1 # Capsular typing method for *Streptococcus agalactiae* using whole

2 genome sequence data

3 Running title: GBS capsular typing using whole genome sequence data

4 Anna E Sheppard,[1#] Alison Vaughan,[1] Nicola Jones,[1] Paul Turner,[2,3] Claudia

5 Turner,[2,3] Androulla Efstratiou,[4,5] Darshana Patel,[4] the Modernising Medical

6 Microbiology (MMM) Informatics Group,[1] A Sarah Walker,[1] James A Berkley,[2,6]

7 Derrick W Crook,[1] Anna C Seale.[2,6#]

8 1. Modernising Medical Microbiology Consortium, Nuffield Department of Clinical

9    Medicine, University of Oxford, UK.

10 2. Centre for Tropical Medicine, Nuffield Department of Clinical Medicine, University

11    of Oxford, UK.

12 3. Shoklo Malaria Research Unit, Thailand

13 4. Microbiology Reference Division, Public Health England, London, UK

14 5. Imperial College, London, UK

15 6. KEMRI-Wellcome Trust Research Programme, Kilifi, Kenya.

16

17 [#]Corresponding authors: Anna E Sheppard (anna.sheppard@ndm.ox.ac.uk) or Anna

18 C Seale (aseale@nhs.net).

19

20 Word count:

21 Abstract 47/50

22 Main text 933/1000

1

23 Abstract

24 Group B streptococcus (GBS) capsular serotype is a major determinant of virulence,

25 and affects potential vaccine coverage. Here we report a whole genome sequencing-

26 based method for GBS serotype assignment. This shows high agreement

27 (kappa=0.92) with conventional methods, and increased serotype assignment

28 (100%) to all ten capsular types.

29

30 Main text

31 *Streptococcus agalactiae,* or Group B Streptococcus (GBS), is an important

32 pathogen in neonates (1-3), with early infection acquired from the maternal genito-

33 urinary tract (4). In addition, GBS is now recognised as an increasingly important

34 pathogen in high-income regions in immunosuppressed and elderly individuals (5, 6).

35 GBS expresses a capsular polysaccharide, which is involved in virulence and

36 immune evasion. Ten different variants, or serotypes, have been described (Ia, Ib, II,

37 III, IV, V VI, VII, VIII and IX), which differ in their disease-causing ability. Conjugate

38 vaccines targeting the most common disease-causing serotypes are currently in

39 development (7). Establishing vaccine serotype coverage is important, as is

40 surveillance post-introduction to monitor for potential serotype replacement, as has

41 been seen following the introduction of other conjugate vaccines (8).

42 Current methods for GBS serotype allocation rely on latex agglutination assays or

43 PCR (9). Recent advances in whole genome sequencing (WGS) have enabled the

44 development of approaches that can be used in place of traditional microbiological

45 methods, such as strain typing and antibiotic susceptibility profiling (10-12). A major

46 advantage of this approach is that the cost of sequencing can be mitigated by the

47  ability to use the same data to generate multiple outputs. Given the decreasing cost

48  of WGS (13), and the ongoing increase in WGS data generation, we sought to

49  establish and validate a WGS-based method for GBS capsular typing.

50  We developed an algorithm for serotype assignment on the basis of sequence

51  similarity between a given *de novo* assembly and capsular gene sequences of the

52  ten GBS serotypes. For nine serotypes, published sequences were used as

53  references (Table 1), while for serotype IX, only a partial capsular locus sequence

54  has been published (14). A suitable reference for the full capsular locus region was

55  therefore determined by WGS of a serotype IX isolate obtained from the Statens

56  Serum Institute, Denmark.

57  To assign serotype for a given isolate, a BLAST database was generated from the

58  *de novo* assembly and queried with the variable region of the capsular locus

59  sequence for each serotype (*cpsG-cpsK* for serotypes Ia-VII and IX, *cpsR-cpsK* for

60  serotype VIII) using BLASTn with an evalue threshold of 1e-100 and otherwise

61  default parameters. A serotype was considered as correct if it showed ≥95%

62  sequence identity over ≥90% of the sequence length. These thresholds were chosen

63  on the basis of being stringent enough to provide differentiation between the various

64  reference sequences, while maximising serotype allocation on an initial test set of

65  publicly available GBS WGS data, where serotype information was not available (so

66  we had no way of knowing whether the assigned serotypes were in fact correct).

