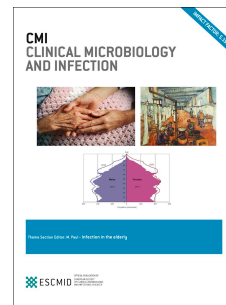


# Journal Pre-proof

“The recent emergence of a highly related virulent *Clostridium difficile* clade with unique characteristics”

Helen Alexandra Shaw, Mark D. Preston, Karuna E.W. Vendrik, Michelle D. Cairns, Hilary P. Browne, Richard A. Stabler, Monique J.T. Crobach, Jeroen Corver, Hanna Pituch, Andre Ingebretsen, Munir Primohammed, Alexandra Faulds-Pain, Esmeralda Valiente, Trevor D. Lawley, Neil F. Fairweather, Ed J. Kuijper, Brendan W. Wren



PII: S1198-743X(19)30489-6

DOI: <https://doi.org/10.1016/j.cmi.2019.09.004>

Reference: CMI 1775

To appear in: *Clinical Microbiology and Infection*

Received Date: 2 June 2019

Revised Date: 6 September 2019

Accepted Date: 7 September 2019

Please cite this article as: Shaw HA, Preston MD, Vendrik KEW, Cairns MD, Browne HP, Stabler RA, Crobach MJT, Corver J, Pituch H, Ingebretsen A, Primohammed M, Faulds-Pain A, Valiente E, Lawley TD, Fairweather NF, Kuijper EJ, Wren BW, “The recent emergence of a highly related virulent *Clostridium difficile* clade with unique characteristics”, *Clinical Microbiology and Infection*, <https://doi.org/10.1016/j.cmi.2019.09.004>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2019 Published by Elsevier Ltd on behalf of European Society of Clinical Microbiology and Infectious Diseases.

1 ORIGINAL ARTICLE

2 Title: "The recent emergence of a highly related virulent *Clostridium difficile* clade with  
3 unique characteristics"

4  
5 Helen Alexandra Shaw<sup>1,2</sup>, Mark D. Preston<sup>1,3</sup>, Karuna E. W. Vendrik<sup>4</sup>, Michelle D. Cairns<sup>1,5</sup>, Hilary P.  
6 Browne<sup>6</sup>, Richard A. Stabler<sup>1</sup>, Monique J. T. Crobach<sup>4</sup>, Jeroen Corver<sup>4</sup>, Hanna Pituch<sup>7</sup>, Andre  
7 Ingebretsen<sup>8,9</sup>, Munir Primohammed<sup>10</sup>, Alexandra Faulds-Pain<sup>1</sup>, Esmeralda Valiente<sup>1</sup>, Trevor D.  
8 Lawley<sup>6</sup>, Neil F. Fairweather<sup>11</sup>, Ed J. Kuijper<sup>4</sup> & Brendan W. Wren<sup>1\*</sup>

9  
10 <sup>1</sup> Department of Pathogen Molecular Biology, London School of Hygiene and Tropical Medicine,  
11 London WC1E 7HT, UK

12 <sup>2</sup> Division of Bacteriology, National Institute for Biological Standards and Controls (NIBSC), Blanche  
13 Lane, South Mimms, Potters Bar, Hertfordshire, EN6 3QG, UK

14 <sup>3</sup> Analytical Biological Service Division, National Institute for Biological Standards and Controls  
15 (NIBSC), Blanche Lane, South Mimms, Potters Bar, Hertfordshire, EN6 3QG, UK

16 <sup>4</sup> National Reference Laboratory for CDI surveillance, Department of Medical Microbiology and RIVM,  
17 Leiden University Medical Centre, Leiden, the Netherlands

18 <sup>5</sup> Public Health Laboratory London, Division of Infection, The Royal London Hospital, London, E1 2ES,  
19 UK

20 <sup>6</sup> Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire, CB10 1SA, UK

21 <sup>7</sup> Department of Medical Microbiology, Medical University of Warsaw, Warsaw, Poland

22 <sup>8</sup> Department of Microbiology, Oslo University Hospital, Oslo, Norway

23 <sup>9</sup> Department of Infection Prevention, Oslo University Hospital, Oslo, Norway

24 <sup>10</sup> Department of Molecular and Clinical Pharmacology, The University of Liverpool, Liverpool, L69

25 3GL, UK

26 <sup>11</sup> CMBI, Imperial College London, South Kensington Campus, London, SW7 2AZ, UK

27

28 \* Correspondence: Professor Brendan Wren, London School of Department of Pathogen Molecular  
29 Biology, London School of Hygiene and Tropical Medicine, London WC1E 7HT, UK

30 Telephone: +44 (0) 20 7927 2288; [brendan.wren@lshtm.ac.uk](mailto:brendan.wren@lshtm.ac.uk)

31

32 Keywords: NGS, *Clostridium difficile*, CDI, RT023

33

34 Abstract word count: 238

35 Text word count: 2638

36

37

38

39

40

41

42

43

44 ABSTRACT:

45 **Objectives**

46 *Clostridium difficile* is a major global human pathogen divided into five clades, of which clade 3 is the  
47 least characterised and consists predominantly of PCR ribotype (RT) 023 strains. Our aim was to  
48 analyse and characterise this clade.

49 **Methods**

50 In this cohort study the clinical presentation of *C. difficile* RT023 infections was analysed in  
51 comparison with known “hypervirulent” and non-hypervirulent strains, using data from the  
52 Netherlands national *C. difficile* surveillance programme. European RT023 strains of diverse origin  
53 were collected and whole-genome sequenced to determine the genetic similarity between isolates.  
54 Distinctive features were investigated and characterised.

