

Manuscript version: Author's Accepted Manuscript

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/154304>

How to cite:

The repository item page linked to above, will contain details on accessing citation guidance from the publisher.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work of researchers of the University of Warwick available open access under the following conditions.

This article is made available under the Creative Commons Attribution 4.0 International license (CC BY 4.0) and may be reused according to the conditions of the license. For more details see: <http://creativecommons.org/licenses/by/4.0/>.



Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk

Major genetic discontinuity and novel toxigenic species in *Clostridioides difficile* taxonomy

Daniel R. Knight^{1,2}, Korakrit Imwattana^{2,3}, Brian Kullin⁴, Enzo Guerrero-Araya^{5,6}, Daniel Paredes-Sabja^{5,6,7}, Xavier Didelot⁸, Kate E. Dingle⁹, David W. Eyre¹⁰, César Rodríguez¹¹, and Thomas V. Riley^{1,2,12,13*}

¹ Medical, Molecular and Forensic Sciences, Murdoch University, Murdoch, Western Australia, Australia. ² School of Biomedical Sciences, the University of Western Australia, Nedlands, Western Australia, Australia. ³ Department of Microbiology, Faculty of Medicine Siriraj Hospital, Mahidol University, Thailand. ⁴ Department of Pathology, University of Cape Town, Cape Town, South Africa. ⁵ Microbiota-Host Interactions and Clostridia Research Group, Facultad de Ciencias de la Vida, Universidad Andrés Bello, Santiago, Chile. ⁶ Millenium Nucleus in the Biology of Intestinal Microbiota, Santiago, Chile. ⁷ Department of Biology, Texas A&M University, College Station, TX, 77843, USA. ⁸ School of Life Sciences and Department of Statistics, University of Warwick, Coventry, UK. ⁹ Nuffield Department of Clinical Medicine, University of Oxford, Oxford, UK; National Institute for Health Research (NIHR) Oxford Biomedical Research Centre, John Radcliffe Hospital, Oxford, UK. ¹⁰ Big Data Institute, Nuffield Department of Population Health, University of Oxford, Oxford, UK; National Institute for Health Research (NIHR) Oxford Biomedical Research Centre, John Radcliffe Hospital, Oxford, UK. ¹¹ Facultad de Microbiología & Centro de Investigación en Enfermedades Tropicales (CIET), Universidad de Costa Rica, San José, Costa Rica. ¹² School of Medical and Health Sciences, Edith Cowan University, Joondalup, Western Australia, Australia. ¹³ Department of Microbiology, PathWest Laboratory Medicine, Queen Elizabeth II Medical Centre, Nedlands, Western Australia, Australia.

*Address correspondence to Professor Thomas V. Riley (thomas.riley@uwa.edu.au), School of Biomedical Sciences, The University of Western Australia, Nedlands, Western Australia, Australia.

Word count (main text): 4776 (methods = 965)
Abstract word count: 149

Abstract

Clostridioides difficile infection (CDI) remains an urgent global One Health threat. The genetic heterogeneity seen across *C. difficile* underscores its wide ecological versatility and has driven the significant changes in CDI epidemiology seen in the last 20 years. We analysed an international collection of over 12,000 *C. difficile* genomes spanning the eight currently defined phylogenetic clades. Through whole-genome average nucleotide identity, and pangenomic and Bayesian analyses, we identified major taxonomic incoherence with clear species boundaries for each of the recently described cryptic clades CI-III. The emergence of these three novel genomospecies predates clades C1-5 by millions of years, rewriting the global population structure of *C. difficile* specifically and taxonomy of the *Peptostreptococcaceae* in general. These genomospecies all show unique and highly divergent toxin gene architecture, advancing our understanding of the evolution of *C. difficile* and close relatives. Beyond the taxonomic ramifications, this work may impact the diagnosis of CDI.

Introduction

The bacterial species concept remains controversial, yet it serves as a critical framework for all aspects of modern microbiology¹. The prevailing species definition describes a genomically coherent group of strains sharing high similarity in many independent phenotypic and ecological properties². The era of whole-genome sequencing (WGS) has seen average nucleotide identity (ANI) replace DNA-DNA hybridization as the next-generation standard for microbial taxonomy^{3,4}. Endorsed by the National Center for Biotechnology Information (NCBI)⁴, ANI provides a precise,

objective and scalable method for delineation of species, defined as monophyletic groups of strains with genomes that exhibit at least 96% ANI^{5,6}.

Clostridioides (Clostridium) difficile is an important gastrointestinal pathogen that places a significant growing burden on health care systems in many regions of the world⁷. In both its 2013⁸ and 2019⁹ reports on antimicrobial resistance (AMR), the US Centers for Disease Control and Prevention rated *C. difficile* infection (CDI) as an urgent health threat, the highest level. Community-associated CDI has become more frequent⁷ and is linked to sources of *C. difficile* in animals and the environment¹⁰. Thus, over the last two decades, CDI has emerged as an important One Health issue¹⁰.

Based on multi-locus sequence type (MLST), there are eight recognised monophyletic groups or ‘clades’ of *C. difficile*¹¹. Strains within these clades show many unique clinical, microbiological, and ecological features¹¹. Critical to the pathogenesis of CDI is the expression of the large clostridial toxins, TcdA and TcdB and, in some strains, binary toxin (CDT), encoded by two separate chromosomal loci, the PaLoc and CdtLoc, respectively¹². Clade 1 (C1) contains over 200 toxigenic and non-toxigenic sequence types (STs) including many of the most prevalent strains causing CDI worldwide e.g., ST2, ST8, and ST17¹¹. Several highly virulent CDT-producing strains, including ST1 (PCR ribotype (RT) 027), a lineage associated with major hospital outbreaks in North America, Europe, and Latin America¹³, are found in clade 2 (C2). Comparatively little is known about clade 3 (C3) although it contains ST5 (RT 023), a toxigenic CDT-producing strain with characteristics that may make laboratory detection difficult¹⁴. *C. difficile* ST37 (RT 017) is found in clade 4 (C4) and, despite the absence of a toxin A gene, is responsible for much of the endemic CDI burden in Asia¹⁵. Clade 5 (C5) contains several CDT-producing strains including ST11 (RTs 078, 126 and others), which are highly prevalent in production animals worldwide¹⁶. The remaining so-called ‘cryptic’ clades (C-I, C-II and C-III), first described in 2012^{17, 18}, contain over 50 STs from clinical and environmental sources^{17, 18, 19, 20, 21}. The evolution of the cryptic clades is poorly understood. Clade C-I strains can cause CDI, however, due to atypical toxin gene architecture, they may not be detected, thus their prevalence may have been underestimated²¹.