67  This sequence-based method for serotype allocation was validated using WGS on a

68  set of 223 colonising or invasive human isolates from Canada, Latin America,

69  Singapore, UK, USA, and Thailand, for which serotype had previously been

70  determined using conventional latex agglutination assays, with PCR used to confirm

71  weak positives or negatives in a subset (15-17). For two rare serotypes (Serotype

72  VIII and IX), one isolate of each was obtained from the Statens Serum Institute,

73  Denmark. GBS isolates stored at -80°C were sub-cultured on Columbia blood agar

74  for 24-48 hours, followed by DNA extraction from a single colony using a commercial

75  kit (QuickGene, Fujifilm, Tokyo, Japan). High throughput sequencing was

76  undertaken at the Wellcome Trust Centre for Human Genetics (Oxford University,

77  UK) using the Illumina HiSeq2500 platform, generating 150 base paired-end reads.

78  *De novo* assembly was performed using Velvet and VelvetOptimiser (18, 19).

79  Agreement between serotype allocations was tested with the Kappa statistic.

80  High quality sequence data were obtained for all 223 GBS isolates (median read

81  number: 2,975,508, range: 1,798,744-13,073,718; median contig number: 46, range

82  16-106; median assembly length: 2.05 Mb, range: 1.94-2.22 Mb). Each isolate was

83  allocated to a single serotype using the WGS data (Table 2). Three isolates that did

84  not have a capsular type assigned by latex agglutination methods had serotypes Ib,

85  VI and VIII assigned. For all previously serotyped GBS isolates with a known capsule

86  type, the kappa statistic (0.92) indicated very high agreement between WGS-

87  predicted and conventional serotype. There were nine discordant isolates. In each

88  case there was strong support for the sequence-allocated serotype, with >98%

89  sequence identity over 100% of the reference length in all nine cases (Figure 1).

90  Across all isolates, differences in relatedness between the capsular locus sequences

91  of the different serotypes led to characteristic relationships between the allocated

92  (best match) serotype and the second-best match. For example, all isolates

93  assigned as serotype Ia had serotype III as the second-best match. In all cases, the

94  second-best match was substantially poorer than the best match, demonstrating that

95  there was no ambiguity in predicted serotype (Figure 1, Table 3).

96    The nine discordant and three non-typeable isolates were retested by latex

97    agglutination (Table 4) and resequenced using the Illumina MiSeq platform, with

98    sequence processing and WGS-based serotype prediction performed as above. In

99    all cases, resequencing was consistent with the initial WGS classification. For 6/9

100   discordant isolates, the new latex agglutination results matched the WGS-based

101   prediction, suggesting that the initial discordance may have resulted from incorrect

102   latex agglutination typing or sample mislabelling. The other three initially discordant

103   isolates, and the three non-typeable isolates, were all classified as non-typeable on

104   retesting.

105   This WGS-based method for GBS serotyping, validated using 223 isolates typed

106   using conventional methods, was therefore highly accurate. Although WGS may not

107   currently be cost-effective for directly replacing traditional serotyping, costs are likely

108   to further decrease. Furthermore, WGS may already be the cheapest option for

109   combined studies, with possibilities to utilise the resulting data for additional

110   analyses such as multi-locus sequence typing, relatedness to other sequenced

111   isolates, and detailed phylogenetic analysis.

112

## Modernising Medical Microbiology (MMM) informatics group

Jim Davies, Charles Crichton, Milind Acharya, Carlos del Ojo Elias

134     References

135     1.      **Baker CJ, Barrett FF, Gordon RC, Yow MD.** 1973. Suppurative meningitis due to streptococci

136             of Lancefield group B: a study of 33 infants. J Pediatr **82:**724-729.