55 **Results**

56 Clinical presentation of *C. difficile* RT023 infections show severe infections akin to those seen with  
57 “hypervirulent” strains from clades 2 (RT027) and 5 (RT078) (35%, 29% and 27% severe CDI  
58 respectively), particularly with significantly more bloody diarrhoea than RT078 and non-  
59 hypervirulent strains (RT023 8%, other RTs 4%,  $p=0.036$ ). The full genome sequence of strain CD305  
60 is presented as a robust reference. Phylogenetic comparison of CD305 and a further 79 previously  
61 uncharacterised European RT023 strains of diverse origin revealed minor genetic divergence with  
62 >99.8% pairwise identity between strains. Analyses revealed distinctive features among clade 3  
63 strains, including conserved PaLoc, CDT and phage insertion toxin genotypes, glycosylation of S-layer  
64 proteins, presence of the RT078 four gene trehalose cluster and an esculinase negative genotype.

65 **Conclusions**

66 Given their recent emergence, virulence and genomic characteristics, the surveillance of clade 3  
67 strains should be more highly prioritised.

68

## 69 INTRODUCTION

70 *C. difficile* remains a major global pathogen; disease severity and relapse incidence have not abated,  
71 and community acquired infections have increased (1). *C. difficile* can be divided into five clades of  
72 virulent strains (2). The most understudied is clade 3, dominated by PCR ribotype (RT) 023 strains  
73 (2). RT023 has been reported primarily in Europe (3) and is amongst the top ten most common *C.*  
74 *difficile* PCR ribotypes in England (4) (CDRN report 2013-2015) and the Netherlands (unpublished  
75 data of the Dutch *C. difficile* Reference Laboratory). RT023 infections are not associated with  
76 increased mortality despite causing a high level of deleterious biomarkers (e.g. neutrophil counts) in  
77 patients and having toxin profiles similar to clade 2 (RT027) and clade 5 (RT078) strains (5, 6).  
78 However, disease severity with RT023 has been reported as similar to “hypervirulent strains”,  
79 particularly in elderly patients (7), and is frequently associated with a relapse of CDI (3).

80 This study investigates the clinical presentation and phylogeny of *C. difficile* clade 3, uncovering and  
81 characterising unique features of these strains.

82

## 83 METHODS

84 **Clinical data collection and analysis**

85 A cohort study was performed. Clinical data from the Dutch national CDI sentinel surveillance from  
86 May 2009 until February 2018 were used to analyse the clinical characteristics of CDI episodes due  
87 to RT023. For this sentinel surveillance all hospitalized patients >2 years old, with clinical signs or  
88 symptoms of CDI in combination with a positive test for *C. difficile* toxins or toxigenic *C. difficile*, in  
89 Dutch participating hospitals, are registered. The indication for testing on CDI and the assay or  
90 algorithm that is used to diagnose CDI is chosen by the local laboratory.  
91 Using classification criteria based on expert opinion that were previously used (8), CDI is classified as  
92 severe if one or more of the following conditions were present; fever (temperature of 38°C or

93 higher) and leucocytosis ( $>15 \times 10^9/L$ ), diarrhoea with hypoalbuminemia ( $<20 \text{ g/L}$ ) and/or  
94 dehydration, pseudomembranous colitis and/or bloody diarrhoea. A complicated course is defined  
95 as the need for surgical procedure, admission to intensive care unit and/or mortality (CDI- or non-  
96 CDI-related) within 30 days after CDI diagnosis (8).

97 Our primary aim was to test the null hypothesis that RT023 causes the same proportion of severe  
98 CDI as non-hypervirulent ribotypes. Therefore, clinical characteristics and 30-day outcome of CDI  
99 episodes due to RT023 were compared to CDI episodes due to other ribotypes (excluding  
100 hypervirulent strains RT027 and RT078/126). Thereafter, the results of the RT023-group were  
101 compared to the results of 4 pre-specified groups; RT027 and RT078/126, which are well-known  
102 hypervirulent strains, and RT001 and RT014/020/295, which are non-hypervirulent strains that are  
103 common in the Netherlands. Each time, results of the RT023-group were compared with the results  
104 of one other group. Some ribotypes were merged into one group since they are hard to distinguish  
105 with PCR ribotyping. Further details are in the web-only Supplementary Material.

106 Data are presented as number of cases (percentage). Age is presented as media [first quartile, third  
107 quartile], because of the skewed distribution. Categorical variables were compared by a Pearson's  
108 Chi square test and numerical variables were compared by a Wilcoxon rank-sum test. To identify the  
109 effect of RT023 on CDI severity, a multivariable logistic regression analysis was performed with age  
110 and sex as covariates. A p-value of  $<0.05$  was considered statistically significant. STATA SE version  
111 12.1 statistical software (StataCorp, Texas, USA) was used for statistical analysis.

112

### 113 **Ethics**

114 This was an observational study, using data that are already collected in the Dutch national CDI  
115 surveillance. This national surveillance program exists since 2009 and collects microbiological and  
116 clinical data from all hospitalized patients with CDI in the participating hospitals in the Netherlands.  
117 The surveillance has been developed by our National Institute of Public Health. There were no

118 additional data or isolates/materials specifically for this study collected and no actions were requested  
119 from patients.

120

### 121 **Whole-genome sequencing**

122 CD305 genomic DNA was sequenced using 454 pyrosequencing (GS-FLX pyrosequencing) to generate  
123 3 kb paired-end libraries and Illumina GAII paired-end libraries of 400 bp insert size and 108 bp read  
124 length. The resulting sequence was assembled using Newbler and Velvet and the assemblies were  
125 combined using Newbler (9, 10). CDS identification and annotation was generated using PROKKA  
126 (11) with a bespoke *C. difficile* library. The assembled and annotated genome is available at  
127 ERS2502454. For 79 study isolates genomic DNA libraries were created using a Nextera XT kit  
128 (Illumina, CA, USA) and data obtained using the MiSeq sequencing system (Illumina, CA, USA).