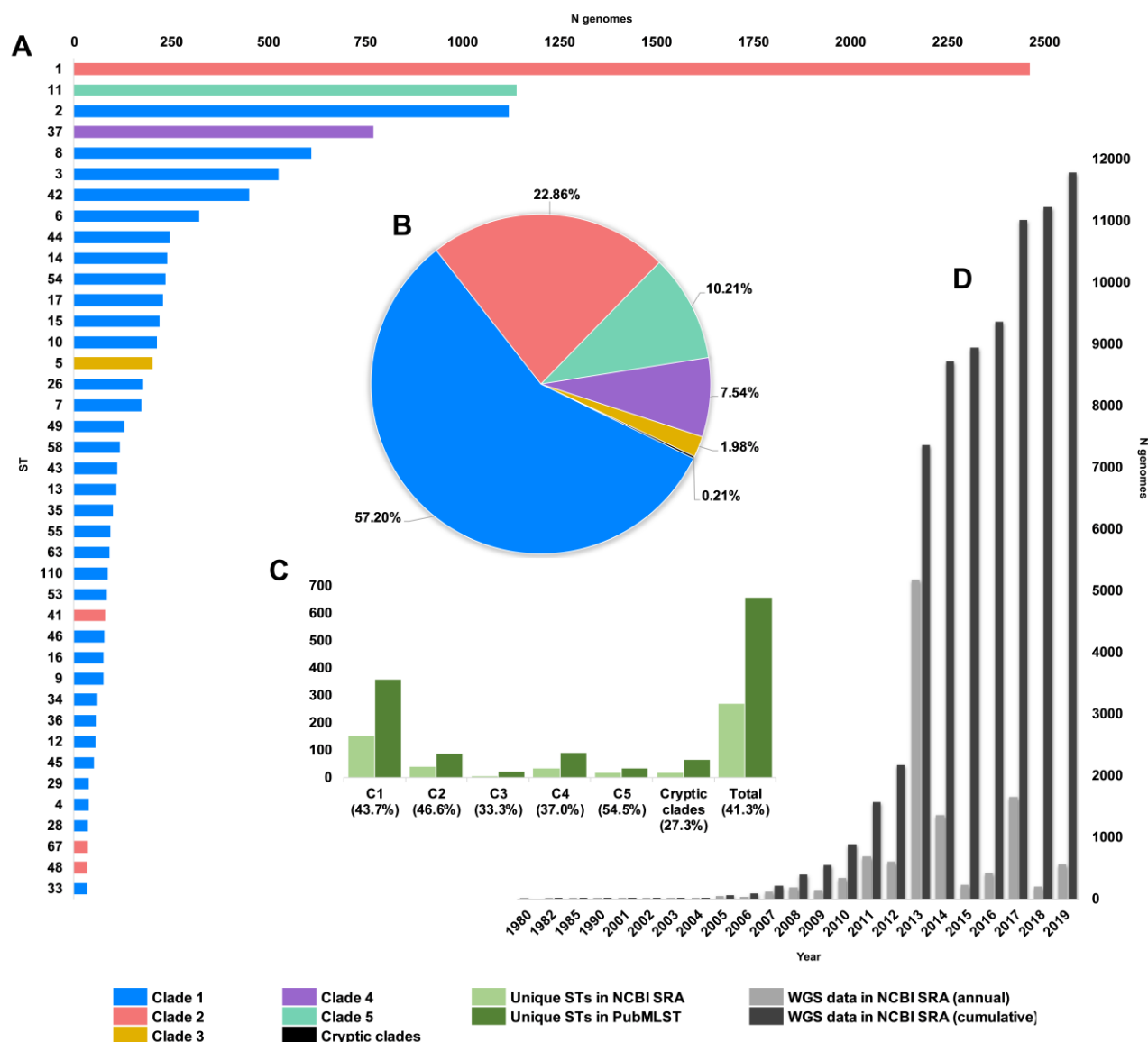
There are over 600 STs currently described and some STs may have access to a gene pool of more than 10,000 genes^{11, 16, 22}. Considering such enormous diversity, and recent contentious taxonomic revisions^{23, 24}, we hypothesise that *C. difficile* comprises a complex of distinct species divided along the major evolutionary clades. In this study, whole-genome ANI, and pangenomic and Bayesian analyses are used to explore an international collection of over 12,000 *C. difficile* genomes, to provide new insights into ancestry, genetic diversity, and evolution of pathogenicity in this enigmatic pathogen.

Results

An updated global population structure based on sequence typing of 12,000 genomes. We obtained and determined the ST and clade for a collection of 12,621 *C. difficile* genomes (taxid ID 1496, Illumina data) existing in the NCBI Sequence Read Archive (SRA) as of 1st January 2020. A total of 272 STs were identified spanning the eight currently described clades, indicating that the SRA contains genomes for almost 40% of known *C. difficile* STs worldwide (n=659, PubMLST, January 2020). C1 STs dominated the database in both prevalence and diversity (**Fig. 1**) with 149 C1 STs comprising 57.2% of genomes, followed by C2 (35 STs, 22.9%), C5 (18 STs, 10.2%), C4 (34 STs, 7.5%), C3 (7 STs, 2.0%) and the cryptic clades C-I, C-II and C-III (collectively 17 STs, 0.2%). The five most prevalent STs represented were ST1 (20.9% of genomes), ST11 (9.8%), ST2 (9.5%), ST37 (6.5%) and ST8 (5.2%), all prominent lineages associated with CDI worldwide¹¹.

Fig. 2 shows an updated global *C. difficile* population structure based on the 659 STs; 27 novel STs were found (an increase of 4%) and some corrections to assignments within C1 and C2 were made, including assigning ST122²⁵ to C1. Based on PubMLST data and bootstraps values of 1.0 in all monophyletic nodes of the cryptic clades (**Fig. 2**), we could confidently assign 25, 9 and 10 STs to cryptic clades I, II and III, respectively. There remained 26 STs spread across the phylogeny that did not fit within a specific clade (defined as outliers). The tree file for **Fig. 2** and

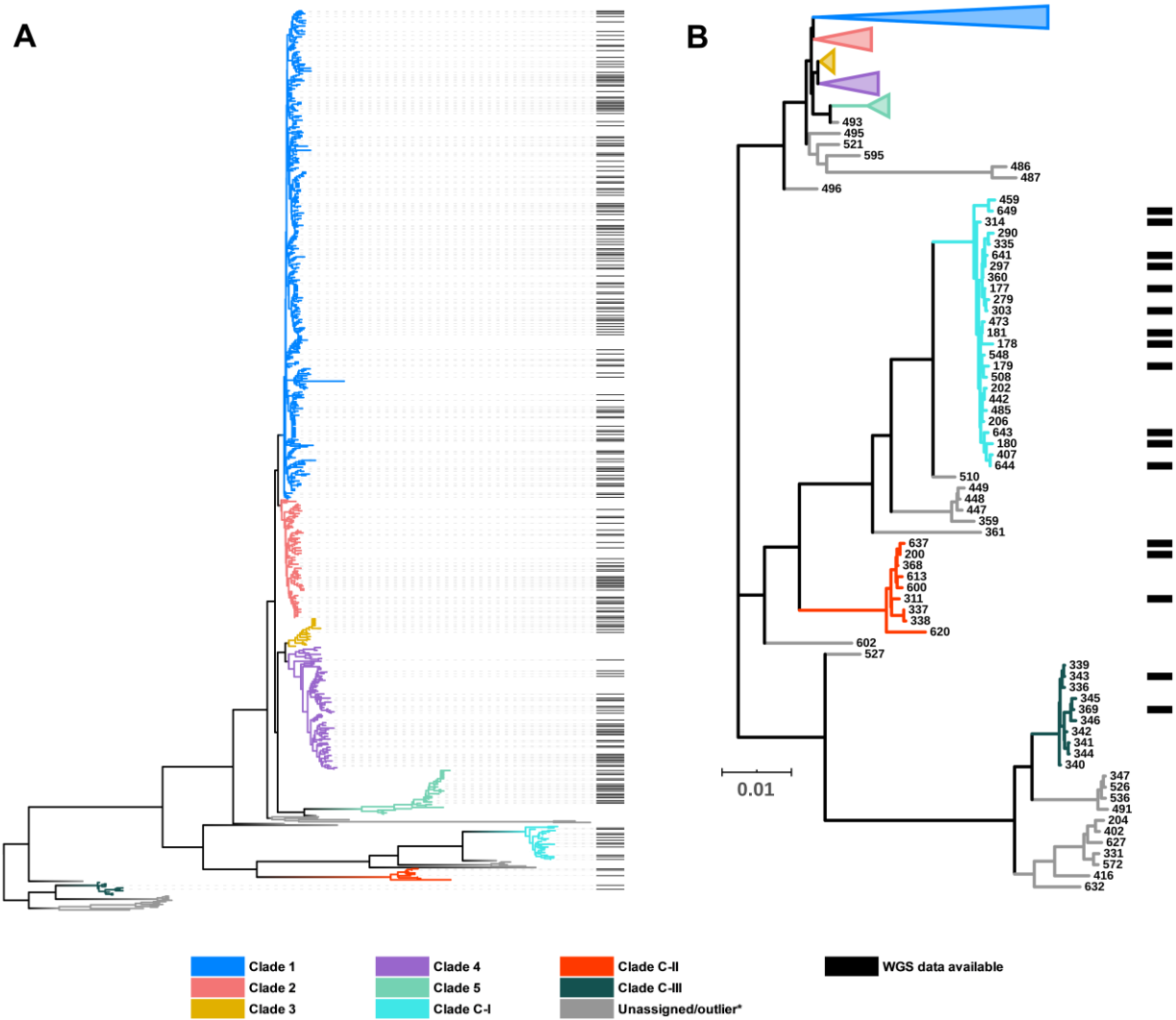
96 full MLST data is available as **Supplementary Files 1a-1d and 2** at
 97 <http://doi.org/10.6084/m9.figshare.12471461>. Representative genomes of each ST present in the
 98 ENA were chosen based on metadata, read depth and assembly quality. This resulted in a final
 99 dataset of 260 STs (C1, n=149; C2, n=35; C3, n=7; C4, n=34; C5, n=18; C-I, n=12; C-II, n=3, C-
 100 III, n=2) used for all subsequent bioinformatics analyses. The list of representative genomes is
 101 available in **Supplementary File 1b**.