137     2.      **Barton LL, Feigin RD, Lins R.** 1973. Group B beta hemolytic streptococcal meningitis in

138             infants. J Pediatr **82:**719-723.

139     3.      **Communicable Disease Surveillance Centre London.** 1985. Neonatal meningitis: a review of

140             routine national data 1975-83. Br Med J (Clin Res Ed) **290:**778-779.

141     4.      **Dillon HC, Jr., Gray E, Pass MA, Gray BM.** 1982. Anorectal and vaginal carriage of group B

142             streptococci during pregnancy. J Infect Dis **145:**794-799.

143     5.      **Schuchat A.** 1998. Epidemiology of group B streptococcal disease in the United States:

144             shifting paradigms. Clin Microbiol Rev **11:**497-513.

145     6.      **Phares CR, Lynfield R, Farley MM, Mohle-Boetani J, Harrison LH, Petit S, Craig AS, Schaffner**

146             **W, Zansky SM, Gershman K, Stefonek KR, Albanese BA, Zell ER, Schuchat A, Schrag SJ,**

147             **Active Bacterial Core surveillance/Emerging Infections Program N.** 2008. Epidemiology of

148             invasive group B streptococcal disease in the United States, 1999-2005. JAMA **299:**2056-

149             2065.

150     7.      **Madhi SA, Dangor Z, Heath PT, Schrag S, Izu A, Sobanjo-Ter Meulen A, Dull PM.** 2013.

151             Considerations for a phase-III trial to evaluate a group B Streptococcus polysaccharide-

152             protein conjugate vaccine in pregnant women for the prevention of early- and late-onset

153             invasive disease in young-infants. Vaccine **31 Suppl 4:**D52-57.

154     8.      **Mulholland K, Satzke C.** 2012. Serotype replacement after pneumococcal vaccination.

155             Lancet **379:**1387; author reply 1388-1389.

156     9.      **Imperi M PM, Alfarone G, Baldassarri L, Orefici G, Creti R.** 2010. A multiplex PCR assay for

157             the direct identification of the capsular type (Ia to IX) of Streptococcus agalactiae. J

158             Microbiol Methods **80:**212-214.

159    10.    **Gordon NC, Price JR, Cole K, Everitt R, Morgan M, Finney J, Kearns AM, Pichon B, Young B,**

160            **Wilson DJ, Llewelyn MJ, Paul J, Peto TE, Crook DW, Walker AS, Golubchik T.** 2014.

161            Prediction of Staphylococcus aureus antimicrobial resistance by whole-genome sequencing.

162            J Clin Microbiol **52:**1182-1191.

163    11.    **Stoesser N, Batty EM, Eyre DW, Morgan M, Wyllie DH, Del Ojo Elias C, Johnson JR, Walker**

164            **AS, Peto TE, Crook DW.** 2013. Predicting antimicrobial susceptibilities for Escherichia coli

165            and Klebsiella pneumoniae isolates using whole genomic sequence data. J Antimicrob

166            Chemother **68:**2234-2244.

167    12.    **Walker TM, Kohl TA, Omar SV, Hedge J, Del Ojo Elias C, Bradley P, Iqbal Z, Feuerriegel S,**

168            **Niehaus KE, Wilson DJ, Clifton DA, Kapatai G, Ip CL, Bowden R, Drobniewski FA, Allix-**

169            **Beguec C, Gaudin C, Parkhill J, Diel R, Supply P, Crook DW, Smith EG, Walker AS, Ismail N,**

170            **Niemann S, Peto TE, Modernizing Medical Microbiology Informatics G.** 2015. Whole-

171            genome sequencing for prediction of Mycobacterium tuberculosis drug susceptibility and

172            resistance: a retrospective cohort study. Lancet Infect Dis doi:10.1016/S1473-

173            3099(15)00062-6.