129

### 130 **Whole-genome bioinformatics analysis**

131 The sequence data were processed according to a standard protocol as previously described (12)  
132 (Detail in web-only Supplementary Materials). SNP loci were identified with a samtools Q-score  $\geq$   
133 30, coverage  $\geq$  10 and 80% of contributing reads. Pipeline, phylogenetic and post-analyses were  
134 carried out using Perl, R and RAxML (13).

135

### 136 **Glycoprotein detection**

137 Glycosylated proteins were detected using Pierce™ Glycoprotein Staining Kit according to the  
138 manufacturer's instructions (Detail in web-only Supplementary Material).

139

140 RESULTS

141 **CDI in hospitalised patients due to RT023 strains is severe comparable with RT027 and RT078**  
142 **strains**

143 Between May 2009 and February 2018, 5359 samples from hospitalised patients in twenty-four  
144 hospitals in the Netherlands were PCR-ribotyped within the context of the national *C. difficile*  
145 surveillance program. Clinical data were complete in 4387 cases. RT023 accounted for 141 cases of  
146 CDI, a mean proportion of 2.4% (95% CI 2.0-2.8), which remained consistent within the study period.  
147 Demographic data, clinical characteristics and 30-day outcome of patients with CDI due to RT023  
148 were compared to data of five other pre-specified ribotype groups, shown in Table 1. There were no  
149 significant differences in age and sex between the RT023 group and the other groups, except for  
150 higher age in the RT001 group.

151 The primary question was whether CDI due to RT023 was more severe when compared to all non-  
152 hypervirulent ribotypes, which was confirmed by our results ( $p=0.000$ : 35% (27-44) vs 22% (21-23)),  
153 also after correcting for sex and age. No significant differences of severity were found when RT023  
154 was compared to “hypervirulent” strains RT027 and RT078/126 ( $p=0.310$  and  $p=0.065$  respectively,  
155 RT023: 35% (27-44), RT027: 29% (20-38), RT078/126: 27% (24-31)), also not after correction for sex  
156 and age. Of note, bloody diarrhoea was more frequently reported in RT023 infections compared  
157 with RT078/126 infections ( $p=0.031$ ), RT014/020/295 ( $p=0.036$ ) and RT001 ( $p=0.037$ ) (RT023: 8% (3-  
158 13), RT078/126, 4% (2-5), RT014/020/295 4% (3-5), RT001 4% (2-5)). When compared to non-  
159 hypervirulent RT001 and RT014/020/295 isolates, with or without correcting for sex and age, RT023  
160 presented with significantly more severe symptoms ( $p=0.000$  for both, RT023: 35% (27-44), RT001:  
161 16% (13-19), RT014/020/295: 21% (18-23)), such as more frequent diarrhoea with dehydration  
162 and/or hypoalbuminemia. However, the outcomes of CDI due to RT023 in terms of a complicated  
163 course, including mortality, were comparable with outcomes of CDI due to RT001, RT014/020/295  
164 and all non-hypervirulent ribotypes. RT027 and RT078/126 infections showed higher overall  
165 mortality than RT023 ( $p=0.032$ ,  $p=0.049$  respectively, RT023: 9% (4-14), RT027: 19% (11-27),



166 RT078/126: 16% (13-19)) but CDI attributable mortality was similar between these groups ( $p=0.293$ ,  
167  $p=0.152$  respectively, RT023: 2%(-1-4), RT027: 4% (0-8), RT078/126: 5%(3-7)). There were  
168 significantly more complicated courses in patients with CDI due to RT027 compared to RT023  
169 ( $p=0.038$ , 23% (14-31) vs 12% (6-18) respectively), but no significant differences were observed  
170 between RT078/126 and RT023 ( $p=0.144$ , RT078/126 17% (14-20)).

171 Comparison of RT023 with all groups in this study revealed that the onset of symptoms of CDI due to  
172 RT023 was more frequently at home and less often in healthcare facilities ( $p=0.000$  compared to all  
173 other groups). Subgroup analysis of community and hospital onset CDI can be found in the web-only  
174 Supplementary Material. The number of episodes that were recurrences of a previous CDI episode 2-  
175 8 weeks earlier was the same in RT023 episodes compared to all other groups (Table 1).

176

### 177 **Clade 3 strains are highly related**

178 A high-quality (14) draft genome of strain CD305 (RT023) was generated and is presented here as a  
179 robust reference for this lineage. Further strains were sourced from across Europe (Supplementary  
180 Table S1), with this study comprising 86 strains: CD305 (reference); 79 (out of 170 WGS strains); and  
181 6 published clade 3 strains (15, 16) (Supplementary Table S2), the largest RT023 genomic collection.  
182 MLST were identified *in silico* from *de novo* assemblies. The six published strains matched their  
183 published MLST with new strains composed of 68 ST005, 10 ST022, and one novel ST (strain  
184 OUS23024) (Figure 1).

185 The 79 core strains were aligned to the CD305 reference strain and a set of 19,262 (<0.5% of the 4.2  
186 Mbp genome) high quality SNP loci identified. The individual strains were very closely related with  
187 only between 58 and 7,876 pair-wise SNP differences, with a mean of 1,767 SNPs (mean: 9.2% of  
188 19262 SNPs; max: 40.9%) equating to >99.8% pairwise identity between strains. A phylogeny was  
189 created from all 86 strain's SNPs that reinforces the conclusion of little genetic diversity within clade

190 3 strains (Figure 1). From our 80 strains there are two outliers: strains 91 and 108698, which are not  
191 RT023 (Figure 1A, Supplementary Figure S3, Supplementary Text). The unassigned MLST strain  
192 (OUS23024) diverged slightly from the main population (Figure 1B). No significant relationship was  
193 found with any phenotypes including the infection date (2007-2014) or geographic origin  
194 (Supplementary Table S1, Supplementary Figure S1). Detail on MLST and ribotype divergence can be  
195 found in the web-only Supplementary Material.