102 **Figure 1. Composition of *C. difficile* genomes in the NCBI SRA.** Snapshot obtained 1st January 2020;
 103 12,304 strains, [taxid ID 1496]. (A) Top 40 most prevalent STs in the NCBI SRA coloured by clade. (B) The
 104 proportion of genomes in ENA by clade. (C) Number/ proportion of STs per clade found in the SRA/present
 105 in the PubMLST database. (D) Annual and cumulative deposition of *C. difficile* genome data in ENA.

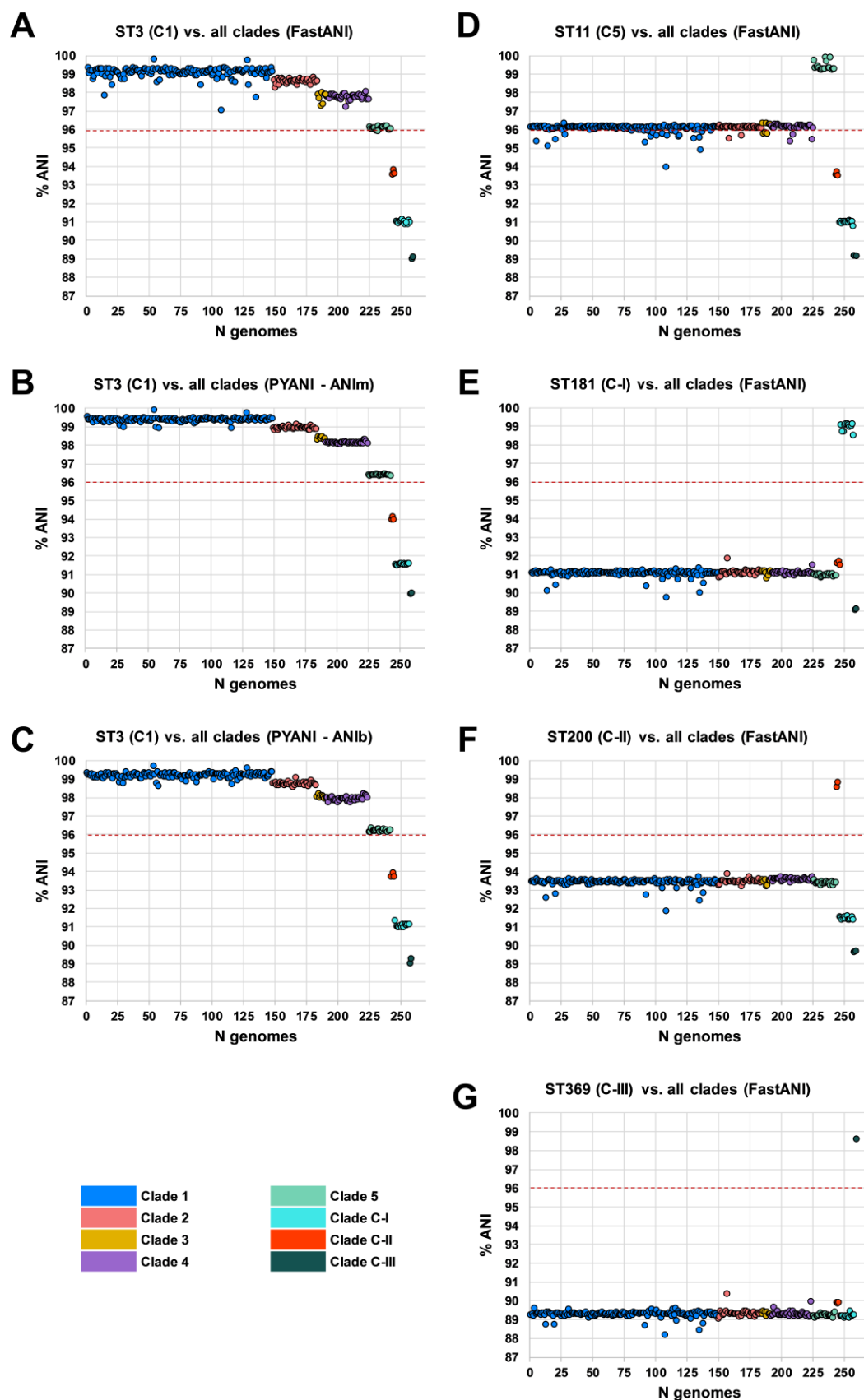
106 **Whole-genome ANI analysis reveals clear species boundaries.** Whole-genome ANI analyses
 107 were used to investigate genetic discontinuity across the *C. difficile* species (Fig. 3 and
 108 **Supplementary File 1f**). Whole-genome ANI values were determined for the final set of 260 STs
 109 using three independent ANI algorithms (FastANI, ANIm and ANIb, see *Methods*). All 225

110 genomes belonging to clades C1-4 clustered within an ANI range of 97.1-99.8% (median FastANI
 111 values of 99.2, 98.7, 97.9 and 97.8%, respectively, **Fig. 3A-C**).



