, Jefferies JMC, Gladstone RA, Faust SN, Christodoulides M, et al. Current  
206 methods for capsular typing of *Streptococcus pneumoniae*. Vol. 113, *Journal of*  
207 *Microbiological Methods*. 2015. p. 41–9.

208 Kapatai G, Sheppard CL, Al-Shahib A, Litt DJ, Underwood AP, Harrison TG, et al. Whole genome  
209 sequencing of *Streptococcus pneumoniae*: development, evaluation and verification of  
210 targets for serogroup and serotype prediction using an automated pipeline. *PeerJ* [Internet].  
211 2016a Sep;4:e2477–e2477. Available from: <https://peerj.com/articles/2477>

212 Ko KS, Baek JY, Song J-H. Capsular gene sequences and genotypes of “serotype 6E”  
213 *Streptococcus pneumoniae* isolates. *J Clin Microbiol*. 2013 Oct;51(10):3395–9.

214 Kokot M, Długosz M, Deorowicz S. KMC 3: counting and manipulating k-mer statistics. 2017 Jan 27;

215 Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and open  
216 software for comparing large genomes. *Genome Biol*. 2004;5(2):R12.

217 Lang ALS, McNeil SA, Hatchette TF, Elsharif M, Martin I, LeBlanc JJ. Detection and prediction of  
218 *Streptococcus pneumoniae* serotypes directly from nasopharyngeal swabs using PCR. *J Med*  
219 *Microbiol*. 2015 Aug;64(8):836–44.

220 Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*. 2012;9(4):357–9.

- 221 Leung MH, Bryson K, Freystatter K, Pichon B, Edwards G, Charalambous BM, et al. Sequotyping:  
222 Serotyping *Streptococcus pneumoniae* by a single PCR sequencing strategy. *J Clin Microbiol.*  
223 2012;50(7):2419–27.
- 224 Menezes AP de O, Campos LC, dos Santos MS, Azevedo J, dos Santos RCN, Carvalho M da GS, et al.  
225 Serotype distribution and antimicrobial resistance of *Streptococcus pneumoniae* prior to  
226 introduction of the 10-valent pneumococcal conjugate vaccine in Brazil, 2000–2007. *Vaccine.*  
227 2011;29(6):1139–44.
- 228 Metcalf BJ, Gertz RE, Gladstone RA, Walker H, Sherwood LK, Jackson D, et al. Strain features and  
229 distributions in pneumococci from children with invasive disease before and after 13-valent  
230 conjugate vaccine implementation in the USA. *Clin Microbiol Infect* [Internet]. 2016  
231 Jan;22(1):60.e9-60.e29. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/26363404>
- 232 O'Brien KL, Wolfson LJ, Watt JP, Henkle E, Deloria-Knoll M, McCall N, et al. Burden of disease caused  
233 by *Streptococcus pneumoniae* in children younger than 5 years: global estimates. *Lancet*  
234 [Internet]. 2009;374(9693):893–902. Available from:  
235 <http://www.ncbi.nlm.nih.gov/pubmed/19748398>
- 236 Park IH, Pritchard DG, Cartee R, Brandao A, Brandileone MCC, Nahm MH. Discovery of a new  
237 capsular serotype (6C) within serogroup 6 of *Streptococcus pneumoniae*. *J Clin Microbiol.*  
238 2007 Apr;45(4):1225–33.
- 239 Salter SJ, Hinds J, Gould KA, Lambertsen L, Hanage WP, Antonio M, et al. Variation at the capsule  
240 locus, *cps*, of mistyped and non-typable *Streptococcus pneumoniae* isolates. 2017;1231.
- 241 Selva L, del Amo E, Brotons P, Muñoz-Almagro C. Rapid and easy identification of capsular serotypes  
242 of *Streptococcus pneumoniae* by use of fragment analysis by automated fluorescence-based  
243 capillary electrophoresis. *J Clin Microbiol.* 2012 Nov;50(11):3451–7.
- 244 Van Tonder AJ, Bray JE, Quirk SJ, Haraldsson G, Jolley KA, Maiden MCJ, et al. Putatively novel  
245 serotypes and the potential for reduced vaccine effectiveness: capsular locus diversity  
246 revealed among 5405 pneumococcal genomes. 2017;
- 247 Wahl B, O'Brien KL, Greenbaum A, Liu L, Chu Y, Black R, et al. Global burden of *Streptococcus*  
248 *pneumoniae* in children younger than 5 years in the pneumococcal conjugate vaccines (PCV)  
249 era: 2000-2015. ISPPD-10 [Internet]. 2016 Dec 15; Available from:  
250 <http://beta.bib.irb.hr/850035>