196 There is high conservation in all 86 strains of larger clade-specific genetic features such as the  
197 pathogenicity locus (PaLoc), binary toxin CDT, PaLoc phage insertion and type B flagella glycosylation  
198 cluster (Supplementary Tables S2 and S3). The only common antibiotic resistance marker is *gyrB*  
199 (V426D) related to fluoroquinolone resistance. Analysis of twelve Polish RT023 strains for  
200 fluoroquinolone resistance revealed resistance to ciprofloxacin but sensitivity to moxifloxacin  
201 (Supplementary Table S4).

202

### 203 **A unique trehalose metabolism genotype is present in clade 3 strains**

204 Analysis of clade 3 strains for two trehalose clusters described to be important in global  
205 dissemination and virulence of *C. difficile* (17) showed a trehalose genotype unique to these strains.  
206 The primary cluster, in which SNP L172I defines increased metabolism in RT027 (clade 2) (Figure 2a),  
207 was absent from all clade 3 genomes analysed. This coincides with polymorphisms and a large  
208 deletion in sugar metabolism genes in clade 3, including beta-glucosidase genes (Supplementary  
209 Text). However, the RT078 (clade 5) second cluster (Figure 2b) was observed in all strains.  
210 Polymorphisms exist between the RT078 cluster in M120 cluster and RT023 CD305, with the most  
211 significant difference being a truncation of *treX* (Figure 2c). Between clade 3 strains there are only a  
212 small number of SNPs, predominantly in strain 91 (Supplementary Text).

213

**214 Clade 3 have a glycosylated surface**

215 SlpA is the major surface protein of *C. difficile* comprised of high and low molecular proteins (HMW  
216 and LMW SLP) (18). A putative glycosylation cluster within the *slp* gene island (Figure 3a) for S-layer  
217 cassette type 11, SLCT11 (18), has been previously reported (19). 83 of the 86 strains contain this  
218 feature (Supplementary Table S2, Supplementary Figure S2). Strains 91, Ox2183 and WCHCD103  
219 from which this feature is absent are genetically distinct from other strains within this clade, with  
220 alternate *slpA* genes. In RT023 the *slpA* gene encodes a smaller LMW SLP than in other clades,  
221 predicted at approximately 18 kDa (Figure 3b). S-layer extracts of representative strains from each of  
222 the five clades of *C. difficile* show two distinct bands of equimolar ratio representing the HMW and  
223 LMW SLPs in clades 1, 2, 4 and 5 by Coomassie brilliant blue staining (Figure 3c). Strain Ox247  
224 (RT005, clade 1) containing SLCT11 (20) along with S-layer preparations from three representative  
225 RT023 strains show an alternative pattern of SLPs. HMW SLP migrates at its expected molecular  
226 weight, but a band at 18 kDa for LMW SLP is absent. A periodic acid-Schiff assay to stain for glycans  
227 on S-layer preparations showed glycosylated proteins at ~45 kDa only in strains containing the  
228 glycosylation cluster, demonstrating the presumed functionality of the cluster and glycosylation of S-  
229 layer proteins.

230

**231 DISCUSSION**

232 This study provides a comprehensive analysis of clade 3 strains of *C. difficile* with an extensive report  
233 of RT023 CDI and detailed WGS analysis. The clinical characteristics of hospitalised patients with CDI  
234 due to RT023 showed CDI severity similar to the “hypervirulent” RT027 and RT078/126, with  
235 comparable CDI-related mortality, though overall mortality was lower in RT023 as previously  
236 reported (6). The phylogeny of clade 3 strains is compact, barring six distinct outliers. In contrast to  
237 clade 2 strains (RT027), clade 3 strains show great similarity consistent with a recently emerged  
238 clade under little selective pressure to evolve (21). WGS analysis revealed a unique trehalose

239 genotype and conserved incorporation of a glycosylation cassette into the clade 3 genomes which  
240 was demonstrated to glycosylate the S-layer.

241 Considering previous investigations, the severity of disease is likely due to the production of binary  
242 toxin and the TcdC stop codon in RT023 (5). Recurrent infections due to RT023 were similar to other  
243 ribotypes. This contrasts with an earlier study, where RT023 was dominating among recurrent cases  
244 (22). We also observed more community acquisition of RT023 symptoms, but current reports cannot  
245 explain this observation. Circulating strains unlikely to be the source of RT023 with no  
246 representation of RT023 in a small group of *C. difficile* carriers (23) and a low representation in *C.*  
247 *difficile* infections in the community (24). The low proportion (2.4%) of CDI due to RT023 observed in  
248 this study in the Netherlands is consistent with a previous study on CDI in Europe (3).

249 Strengths of this study are the high sample size, multicenter design with high number of hospitals in  
250 different geographic regions, and 10 years of available data, making the data generalizable for  
251 hospitalized patients. Similarly, a sample size of over 80 strains across 8 years from a variety of pan-  
252 European sources for WGS, as well as published strains including Chinese strains, enabled us to  
253 understand the phylogeny of clade 3 in much greater detail. Limitations of the clinical data include  
254 the location of symptoms onset being documented but not the location of *C. difficile* acquisition.  
255 Furthermore, there was no data available regarding comorbidity, which might affect the outcome.  
256 Regarding severity of disease, occasionally not all laboratory parameters needing lab results were  
257 measured and included.