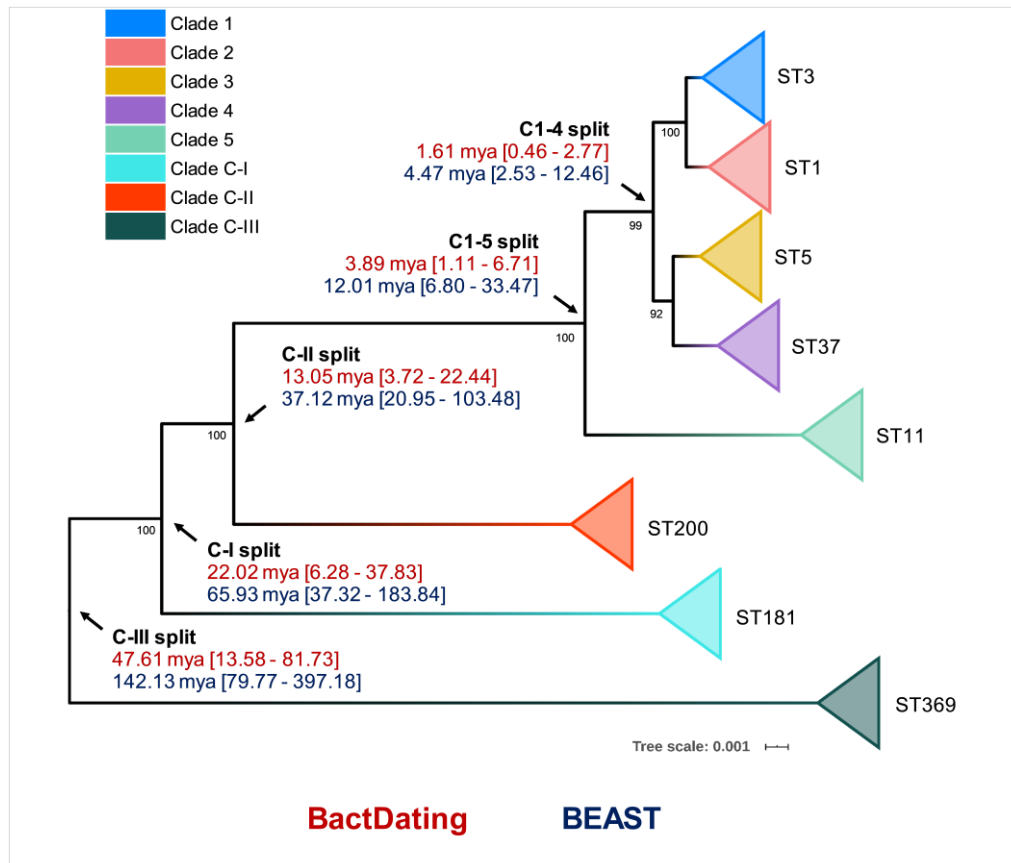
112 **Figure 2. *C. difficile* population structure.** (A) NJ phylogeny of 659 aligned, concatenated, multilocus
 113 sequence type allele combinations coloured by current PubMLST clade assignment. Black bars indicate
 114 WGS available for ANI analysis (n=260). (B) A subset of the NJ tree showing cryptic clades C-I, C-II and C-
 115 III. Again, black bars indicate WGS available for ANI analysis (n=17).

116 These ANI values are above the 96% species demarcation threshold used by the NCBI⁴ and indicate
 117 that strains from these clades belong to the same species. ANI values for all 18 genomes belonging
 118 to C5 clustered on the borderline of the species demarcation threshold (FastANI range 95.9-96.2%,
 119 median 96.1%). ANI values for all three cryptic clades fell well below the species threshold; C-I
 120 (FastANI range 90.9-91.1%, median 91.0%), C-II (FastANI range 93.6-93.9%, median 93.7%) and
 121 C-III (FastANI range 89.1-89.1%, median 89.1%). All results were corroborated across the three
 122 independent ANI algorithms (**Fig. 3A-C**). *C. difficile* strain ATCC 9689 (ST3, C1) was defined by
 123 Lawson *et al.* as the type strain for the species²³, and used as a reference in all the above analyses.
 124 To better understand the diversity among the divergent clades themselves, FastANI analyses were
 125 repeated using STs 11, 181, 200 and 369 as reference archetypes of clades C5, C-I, C-II and C-III,
 126 respectively. This approach confirmed that C5 and the three cryptic clades were as distinct from
 127 each other as they were collectively from C1-4 (**Fig. 3D-G**).



128 **Figure 3. Species-wide ANI analysis.** Panels A-C show ANI plots for ST3 (C1) vs. all clades (260 STs)
 129 using FastANI, ANIm and ANIb algorithms, respectively. Panels D-G show ANI plots for ST11 (C5),
 130 ST181 (C-I), ST200 (C-II) and ST369 (C-III) vs all clades (260 STs), respectively. NCBI species
 131 demarcation of 96% indicated by red dashed line⁴.

132 **Taxonomic placement of cryptic clades predates *C. difficile* emergence by millions of years.**
 133 Previous studies using BEAST have estimated the common ancestor of C1-5 existed between 1 to
 134 85 or 12 to 14 million years ago (mya)^{26, 27}. Here, we used an alternative Bayesian approach,
 135 BactDating, to estimate the age of all eight *C. difficile* clades currently described. The last common
 136 ancestor for *C. difficile* clades C1-5 was estimated to have existed between 1.11 to 6.71 mya. In
 137 contrast, all three cryptic clades were estimated to have emerged millions of years prior to the
 138 common ancestor of C1-5 (Fig. 4). Independent analysis with BEAST, using a smaller core gene
 139 dataset (see *Methods*), provided temporal estimates of clade emergence that were of the same order
 140 of magnitude and, importantly, supported the same branching order for all clades (Fig. 4).



141 **Figure 4. Bayesian analysis of species and clade divergence.** BactDating and BEAST estimates of the age
 142 of major *C. difficile* clades. Node dating ranges for both Bayesian approaches are transposed onto an ML
 143 phylogeny built from concatenated MLST alleles of a dozen STs from each clade. Archetypal STs in each
 144 evolutionary clade are indicated. The tree is midpoint rooted and bootstrap values are shown (all
 145 bootstrapping values of the cryptic clade branches are 100%). Scale bar indicates the number of substitutions
 146 per site. BactDating estimates the median time of the most recent common ancestor of C1-5 at 3.89 million
 147 years ago (mya) [95% credible interval (CI), 1.11-6.71 mya]. Of the cryptic clades, C-II shared the most
 148 recent common ancestor with C1-5 (13.05 mya, 95% CI 3.72-22.44 mya), followed by C-I (22.02 mya, 95%
 149 CI 6.28-37.83 mya), and C-III (47.61 mya, 95% CI 13.58-81.73 mya). Comparative temporal estimates from
 150 BEAST show the same order of magnitude and support the same branching order [clades C1-5 (12.01 mya,
 151 95% CI 6.80-33.47 mya); C-II (37.12 mya, 95% CI 20.95-103.48 mya); C-I (65.93 mya, 95% CI 37.32-
 152 183.84 mya); C-III (142.13 mya, 95% CI 79.77-397.18 mya)].

153 Next, to identify their true taxonomic placement, ANI was determined for ST181 (C-I), ST200 (C-
 154 II) and ST369 (C-III) against two reference datasets. The first dataset comprised 25 species
 155 belonging to the *Peptostreptococcaceae* as defined by Lawson *et al.*²³ in their 2016 reclassification
 156 of *Clostridium difficile* to *Clostridioides difficile*. The second dataset comprised 5,895 complete
 157 genomes across 21 phyla from the NCBI RefSeq database (accessed 14th January 2020), including
 158 1,366 genomes belonging to *Firmicutes*, 92 genomes belonging to 15 genera within the

159 *Clostridiales* and 20 *Clostridium* and *Clostridioides* species. The nearest ANI matches to species
 160 within the *Peptostreptococcaceae* dataset were *C. difficile* (range 89.3-93.5% ANI),
 161 *Asaccharospora irregularis* (78.9-79.0% ANI) and *Romboutsia lituseburensis* (78.4-78.7% ANI).
 162 Notably, *Clostridioides manganotii*, the only other known member of *Clostridioides*, shared only
 163 77.2-77.8% ANI with the cryptic clade genomes (Table 1).

