

# 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>, Martin Hunt<sup>3</sup>, Andries J. van Tonder<sup>4</sup>, Rebecca A. Gladstone<sup>4</sup>, The  
6 Global Pneumococcal Sequencing consortium, Stephen D. Bentley<sup>4</sup>, Andrew J. Page<sup>\*1,+</sup>, Jacqueline  
7 A. 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>Nuffield Department of Medicine, University of Oxford, Oxford, UK, OX3 7BN.

12 <sup>4</sup>Infection Genomics, Wellcome Trust Sanger Institute, Hinxton, Cambs, UK, CB10 1SA.

13

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

15 + joint corresponding authors

16

---

## 17 **ABSTRACT**

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

31

32

---

## 33 **DATA SUMMARY**

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

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

42

---

## 43 **IMPACT STATEMENT**

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

48

49

---

## 50 **INTRODUCTION**

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

61

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

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

76

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

85

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

95

96 To address these problems, we developed SeroBA, which makes efficient use of computational  
97 resources and can accurately detect the *cps* locus even at low coverage, and thus predict serotypes  
98 from WGS data using a database adapted from PneumoCAT (Kapatai et al. 2016). This accuracy was  
99 evaluated by comparing the results to a standard, validated dataset from PHE (Kapatai et al. 2016).  
100 We showed that it is scalable and robust by calculating the serotypes of 9,886 samples from the GPS  
101 project, an ongoing global pneumococcal sequencing project, on commodity hardware. Simulated  
102 read data, with varying coverage over a known reference sequence, was used to show the minimum  
103 depth of coverage required to call a serotype.

104

105

---

## 106 THEORY AND IMPLEMENTATION

107

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

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

124

## 125 **VALIDATION DATASET**

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

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

142

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

## 146 **EVALUATION USING A LARGE DATASET**

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

## 153 **IMPACT OF DEPTH OF COVERAGE**

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

169

170

---

## 171 **CONCLUSION**

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

179

---

## 180 **AUTHOR STATEMENTS**

181

182 [REMOVED FOR BLIND REVIEW]

183

184

---

## 185 **ABBREVIATIONS**

186 SNP: Single nucleotide polymorphism

187 WGS: Whole genome sequencing

188 CTV: Capsular Type Variant database

189 CPS: Capsular polysaccharide biosynthesis

190 GPS: The Global Pneumococcal Sequencing

191

192

---

## 193 REFERENCES

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

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

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

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

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

209 Kapatai G, Sheppard CL, Al-Shahib A, Litt DJ, Underwood AP, Harrison TG, et al. Whole genome  
210 sequencing of *Streptococcus pneumoniae*: development, evaluation and verification of  
211 targets for serogroup and serotype prediction using an automated pipeline. *PeerJ* [Internet].  
212 2016a Sep;4:e2477–e2477. Available from: <https://peerj.com/articles/2477>

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

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

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

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

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

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

252

253

---

## 254 DATA BIBLIOGRAPHY

255 1. <https://github.com/sanger-pathogens/seroba>

256 2. Lennard Epping, figshare. DOI: <https://doi.org/10.6084/m9.figshare.5086054.v1>