258 It has recently been shown that S-layer glycosylation is important for adherence to Caco-2 intestinal  
259 epithelial cells but not biofilm formation (20). Therefore, glycosylation of the S-layer in clade 3 may  
260 be important for colonisation but not persistence, explaining a low level of carriage and recurrence  
261 of these strains. Despite severe clinical presentation this clade is not as widely disseminated as other  
262 clades. The emergence of RT027 and RT078 strains has been linked to an increased ability to  
263 metabolise the food additive trehalose (17). RT023 strains contain the second four gene cluster,

264 corroborated by a recent study of trehalose genes in all clades of *C. difficile*. The presence of only the  
265 secondary cluster and the SNPs between RT023 and RT078 may result in a difference in uptake and  
266 metabolism of trehalose between these strains, which could explain the relatively reduced  
267 prevalence of RT023 strains compared with RT078 and RT027 strains globally. No link between  
268 trehalose and adverse disease outcomes has been suggested (25). Meanwhile, the emergence of  
269 epidemic clade 2 strains has also been linked to environmental spore contamination and the  
270 acquisition of fluoroquinolone resistance, which is less pronounced for clade 3 strains (21). More  
271 analysis on sporulation in clade 3 is required as reduced sporulation efficiency and survival outside  
272 the human host has been reported (26), however, a recent study highlighted a clade 3 strain in China  
273 which had a high sporulation and germination rate (27).

274 It remains to be determined why evolutionary distinct clades of *C. difficile* are emerging  
275 simultaneously to cause disease in human populations, or if *C. difficile* is evolving into subspecies  
276 (28). Our study suggests that a heightened awareness and continued surveillance of RT023 strains  
277 globally should be a current imperative.

278

#### 279 DATA AVAILABILITY

280 Sequence data that supports the findings of this study have been deposited in EMBL Nucleotide  
281 Sequence Database (ENA) with accession code PRJEB26893 and CD305 reference genome  
282 ERS2502454.

283

#### 284 TRANSPARENCY DECLARATION

285 The authors declare no conflicts of interest. The work was supported by The Wellcome Trust (Grant  
286 Reference 102979/Z/13/Z and 098051) and the Medical Research Council (Grant Reference  
287 MR/K000551/1).

288

## 289 ACKNOWLEDGEMENTS

290 We thank Ed Kuijper, Andre Ingebretsen, Hanna Pituch, Munir Primohammed, Paul Roberts (Royal  
291 Liverpool Hospital) and Neil Fairweather for the supply of strains to this study. We acknowledge the  
292 London *C. difficile* Ribotyping Laboratory for help with PCR ribotyping and supply of strains and Dr.  
293 Piotr Obuch-Woszczatynski for assistance with Polish RT023 strains. We thank all laboratories for  
294 helping to collect the data for the National Surveillance program in the Netherlands.

295

## 296 AUTHOR CONTRIBUTIONS

297 Concept and design of study: H.A.S., M.D.C. and B.W.W. Genomic assembly and annotation: M.D.P.,  
298 H.P.B. and R.A.S. Genomic analysis: H.A.S. and M.D.P. Phenotypic experiments: H.A.S.  
299 Fluoroquinolone testing: H.P. Clinical analysis: K.E.W.V., M.J.T.C and E.J.K. The manuscript was  
300 drafted by H.A.S, M.D.P., K.E.W.V. and B.W.W., and revised by all authors.

301

## 302 REFERENCES

303

- 304 1. Khanna S, Pardi DS, Aronson SL, Kammer PP, Orenstein R, St Sauver JL, et al. The  
305 epidemiology of community-acquired *Clostridium difficile* infection: a population-based study. The  
306 American journal of gastroenterology. 2012;107(1):89-95.
- 307 2. Stabler RA, Dawson LF, Valiente E, Cairns MD, Martin MJ, Donahue EH, et al. Macro and  
308 micro diversity of *Clostridium difficile* isolates from diverse sources and geographical locations. PloS  
309 one. 2012;7(3):e31559.
- 310 3. Bauer MP, Notermans DW, van Benthem BH, Brazier JS, Wilcox MH, Rupnik M, et al.  
311 *Clostridium difficile* infection in Europe: a hospital-based survey. Lancet (London, England).  
312 2011;377(9759):63-73.
- 313 4. Wilcox MH, Shetty N, Fawley WN, Shemko M, Coen P, Birtles A, et al. Changing epidemiology  
314 of *Clostridium difficile* infection following the introduction of a national ribotyping-based  
315 surveillance scheme in England. Clinical infectious diseases : an official publication of the Infectious  
316 Diseases Society of America. 2012;55(8):1056-63.
- 317 5. Dingle KE, Griffiths D, Didelot X, Evans J, Vaughan A, Kachrimanidou M, et al. Clinical  
318 *Clostridium difficile*: clonality and pathogenicity locus diversity. PloS one. 2011;6(5):e19993.