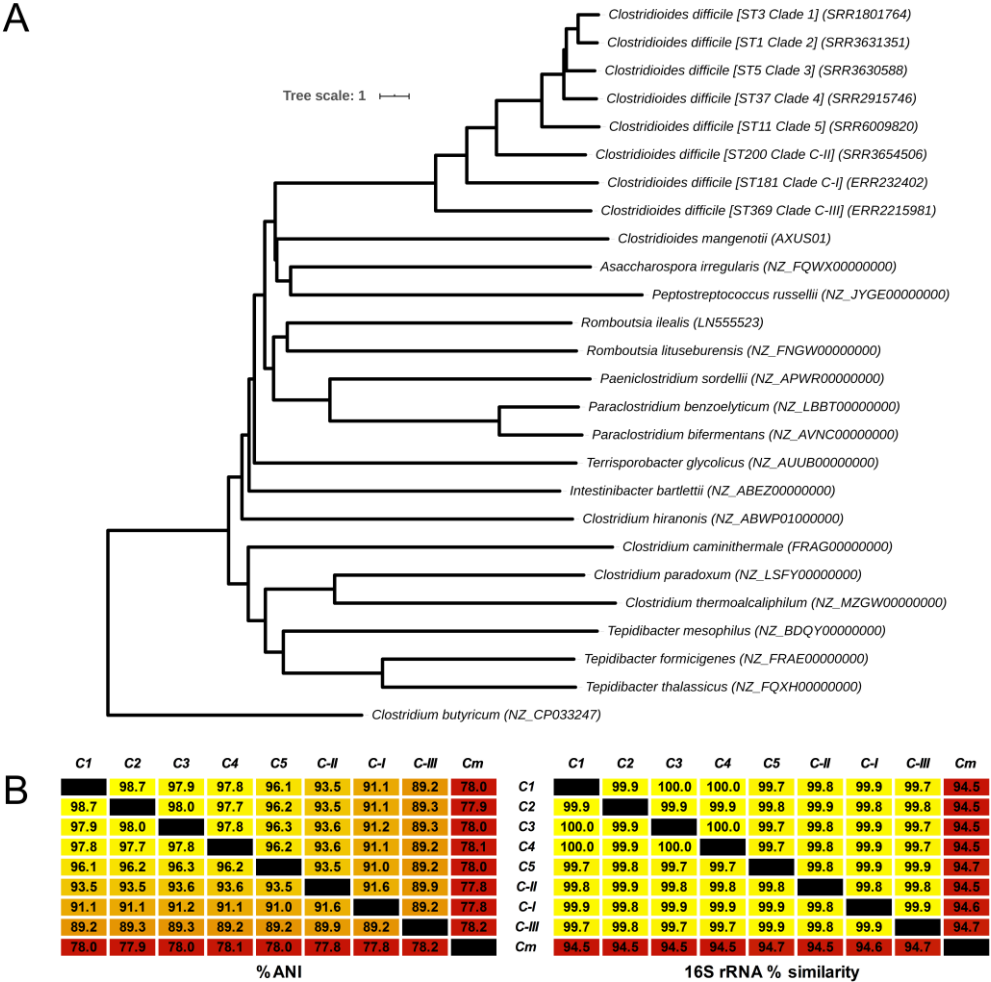
164 **Table 1 Whole-genome ANI analysis of cryptic clades vs. 25 *Peptostreptococcaceae* species**
 165 **from Lawson *et al*²³.**

Species	NCBI accession	ANI %		
		ST181 (C-I)	ST200 (C-II)	ST369 (C-III)
<i>Clostridioides difficile</i> (ST3)	AQWV000000000.1	91.11	93.54	89.30
<i>Asaccharospora irregularis</i>	NZ_FQWX000000000	78.94	78.87	78.91
<i>Romboutsia lituseburensis</i>	NZ_FNGW000000000.1	78.51	78.36	78.66
<i>Romboutsia ilealis</i>	LN555523.1	78.45	78.54	78.44
<i>Paraclostridium benzoelyticum</i>	NZ_LBBT000000000.1	77.92	77.71	78.14
<i>Paraclostridium bifermentans</i>	NZ_AVNC000000000.1	77.89	77.89	78.06
<i>Clostridioides manganotii</i>	GCA_000687955.1	77.82	77.84	78.15
<i>Paeniclostridium sordellii</i>	NZ_APWR000000000.1	77.73	77.59	77.86
<i>Clostridium hiranonis</i>	NZ_ABWP01000000	77.52	77.42	77.59
<i>Terrisporobacter glycolicus</i>	NZ_AUUB000000000.1	77.47	77.53	77.53
<i>Intestinibacter bartlettii</i>	NZ_ABEZ000000000.2	77.29	77.52	77.48
<i>Clostridium paradoxum</i>	NZ_LSFY000000000.1	76.60	76.65	76.93
<i>Clostridium thermoalcaliphilum</i>	NZ_MZGW000000000.1	76.49	76.61	76.85
<i>Tepidibacter formicigenes</i>	NZ_FRAE000000000.1	76.41	76.47	76.38
<i>Tepidibacter mesophilus</i>	NZ_BDQY000000000.1	76.38	76.44	76.22
<i>Tepidibacter thalassicus</i>	NZ_FQXH000000000.1	76.34	76.31	76.46
<i>Peptostreptococcus russellii</i>	NZ_JYGE000000000.1	76.30	76.08	76.38
<i>Clostridium formicaceticum</i>	NZ_CP020559.1	75.18	75.26	75.62
<i>Clostridium caminithermale</i>	FRAG000000000	74.97	75.07	75.03
<i>Clostridium acetivum</i>	NZ_JYHU000000000.1	≤70.00	≤70.00	≤70.00
<i>Clostridium litorale</i>	FSRH01000000	≤70.00	≤70.00	≤70.00
<i>Eubacterium acidaminophilum</i>	NZ_CP007452.1	≤70.00	≤70.00	≤70.00
<i>Filifactor alocis</i>	NC_016630.1	≤70.00	≤70.00	≤70.00
<i>Peptostreptococcus anaerobius</i>	ARMA01000000	≤70.00	≤70.00	≤70.00
<i>Peptostreptococcus stomatis</i>	NZ_ADGQ000000000.1	≤70.00	≤70.00	≤70.00

166 Similarly, the nearest ANI matches to species within the RefSeq dataset were several *C. difficile*
 167 strains (range C-I: 90.9-91.1%; C-II: 93.4-93.6%; and C-III: 89.2-89.4%) and *Paeniclostridium*
 168 *sordellii* (77.7-77.9%). A low ANI (range ≤70-75%) was observed between the cryptic clade
 169 genomes and 20 members of the *Clostridium* including *C. tetani*, *C. botulinum*, *C. perfringens* and
 170 *C. butyricum*, the type strain of the *Clostridium* genus *sensu stricto*. An updated ANI-based
 171 taxonomy for the *Peptostreptococcaceae* is shown in Fig. 5A. The phylogeny places C-I, C-II and
 172 C-III between *C. manganotii* and *C. difficile* C1-5, suggesting that they should be assigned to the
 173 *Clostridioides* genus, distinct from both *C. manganotii* and *C. difficile*. Comparative analysis of
 174 ANI and 16S rRNA values for the eight *C. difficile* clades and *C. manganotii* shows significant
 175 incongruence between the data generated by the two approaches (Fig. 5B). The range of 16S rRNA
 176 % similarity between *C. difficile* C1-4, cryptic clades I-III and *C. manganotii* was narrower (range
 177 94.5-100) compared to the range of ANI values (range 77.8-98.7). Curiously, *C. manganotii* and *C.*
 178 *difficile* shared 94.5-94.7% similarity in 16s rRNA sequence identity, yet only 77.8-78.2% ANI,
 179 indicating they should not even be considered within the same genus, as proposed by Lawson *et*
 180 *al.*²³