251

252

---

## 253 DATA BIBLIOGRAPHY

- 254 1. <https://github.com/sanger-pathogens/seroba>
- 255 2. Lennard Epping, figshare. DOI: <https://doi.org/10.6084/m9.figshare.5086054.v1>

256 3. *Croucher, N. J., Streptococcus pneumoniae* ATCC 700669. NCBI. , FM211187

257

---

258 **FIGURES AND TABLES**

259 Table 1: Performance of SeroBA and PneumoCat on the validation set

Tool	Mean Wall Clock Time (m)	Mean RAM usage (MB)
PneumoCat	65.84	922.89
SeroBA	4.53	187.82

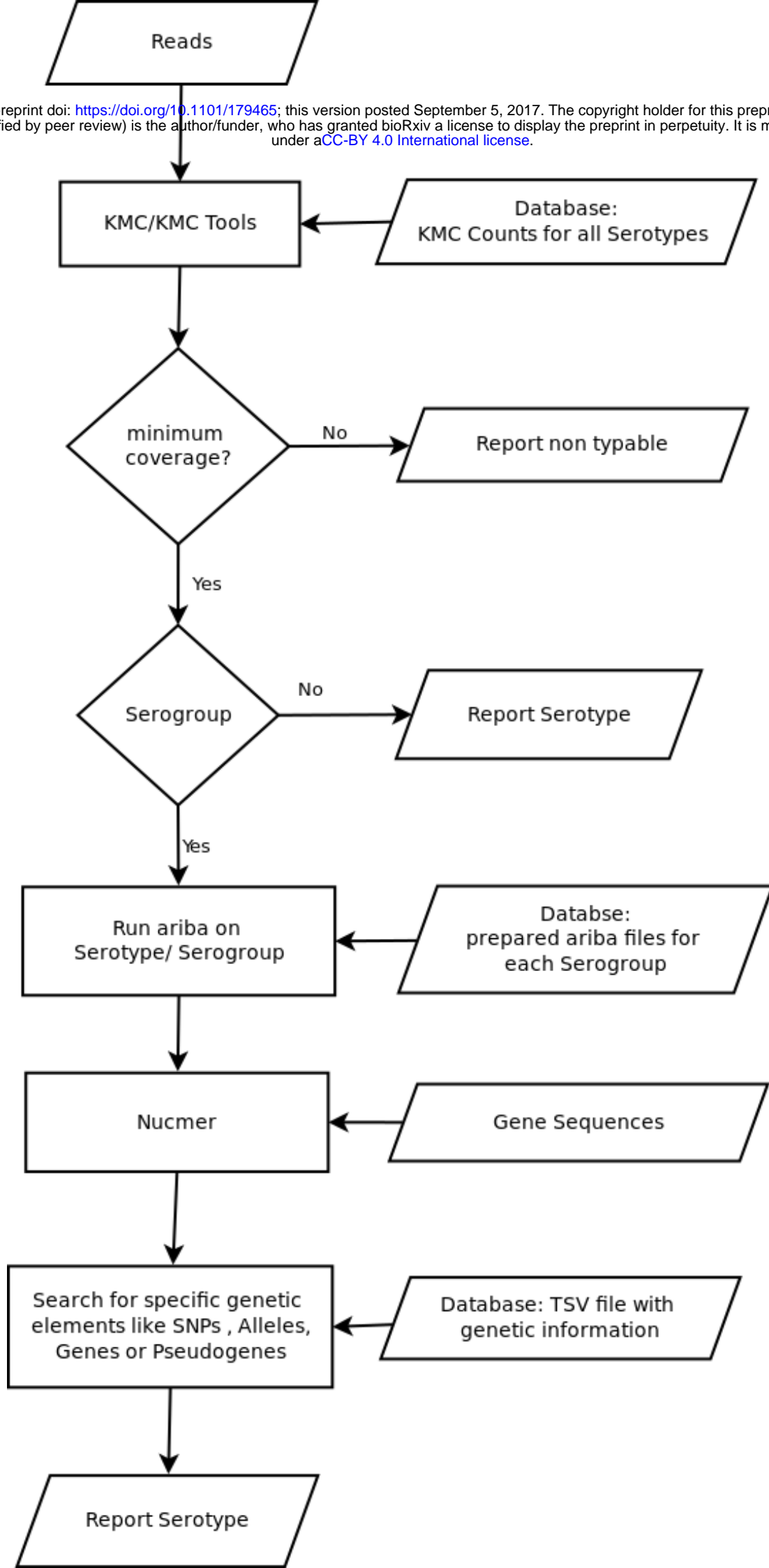
260