174    13.    **Loman NJ, Constantinidou C, Chan JZ, Halachev M, Sergeant M, Penn CW, Robinson ER,**

175            **Pallen MJ.** 2012. High-throughput bacterial genome sequencing: an embarrassment of

176            choice, a world of opportunity. Nat Rev Microbiol **10:**599-606.

177    14.    **Slotved HC, Kong F, Lambertsen L, Sauer S, Gilbert GL.** 2007. Serotype IX, a Proposed New

178            Streptococcus agalactiae Serotype. J Clin Microbiol **45:**2929-2936.

179    15.    **Bisharat N, Crook DW, Leigh J, Harding RM, Ward PN, Coffey TJ, Maiden MC, Peto T, Jones**

180            **N.** 2004. Hyperinvasive neonatal group B streptococcus has arisen from a bovine ancestor. J

181            Clin Microbiol **42:**2161-2167.

182    16.    **Jones N, Oliver K, Jones Y, Haines A, Crook D.** 2006. Carriage of group B streptococcus in

183            pregnant women from Oxford, UK. J Clin Pathol **59:**363-366.

184    17.    **Davies HD, Jones N, Whittam TS, Elsayed S, Bisharat N, Baker CJ.** 2004. Multilocus sequence

185            typing of serotype III group B streptococcus and correlation with pathogenic potential. J

186            Infect Dis **189:**1097-1102.

187    18.    **Gladman S, Seeman T.** 2012. VelvetOptimiser.

188            http://bioinformatics.net.au/software.velvetoptimiser.shtml

189    19.    **Zerbino DR, Birney E.** 2008. Velvet: algorithms for de novo short read assembly using de

190            Bruijn graphs. Genome Res **18:**821-829.

191    20.    **Yamamoto S, Miyake K, Koike Y, Watanabe M, Machida Y, Ohta M, Iijima S.** 1999.

192            Molecular characterization of type-specific capsular polysaccharide biosynthesis genes of

193            Streptococcus agalactiae type Ia. J Bacteriol **181:**5176-5184.

194    21.    **Watanabe M, Miyake K, Yanae K, Kataoka Y, Koizumi S, Endo T, Ozaki A, Iijima S.** 2002.

195            Molecular characterization of a novel beta1,3-galactosyltransferase for capsular

196            polysaccharide synthesis by Streptococcus agalactiae type Ib. J Biochem **131:**183-191.

197    22.    **Martins ER, Melo-Cristino J, Ramirez M.** 2007. Reevaluating the serotype II capsular locus of

198            Streptococcus agalactiae. J Clin Microbiol **45:**3384-3386.

199    23.    **Chaffin DO, Beres SB, Yim HH, Rubens CE.** 2000. The serotype of type Ia and III group B

200            streptococci is determined by the polymerase gene within the polycistronic capsule operon.

201            J Bacteriol **182:**4466-4477.

202    24.    **Cieslewicz MJ, Chaffin D, Glusman G, Kasper D, Madan A, Rodrigues S, Fahey J, Wessels

203            MR, Rubens CE.** 2005. Structural and genetic diversity of group B streptococcus capsular

204            polysaccharides. Infect Immun **73:**3096-3103.

205

206  **Table 1.** Reference sequences used for sequence-based serotype allocation

| Serotype | Accession | Region | Reference |
|---|---|---|---|
| Ia | AB028896.2 | 6982-11695 | Yamamoto et al.(20) |
| Ib | AB050723.1 | 2264-6880 | Watanabe et al.(21) |
| II | EF990365.1 | 1915-8221 | Martins et al.(22) |
| III | AF163833.1 | 6592-11193 | Chaffin et al.(23) |
| IV | AF355776.1 | 6417-11656 | Cieslewicz et al.(24) |
| V | AF349539.1 | 6400-12547 | Cieslewicz et al.(24) |
| VI | AF337958.1 | 6437-10913 | Cieslewicz et al.(24) |
| VII | AY376403.1 | 3403-8666 | Cieslewicz et al.(24) |
| VIII | AY375363.1 | 2971-7340 | Cieslewicz et al.(24) |
| IX | NA | NA | This study |