181 We also extended our approach to five other medically important clostridia available on the
 182 NCBI database; *Clostridium botulinum* (n=783), *Clostridium perfringens* (n=358), *Clostridium*
 183 *sporogenes* (n=100), *Clostridium tetani* (n=32), and *Paeniclostridium sordellii* (formerly
 184 *Clostridium sordellii*, n=46). We found that three out of the five species (*C. perfringens*, *C.*
 185 *sporogenes*, and *C. botulinum*) showed evidence of taxonomic discontinuity similar to that observed

186 for *C. difficile* (e.g., a proportion of strains with pairwise ANI below the 96% demarcation
 187 threshold). This was most notable for *C. sporogenes* and *C. botulinum*, where there were many
 188 sequenced strains with a pairwise ANI below 90% (8% and 31% of genomes, respectively,
 189 **Supplementary File 1i**).



190 **Figure 5. Revised taxonomy for the *Peptostreptococcaceae*.** (A) ANI-based minimum evolution tree
 191 showing evolutionary relationship between eight *C. difficile* ‘clades’ along with 17 members of the
 192 *Peptostreptococcaceae* (from Lawson *et al*²³) as well as *Clostridium butyricum* as the outgroup and type
 193 strain of the *Clostridium* genus *sensu stricto*. To convert the ANI into a distance, its complement to 1 was
 194 taken. (B) Matrices showing pairwise ANI and 16S rRNA values for the eight *C. difficile* clades and *C.*
 195 *manganotii*, the only other known member of *Clostridioides*.

196 **Evolutionary and ecological insights from the *C. difficile* species pangenome.** Next, we sought
 197 to quantify the *C. difficile* species pangenome and identify genetic loci that are significantly
 198 associated with the taxonomically divergent clades. With Panaroo, the *C. difficile* species
 199 pangenome comprised 17,470 genes, encompassing an accessory genome of 15,238 genes and a
 200 core genome of 2,232 genes, just 12.8% of the total gene repertoire (**Fig 6**). The size of the
 201 pangenome reduced by 2,082 genes with the exclusion of clades CI-III, and a further 519 genes
 202 with the exclusion of C5. Compared to Panaroo, Roary overestimated the size of the pangenome
 203 (32,802 genes, 87.7% overestimation), resulting in markedly different estimates of the percentage
 204 core genome, 3.9 and 12.8%, respectively ($\chi^2=1,395.3$, $df=1$, $p<0.00001$). The overestimation of
 205 pangenome was less pronounced when the identity threshold was decreased to 90% (42.0%

overestimation) and when the paralogs were merged (28.7% overestimation). Panaroo can account for errors introduced during assembly and annotation, thus polishing the 260 Prokka-annotated genomes with Panaroo resulted in a significant reduction in gene content per genome (median 2.48%; 92 genes, range 1.24-12.40%; 82-107 genes, $p < 0.00001$). The *C. difficile* species pangenome was determined to be open²⁸ (Fig 6).

Pan-GWAS analysis with Scoary revealed 142 genes with significant clade specificity. Based on KEGG orthology, these genes were classified into four functional categories: environmental information processing (7), genetic information processing (39), metabolism (43), and signaling and cellular processes (53). We identified several uniquely present, absent, or organised gene clusters associated with ethanolamine catabolism (C-III), heavy metal uptake (C-III), polyamine biosynthesis (C-III), fructosamine utilization (C-I, C-III), zinc transport (C-II, C5) and folate metabolism (C-I, C5). A summary of the composition and function of these major lineage-specific gene clusters is given in Table 2, and a comparative analysis of their respective genetic architecture can be found in Supplementary File 1h.