261 Figure 1: Flowchart outlining the main steps of the SeroBA algorithm

262 Figure 2: Agreement of serotyping results between different methods

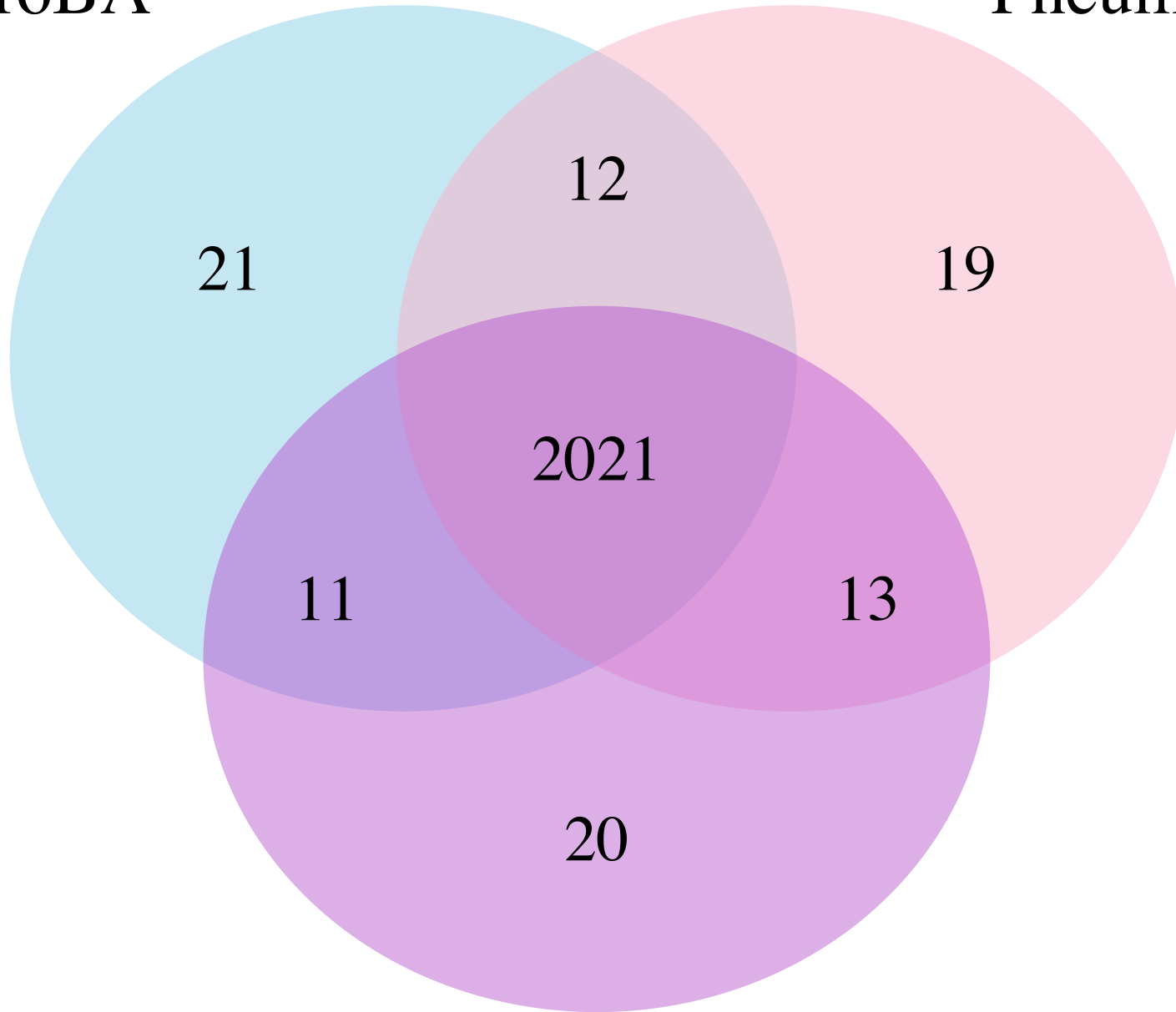
263 Figure 3: a) mean CPU time in seconds used by SeroBA and PneumoCat when varying the coverage  
264 from 1x to 350x; b) maximum memory allocation of SeroBA and PneumoCat when varying the  
265 coverage from 1x to 350x. Each data point represents the mean value of ten identical experiments.





SeroBA

PneumoCat



experimental method