- 319 6. Walker AS, Eyre DW, Wyllie DH, Dingle KE, Griffiths D, Shine B, et al. Relationship between  
320 bacterial strain type, host biomarkers, and mortality in *Clostridium difficile* infection. *Clinical*  
321 *infectious diseases : an official publication of the Infectious Diseases Society of America.*  
322 2013;56(11):1589-600.
- 323 7. Vanek J, Hill K, Collins J, Berrington A, Perry J, Inns T, et al. Epidemiological survey of  
324 *Clostridium difficile* ribotypes in the North East of England during an 18-month period. *The Journal of*  
325 *hospital infection.* 2012;81(3):209-12.
- 326 8. Goorhuis A, Bakker D, Corver J, Debast SB, Harmanus C, Notermans DW, et al. Emergence of  
327 *Clostridium difficile* infection due to a new hypervirulent strain, polymerase chain reaction ribotype  
328 078. *Clinical infectious diseases : an official publication of the Infectious Diseases Society of America.*  
329 2008;47(9):1162-70.
- 330 9. Bonfield JK, Smith K, Staden R. A new DNA sequence assembly program. *Nucleic acids*  
331 *research.* 1995;23(24):4992-9.
- 332 10. Zerbino DR. Using the Velvet de novo assembler for short-read sequencing technologies.  
333 *Current protocols in bioinformatics.* 2010;Chapter 11:Unit 11.5.
- 334 11. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics (Oxford, England).*  
335 2014;30(14):2068-9.
- 336 12. Cairns MD, Preston MD, Hall CL, Gerding DN, Hawkey PM, Kato H, et al. Comparative  
337 Genome Analysis and Global Phylogeny of the Toxin Variant *Clostridium difficile* PCR Ribotype 017  
338 Reveals the Evolution of Two Independent Sublineages. *Journal of clinical microbiology.*  
339 2017;55(3):865-76.
- 340 13. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
341 phylogenies. *Bioinformatics (Oxford, England).* 2014;30(9):1312-3.
- 342 14. Chain PS, Grafham DV, Fulton RS, Fitzgerald MG, Hostetler J, Muzny D, et al. Genomics.  
343 Genome project standards in a new era of sequencing. *Science (New York, NY).* 2009;326(5950):236-  
344 7.
- 345 15. Dingle KE, Elliott B, Robinson E, Griffiths D, Eyre DW, Stoesser N, et al. Evolutionary history of  
346 the *Clostridium difficile* pathogenicity locus. *Genome biology and evolution.* 2014;6(1):36-52.
- 347 16. Chen R, Feng Y, Wang X, Yang J, Zhang X, Lu X, et al. Whole genome sequences of three  
348 Clade 3 *Clostridium difficile* strains carrying binary toxin genes in China. *Scientific reports.*  
349 2017;7:43555.
- 350 17. Collins J, Robinson C, Danhof H, Knetsch CW, van Leeuwen HC, Lawley TD, et al. Dietary  
351 trehalose enhances virulence of epidemic *Clostridium difficile*. *Nature.* 2018.
- 352 18. Kirk JA, Banerji O, Fagan RP. Characteristics of the *Clostridium difficile* cell envelope and its  
353 importance in therapeutics. *Microbial biotechnology.* 2017;10(1):76-90.
- 354 19. Dingle KE, Didelot X, Ansari MA, Eyre DW, Vaughan A, Griffiths D, et al. Recombinational  
355 switching of the *Clostridium difficile* S-layer and a novel glycosylation gene cluster revealed by large-  
356 scale whole-genome sequencing. *The Journal of infectious diseases.* 2013;207(4):675-86.
- 357 20. Richards E, Bouche L, Panico M, Arbeloa A, Vinogradov E, Morris H, et al. The S-layer protein  
358 of a *Clostridium difficile* SLCT-11 strain displays a complex glycan required for normal cell growth and  
359 morphology. *The Journal of biological chemistry.* 2018;293(47):18123-37.
- 360 21. He M, Miyajima F, Roberts P, Ellison L, Pickard DJ, Martin MJ, et al. Emergence and global  
361 spread of epidemic healthcare-associated *Clostridium difficile*. *Nature genetics.* 2013;45(1):109-13.
- 362 22. Sandell S, Rashid MU, Jorup-Ronstrom C, Ellstrom K, Nord CE, Weintraub A. *Clostridium*  
363 *difficile* recurrences in Stockholm. *Anaerobe.* 2016;38:97-102.
- 364 23. Zomer TP, E VAND, Wielders CCH, Veenman C, Hengeveld P, W VDH, et al. Prevalence and  
365 risk factors for colonization of *Clostridium difficile* among adults living near livestock farms in the  
366 Netherlands. *Epidemiology and infection.* 2017;145(13):2745-9.
- 367 24. Hensgens MP, Dekkers OM, Demeulemeester A, Buiting AG, Bloembergen P, van Benthem  
368 BH, et al. Diarrhoea in general practice: when should a *Clostridium difficile* infection be considered?

- 369 Results of a nested case-control study. *Clinical microbiology and infection : the official publication of*  
370 *the European Society of Clinical Microbiology and Infectious Diseases*. 2014;20(12):O1067-74.
- 371 25. Eyre DW, Didelot X, Buckley AM, Freeman J, Moura IB, Crook DW, et al. *Clostridium difficile*  
372 *trehalose metabolism variants are common and not associated with adverse patient outcomes when*  
373 *variably present in the same lineage*. *EBioMedicine*. 2019;43:347-55.
- 374 26. Connor M, Flynn PB, Fairley DJ, Marks N, Manesiotis P, Graham WG, et al. *Evolutionary clade*  
375 *affects resistance of Clostridium difficile spores to Cold Atmospheric Plasma*. *Scientific reports*.  
376 2017;7:41814.
- 377 27. Li C, Harmanus C, Zhu D, Meng X, Wang S, Duan J, et al. *Characterization of the virulence of a*  
378 *non-RT027, non-RT078 and binary toxin-positive Clostridium difficile strain associated with severe*  
379 *diarrhea*. *Emerg Microbes Infect*. 2018;7(1):211.
- 380 28. Kumar N, Browne HP, Viciani E, Forster SC, Clare S, Harcourt K, et al. *Adaptation of host*  
381 *transmission cycle during Clostridium difficile speciation*. *Nature genetics*. 2019.

382

383



384 **FIGURE LEGENDS**385 **Figure 1: Phylogenetic Tree by MLST**

386 Phylogenetic tree of 86 strains generated from analysis of high-quality SNPs and coloured by MLST.

387 A: full tree, with two cohort outliers (samples 91 and 108676), Ox2183 and three Chinese strains. B:  
388 the large, temporally indistinguishable main cluster, with reference CD305 and novel MLST strain  
389 OUS23024 indicated.