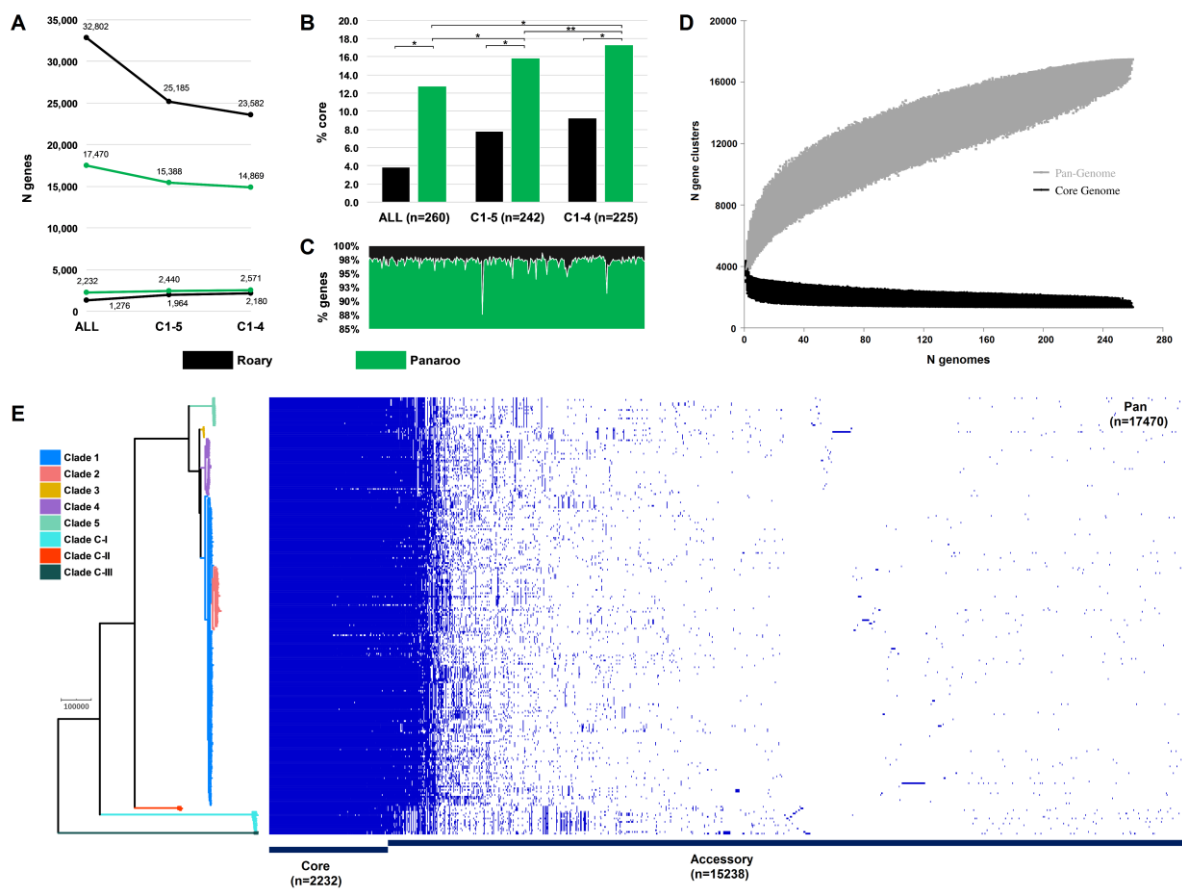


Figure 6. *Clostridioides difficile* species pangenome. (A) Pan and core genome estimates for all 260 STs, clades C1-4 (n=242 STs) and clades C1-5 (n=225 STs). (B) The difference in % core genome and pangenome sizes with Panaroo and Roary algorithms. (*) indicates $\chi^2 p < 0.00001$ and (**) indicates $\chi^2 p = 0.0008$. (C) The proportion of retained genes per genome after polishing Prokka-annotated genomes with Panaroo. (D) The total number of genes in the pan (grey) and core (black) genomes are plotted as a function of the number of genomes sequentially added (n=260). Following the definition of Tettelin *et al.*²⁸, the *C. difficile* species pangenome showed characteristics of an “open” pangenome. First, the pangenome increased in size exponentially with sampling of new genomes. At n=260, the pangenome exceeded more than double the average number of genes found in a single *C. difficile* genome (~3,700) and the curve was yet to reach a plateau or exponentially decay, indicating more sequenced strains are needed to capture the complete species gene repertoire. Second, the number of new ‘strain-specific’ genes did not converge to zero upon sequencing of additional strains, at n=260, an average of 27 new genes were contributed to the gene pool. Finally, according to Heap’s Law, α values of ≤ 1 are representative of open pangenome. Rarefaction analysis of our

233 pangenome curve using a power-law regression model based on Heap's Law²⁸ showed the pangenome was
 234 predicted to be open ($B_{pan} (\approx \alpha^{28}) = 0.47$, curve fit, $r^2=0.999$). (E) Presence absence variation (PAV) matrix
 235 for 260 *C. difficile* genomes is shown alongside a maximum-likelihood phylogeny built from a
 236 recombination-adjusted alignment of core genes from Panaroo (2,232 genes, 2,606,142 sites).

237 **Table 2 Major clade-specific gene clusters identified by pan-GWAS**

Protein	Gene	Clade specificity	Functional insights
Ethanolamine kinase	<i>ETNK, EKI</i>	Unique to C-III and is in addition to the highly conserved <i>eut</i> cluster found in all lineages. Has a unique composition and includes six additional genes that are not present in the traditional CD630 <i>eut</i> operon or any other non-C-III strains.	An alternative process for the breakdown of ethanolamine and its utilization as a source of reduced nitrogen and carbon.
Agmatinase	<i>speB</i>		
1-propanol dehydrogenase	<i>pduQ</i>		
Ethanolamine utilization protein EutS	<i>eutS</i>		
Ethanolamine utilization protein EutP	<i>eutP</i>		
Ethanolamine ammonia-lyase large subunit	<i>eutB</i>		
Ethanolamine ammonia-lyase small subunit	<i>eutC</i>		
Ethanolamine utilization protein EutL	<i>eutL</i>		
Ethanolamine utilization protein EutM	<i>eutM</i>		
Acetaldehyde dehydrogenase	<i>E1.2.1.10</i>		
Putative phosphotransacetylase	<i>K15024</i>		
Ethanolamine utilization protein EutN	<i>eutN</i>		
Ethanolamine utilization protein EutQ	<i>eutQ</i>		
TfoX/Sxy family protein	-	Unique to C-III	Multicomponent transport system with specificity for chelating heavy metal ions.
Iron complex transport system permease protein	<i>ABC.FEV.P</i>		
Iron complex transport system ATP-binding protein	<i>ABC.FEV.A</i>		
Iron complex transport system substrate-binding protein	<i>ABC.FEV.S</i>		
Hydrogenase nickel incorporation protein HypB	<i>hypB</i>		
Putative ABC transport system ATP-binding protein	<i>yxdL</i>		
Class I SAM-dependent methyltransferase	-		
Peptide/nickel transport system substrate-binding protein	<i>ABC.PE.S</i>		
Peptide/nickel transport system permease protein	<i>ABC.PE.P</i>		
Peptide/nickel transport system permease protein	<i>ABC.PE.P1</i>		
Peptide/nickel transport system ATP-binding protein	<i>ddpD</i>		
Oligopeptide transport system ATP-binding protein	<i>oppF</i>		
Class I SAM-dependent methyltransferase	-		
Heterodisulfide reductase subunit D [EC:1.8.98.1]	<i>hdrD</i>	Unique to C-III and is in addition to the highly conserved spermidine uptake cluster found in all other lineages.	Alternative spermidine uptake processes which may play a role in stress response to nutrient limitation. The additional cluster has homologs in <i>Romboutsia</i> , <i>Paraclostridium</i> and <i>Paeniclostridium</i> spp.
CDP-L-myo-inositol myo-inositolphosphotransferase	<i>dipps</i>		
Spermidine/putrescine transport system substrate-binding protein	<i>ABC.SP.S</i>		
Spermidine/putrescine transport system permease protein	<i>ABC.SP.P1</i>		
Spermidine/putrescine transport system permease protein	<i>ABC.SP.P</i>		
Spermidine/putrescine transport system ATP-binding protein	<i>potA</i>		
Sigma -54 dependent transcriptional regulator	<i>gfrR</i>	Present in all lineages except C-I. Cluster found in a different genomic position in C-III.	Mannose-type PTS system essential for utilization of fructosamines such as fructoselysine and glucoselysine, abundant components of rotting fruit and vegetable matter.
Fructoselysine/glucoselysine PTS system EIIB component	<i>gfrB</i>		
Mannose PTS system EIIA component	<i>manXa</i>		
Fructoselysine/glucoselysine PTS system EIIC component	<i>gfrC</i>		
Fructoselysine/glucoselysine PTS system EIID component	<i>gfrD</i>		
SIS domain-containing protein	-		
Fur family transcriptional regulator, ferric uptake regulator	<i>furB</i>	Unique to C-II and C5	Associated with EDTA resistance in <i>E.coli</i> , helping the bacteria survive in Zn-depleted environment.
Zinc transport system substrate-binding protein	<i>znuA</i>		
Fe-S-binding protein	<i>yeiR</i>		
Rrf2 family transcriptional regulator	-		
Putative signalling protein	-	Unique to C-I and C5 STs 163, 280, and 386	In <i>E. coli</i> , AbgAB proteins enable uptake and cleavage of the folate catabolite <i>p</i> -aminobenzoyl-glutamate, allowing the bacterium to survive on exogenous sources of folic acid.
Aminobenzoyl-glutamate utilization protein B	<i>abgB</i>		
MarR family transcriptional regulator	-		

238 **Cryptic clades CI-III possessed highly divergent toxin gene architecture.** Overall, 68.8%
 239 (179/260) of STs harboured *tcdA* (toxin A) and/or *tcdB* (toxin B), the major virulence factors in
 240 *C. difficile*, while 67 STs (25.8%) harboured *cdtA/cdtB* (binary toxin). The most common genotype

was A⁺B⁺CDT⁻ (113/187; 60.4%), followed by A⁺B⁺CDT⁺ (49/187; 26.2%), A⁻B⁺CDT⁺ (10/187; 5.3%), A⁻B⁻CDT⁺ (8/187; 4.3%) and A⁻B⁺CDT⁻ (7/187; 3.7%). Toxin gene content varied across clades (C1, 116/149, 77.9%; C2, 35/35, 100.0%; C3, 7/7, 100.0%; C4, 6/34, 17.6%; C5, 18/18, 100.0%; C-I, 2/12, 16.7%; C-II, 1/3, 33.3%; C-III, 2/2, 100.0%) (**Fig. 7**).