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

258

---

259 **FIGURES AND TABLES**

260 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

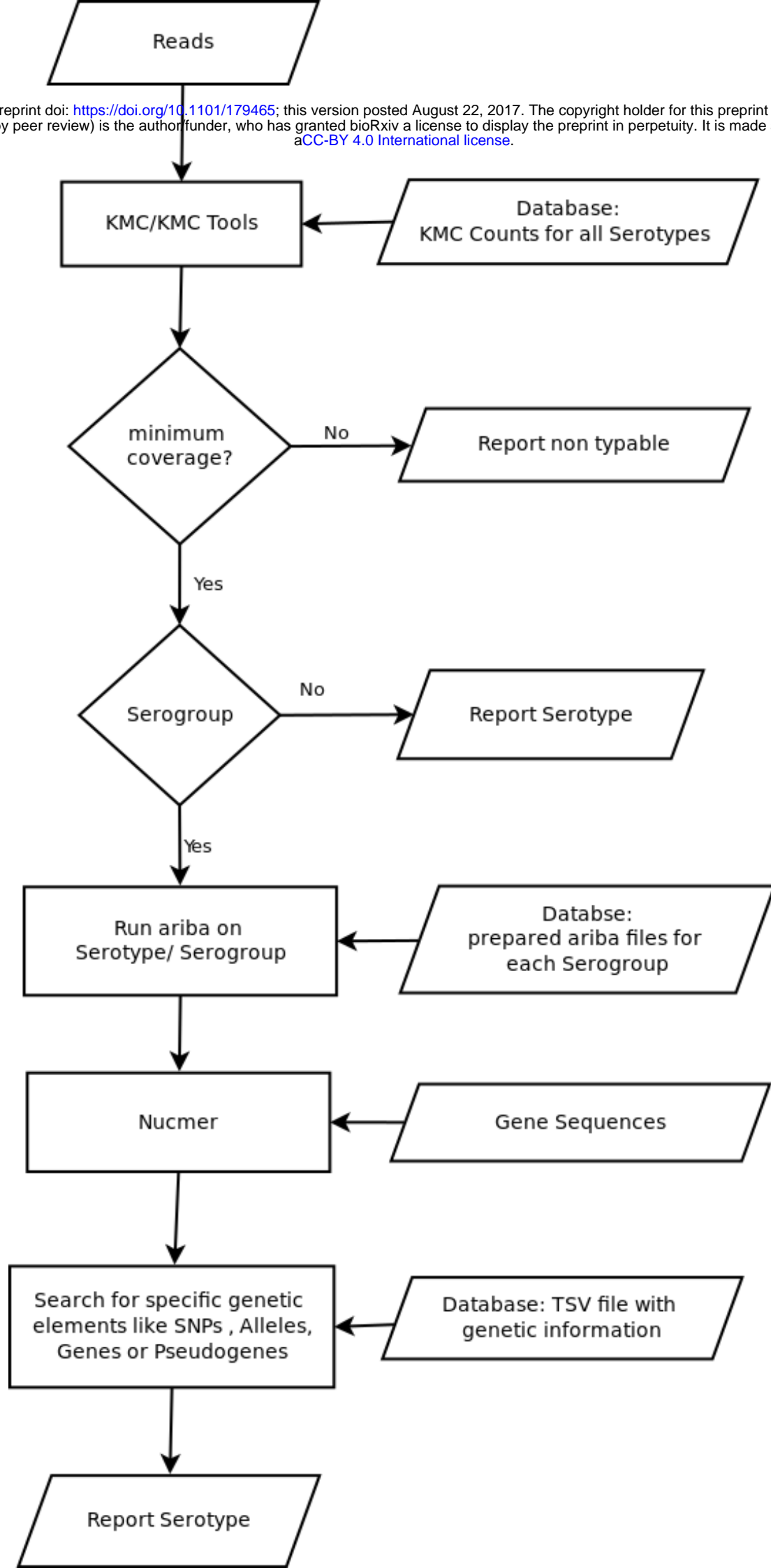
261

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

263 Figure 2: Agreement of serotyping results between different methods

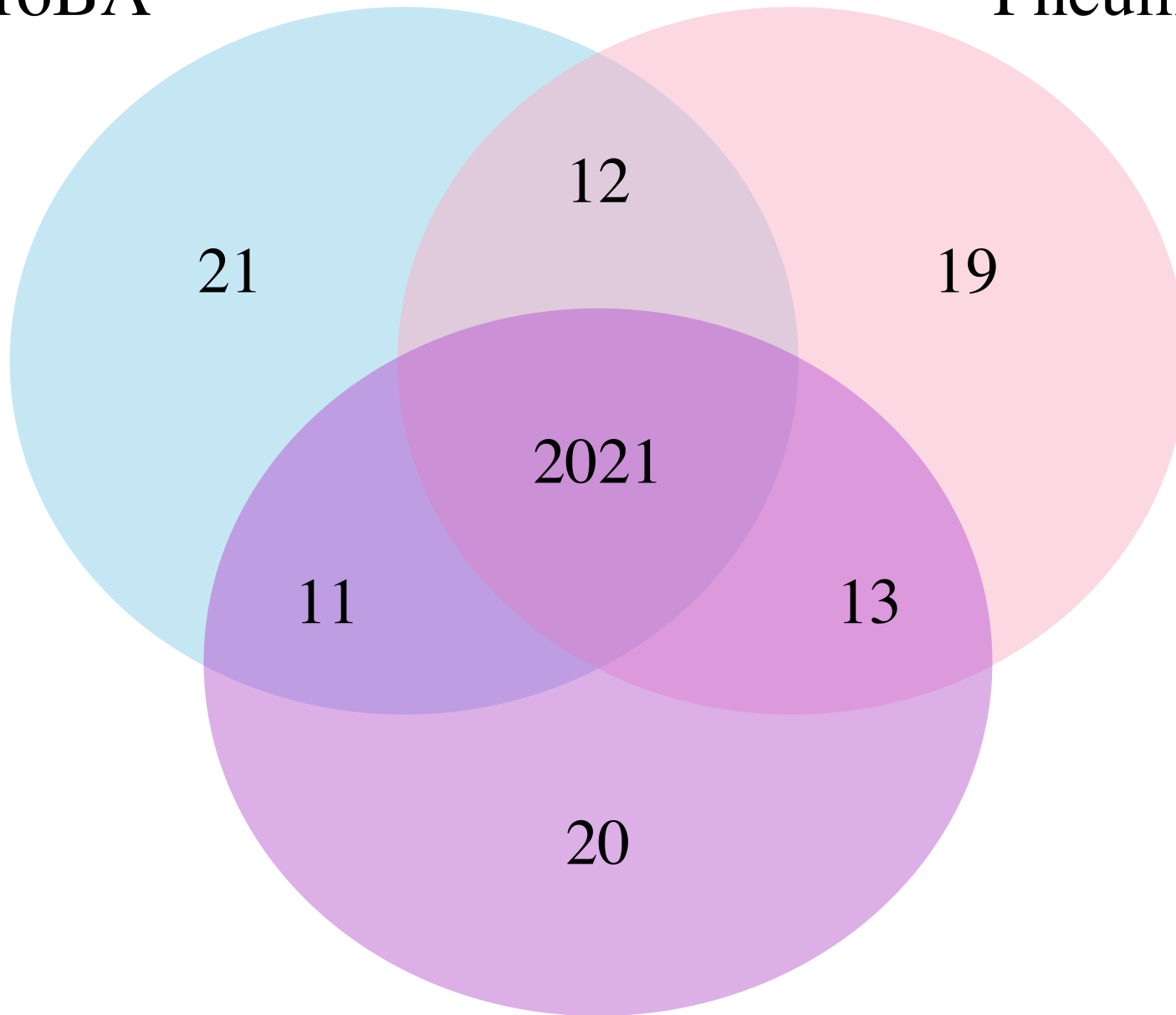
264 Figure 3: a) mean CPU time in seconds used by SeroBA and PneumoCat when varying the coverage  
265 from 1x to 350x; b) maximum memory allocation of SeroBA and PneumoCat when varying the  
266 coverage from 1x to 350x. Each data point represents the mean value of ten identical experiments.





SeroBA

PneumoCat



experimental method