390

391 **Figure 2: Clade 3 show a unique trehalose genotype**

392 Schematic demonstrating the three trehalose metabolism genotypes observed in *C. difficile* with  
393 clade 3 strains lacking the primary trehalose metabolism cluster. A: RT012 630 and RT027 R20291  
394 genotypes of a primary trehalose cluster, with the L172I SNP associated with increased metabolism  
395 of trehalose. B: RT078 M120 genotype with primary and secondary trehalose metabolism gene  
396 clusters observed. C: RT023 CD305 trehalose genotype with only the secondary cluster including a  
397 truncated *treX* gene.

398

399 **Figure 3: Insertion of a glycosylation cluster results in S-layer glycosylation**

400 RT023 contains a glycosylation cluster within the *slp* gene island. A: Genomic organisation of the *slp*  
401 gene island in 630 (Clade 1) and CD305 (Clade 3) showing loss of Cwp2 and acquisition of a gene  
402 cluster comprising putative glycosylation genes (adapted from Kirk *et al* (18)). B: Structure of SlpA in  
403 630 and CD305 showing Cwp84 cleavage sites and truncated LMW (light grey) in CD305. C:  
404 Coomassie staining of S layer protein preparations from representative strains from each clade  
405 showing characteristic double banding for HMW and LMW SLP (grey arrows). D: Periodic acid-Schiff  
406 staining of glycans in S layer preparations.

407

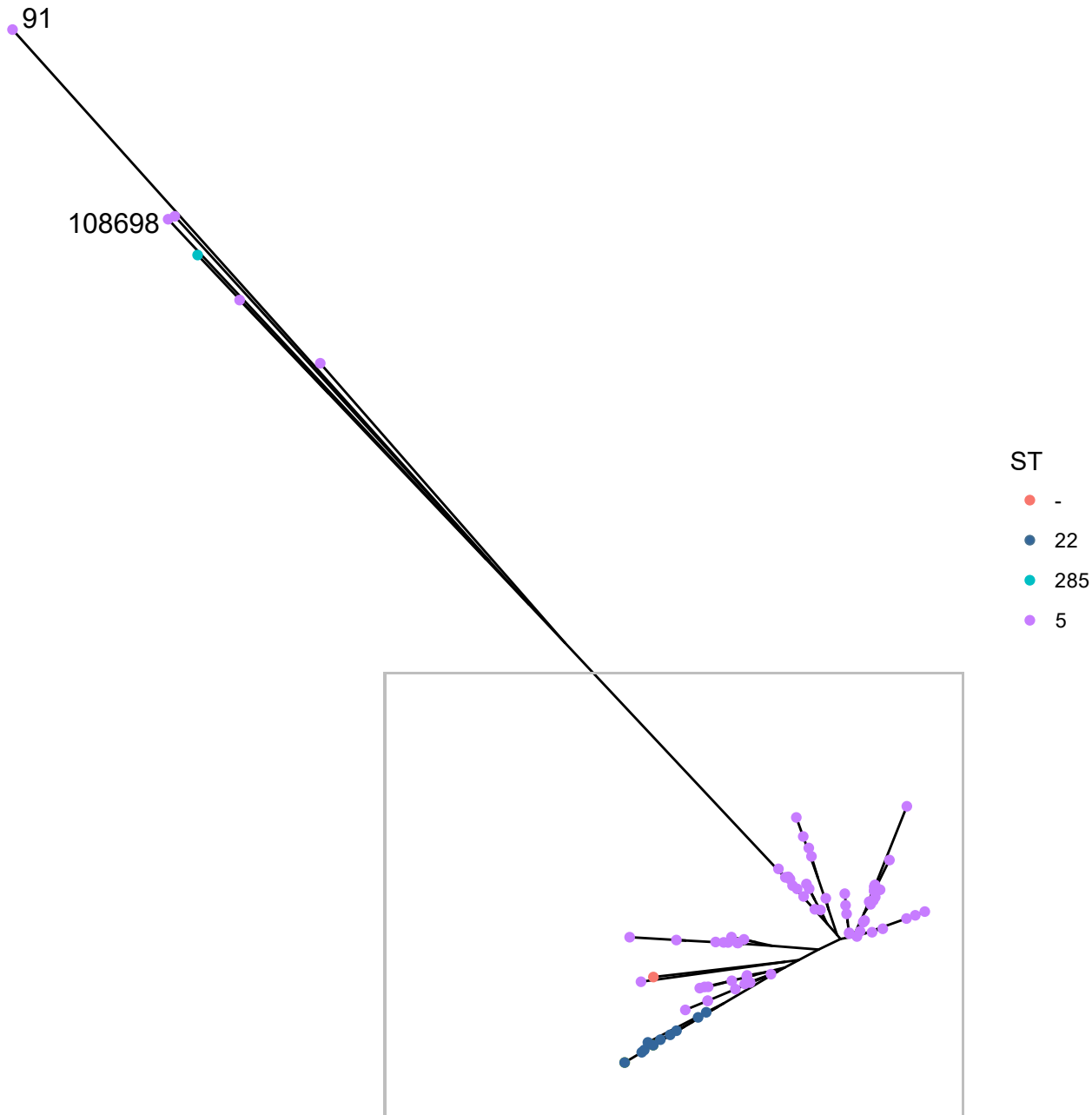
**Table 1. Comparison of clinical characteristics of patients with RT023 versus other ribotypes (excluding RT027 and RT078/126), RT027, RT078/126, RT014/020/295 and RT001.**