207

208  **Table 2.** Serotype allocation by WGS to serotype allocation by latex agglutination

| | | Serotype allocated by WGS | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Ia | Ib | II | III | IV | V | VI | VII | VIII | IX | Total |
| | Ia | **34** | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | **35** |
| | Ib | 0 | **9** | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | **10** |
| | II | 0 | 0 | **25** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | **25** |
| | III | 3 | 0 | 0 | **111** | 0 | 0 | 0 | 0 | 0 | 1 | **115** |
| | IV | 0 | 0 | 0 | 0 | **1** | 0 | 1 | 0 | 0 | 0 | **2** |
| Serotype by latex agglutination | V | 0 | 0 | 0 | 0 | 0 | **16** | 0 | 0 | 0 | 0 | **16** |
| | VI | 0 | 0 | 0 | 0 | 0 | 1 | **8** | 0 | 0 | 0 | **9** |
| | VII | 0 | 0 | 0 | 0 | 0 | 0 | 0 | **5** | 0 | 0 | **5** |
| | VIII | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | **1*** | 0 | **1** |
| | IX | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | **1*** | **2** |
| | Non-typeable | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | **3** |
| | Total | 37 | 11 | 26 | 112 | 1 | 17 | 10 | 6 | 1 | 2 | **223** |

209  *Reference GBS isolates from Statens Serum Institute serotypes VIII and IX

210

211    **Table 3.** Relationship between allocated serotype and second-best match (see also Figure 1)

| Allocated serotype | % match | Second-best serotype | % match |
|---|---|---|---|
| Ia | 93.91-100 | III | 64.56 |
| III | 100 | Ia | 62.98 |
| V | 100 | IX | 36.26 |
| IX | 100 | V | 31.05 |
| VI | 100 | III | 26.68 |
| IV | 100 | Ia | 20.3 |
| Ib | 99.61-100 | VI | 15.55 |
| II | 99.86-100 | IV | 9.45 |
| VII | 100 | Ib | 6.95 |
| VIII | 100 | none | 0 |

212

213

214     **Table 4.** Retyping of discordant and non-typable isolates

| | | Latex agglutination | | WGS | |
|---|---|---|---|---|---|
| **Isolate** | **Reason for retyping** | **Initial** | **Repeat** | **Initial** | **Repeat** |
| CB466 | Discordant | III | Ia | Ia | Ia |
| IW8194 | Discordant | III | IX | IX | IX |
| IW8466 | Discordant | Ia | III | III | III |
| IW8471 | Discordant | III | Ia | Ia | Ia |
| IW7157 | Discordant | Ib | II | II | II |
| SMRU1 | Discordant | VI | V | V | V |
| SMRU25 | Discordant | IV | NT | VI | VI |
| SMRU4 | Discordant | IX | NT | Ib | Ib |
| SMRU59 | Discordant | III | NT | Ia | Ia |
| Z41 | Non-typeable | NT | NT | Ib | Ib |
| UK22 | Non-typeable | NT | NT | VII | VII |
| IW2723 | Non-typeable | NT | NT | VI | VI |
| CB454 | Control | III | III | III | III |
| IW4445 | Control | Ia | Ia | Ia | Ia |
| IW4077 | Control | II | II | II | II |

215

216

**Figure 1 Discordant isolates show high support for sequence-based serotype allocation.**

For each isolate, the percentage of the capsular locus region present (≥95% sequence

identity) for the assigned serotype is shown on the X axis, and that for the serotype showing

the next best match on the Y axis. Isolates showing agreement between sequence-based

and conventional serotyping are shown in grey, those classified as non-typeable by

conventional methods in blue, and discordant isolates in red. Small circles represent single

isolates, the large circle represents 100 isolates. For each serotype, the second-best match is

identical in all cases, leading to the observed horizontal banding (details in Table 3).