	Primary outcome		Hypervirulent strains		Non-hypervirulent strains		All info available
	RT023, n=141	Others, n=4368	RT027, n=116	RT078/126, n=734	RT014/020/295, n=962	RT001, n=699	
<b>Age</b>	71.4 [10.0, 97.7]	71.3 [1.9, 102.3]	73.2 [11.2, 91.5]	70.9 [5.2, 100.7]	70.4 [2.1, 99.2]	76.0 [3.3, 96.7]*	5359/5359
<b>Men</b>	71 (50)	2095 (48)	63 (54)	365 (50)	444 (46)	344 (49)	5356/5359
<b>Severe CDI</b>	45 (35)	880 (22)*	30 (29)	188 (27)	185 (21)*	104 (16)*	4948/5359
<i>Dehydration and/or hypoalbuminemia</i>	25 (20)	450 (11)*	14 (14)	100 (15)	97 (11)*	44 (7)*	4940/5359
<i>Bloody diarrhoea</i>	10 (8)	192 (5)	6 (6)	25 (4)*	34 (4)*	24 (4)*	4948/5359
<i>Pseudomembranous colitis</i>	8 (6)	159 (4)	6 (6)	41 (6)	28 (3)	21 (3)	4948/5359
<i>Fever and leucocytosis</i>	11 (9)	295 (7)	9 (9)	76 (11)	64 (7)	36 (6)	4940/5359
<b>Complicated course</b>	13 (12)	485 (14)	21 (23)*	104 (17)	78 (10)	95 (17)	4387/5359
<i>Overall mortality</i>	10 (9)	428 (12)	18 (19)*	98 (16)*	68 (9)	86 (15)	4387/5359
<i>CDI mortality</i>	2 (2)	104 (3)	4 (4)	29 (5)	16 (2)	27 (5)	4387/5359
<b>Community onset</b>	75 (54)	1545 (36)*	31 (27)*	272 (37)*	356 (38)*	155 (23)*	5283/5359
<b>CDI last 8 weeks</b>	22 (27)	684 (25)	12 (20)	133 (29)	161 (27)	115 (25)	3312/5359

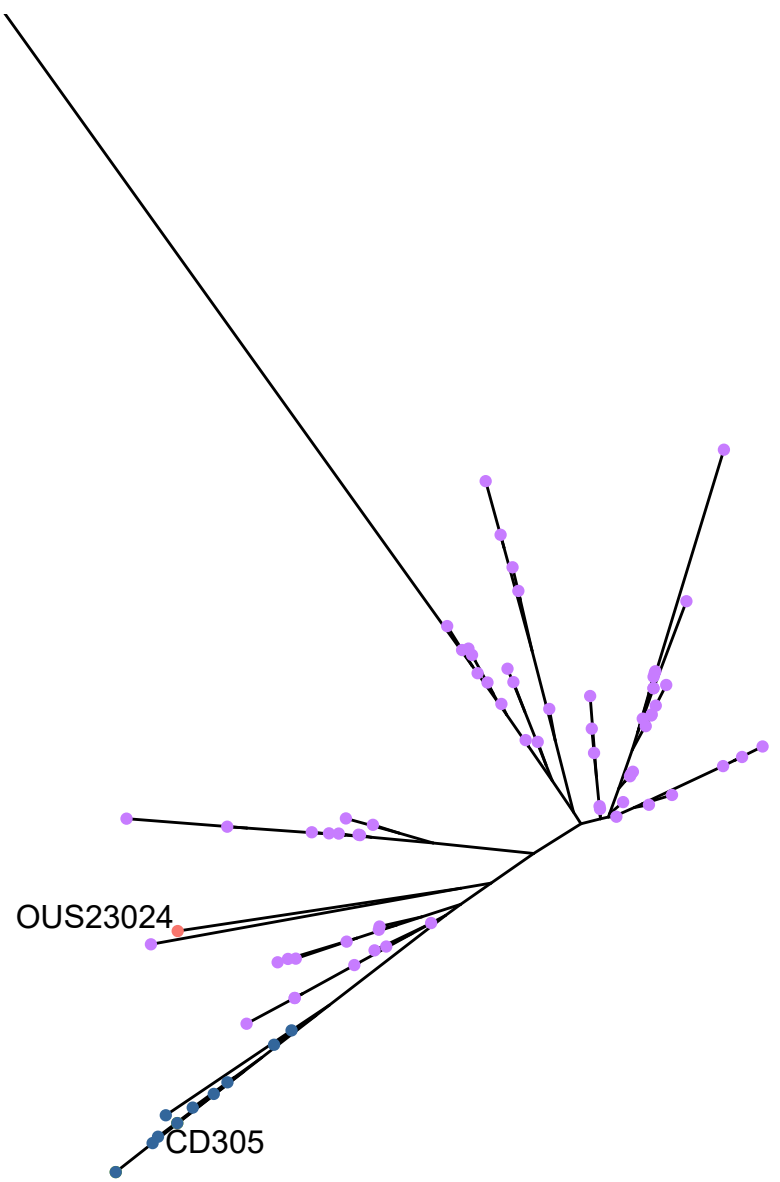
408 Data are presented as number of cases (percentage). Age is presented as median [first quartile, third quartile], because of the skewed distribution. Categorical variables were  
409 compared by a Pearson's Chi square test and numerical variables were compared by a Wilcoxon rank-sum test. An asterisk (\*) represents a  $p$ -value<0.05, when comparing  
410 with RT023. Abbreviations: LTCF: longtermcare facility, HCF: healthcare facility, RT: ribotype, CDI: Clostridium difficile infection

Journal Pre-proof

A

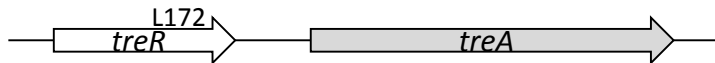


B

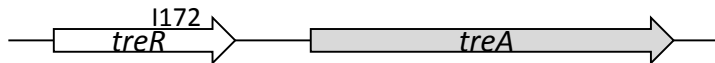


**a**

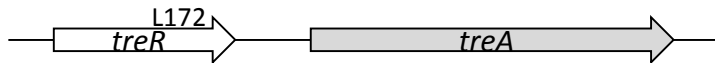
RT012 630



RT027 R20291

**b**

RT078 M120

**c**

RT023 CD305