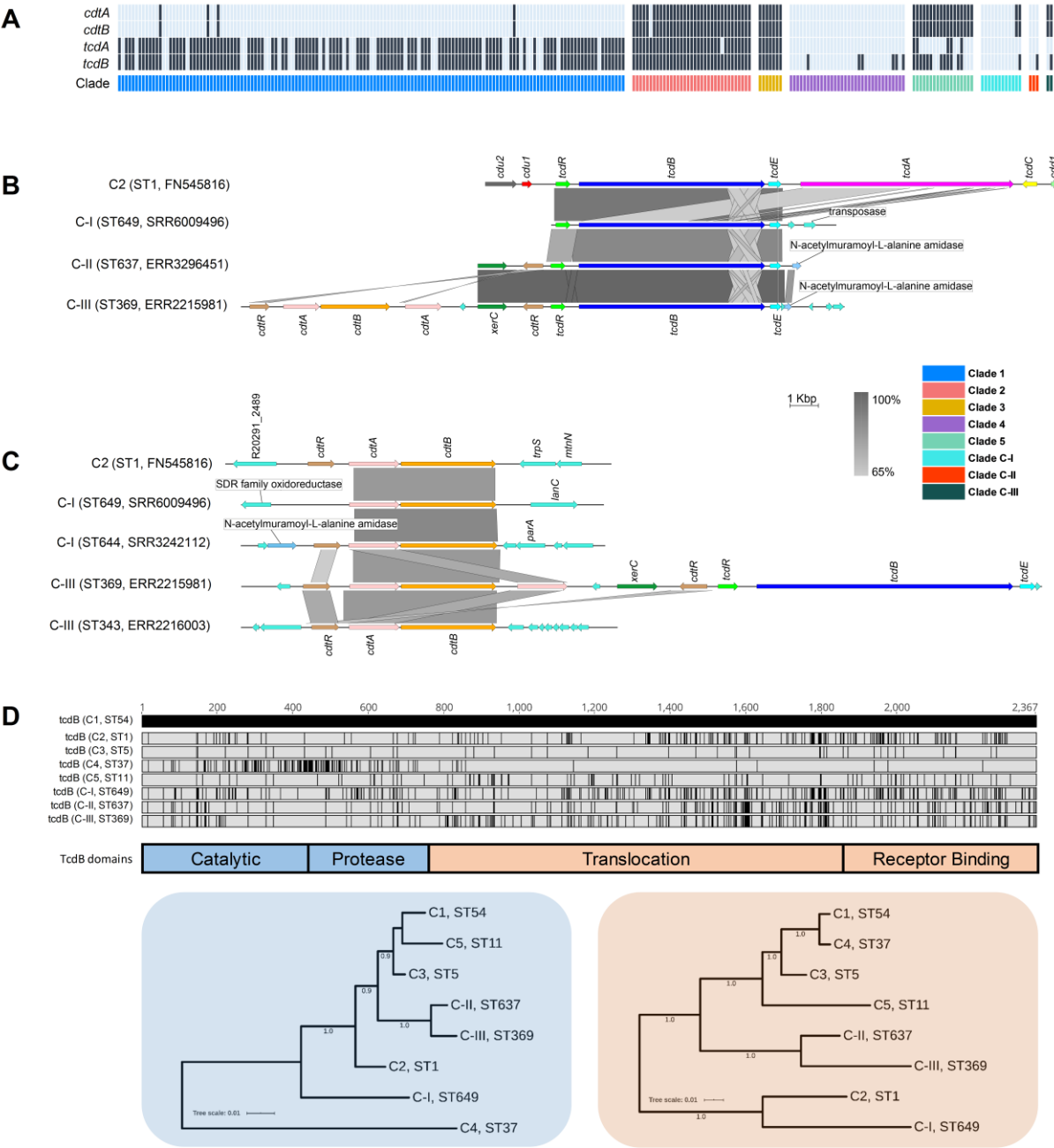


Figure 7. Toxin gene analysis. (A) Distribution of toxin genes across *C. difficile* clades (n=260 STs). Presence is indicated by black bars and absence by light blue bars. (B) Comparison of PaLoc architecture in the chromosome of strain R20291 (C2, ST1) and cognate chromosomal regions in genomes of cryptic STs 649 (C-I), 637 (C-II), and 369 (C-III). All three cryptic STs show atypical ‘monotoxin’ PaLoc structures, with the presence of syntenic *tcdR*, *tcdB*, and *tcdE*, and the absence of *tcdA*, *tcdC*, *cdd1* and *cdd2*. ST369 genome ERR2215981 shows colocalization of the PaLoc and CdtLoc, see below. (C) Comparison of CdtLoc architecture in the chromosome of strain R20291 (C2, ST1) and cognate chromosomal regions in genomes of cryptic STs 649/644 (C-I) and 343/369 (C-III). Several atypical CdtLoc features are observed; *cdtR* is absent in ST649, and an additional copy of *cdtA* is present in ST369, the latter comprising part of a CdtLoc co-located with the PaLoc. (D) Amino acid differences in TcdB among cryptic STs 649, 637, and 369 and reference strains from clades C1-5. Variations are shown as black lines relative to CD630 (C1, ST54). Phylogenies constructed from the catalytic and protease domains (in blue) and translocation and receptor-

binding domains (in orange) of TcdB for the same eight STs included in (D). Scale bar shows the number of amino acid substitutions per site. Trees are mid-point rooted and supported by 500 bootstrap replicates.

Critically, at least one ST in each of clades C-I, C-II and C-III harboured divergent *tcdB* (89-94% identity to *tcdB*_{R20291}) and/or *cdtAB* alleles (60-71% identity to *cdtA*_{R20291}, 74-81% identity to *cdtB*_{R20291}). These genes were located on atypical and novel PaLoc and CdtLoc structures flanked by mediators of lateral gene transfer (Fig. 7). Sequence types 359, 360, 361 and 649 (C-I), 637 (C-II) and 369 (C-III) harboured ‘monotoxin’ PaLocs characterised by the presence of syntenic *tcdR*, *tcdB* and *tcdE*, and complete absence of *tcdA* and *tcdC*. In STs 360 and 361 (C-I), and 637 (C-II), a gene encoding an endolysin with predicted N-acetylmuramoyl-L-alanine amidase activity (*cwlH*) was found adjacent to the phage-derived holin gene *tcdE*.

Remarkably, a full CdtLoc was found upstream of the PaLoc in ST369 (C-III). This CdtLoc was unusual, characterised by the presence of *cdtB*, two copies of *cdtA*, two copies of *cdtR* and *xerC* encoding a site-specific tyrosine recombinase (Fig. 7). Both ST644 (C-I) and ST343 (C-III) were CdtLoc-positive but PaLoc-negative (A⁺B⁺CDT⁺). In ST649 (C-I) *cdtR* was completely absent and, in ST343 (C-III), the entire CdtLoc was contained within the genome of a 56Kbp temperate bacteriophage termed ΦSemix9P1²⁹. Toxin regulators TcdR and CdtR are highly conserved across clades C1-5²¹. In contrast, the CdtR of STs 644 (C-I), 343 (C-III) and 369 (C-III) shared only 46-54% amino acid identity (AAI) with CdtR of strain R20291 from clade 2 and ~40% AAI to each other. Similarly, the TcdR of ST 369 shared only 82.1% AAI compared to R20291 (Supplementary File 1m).

Compared to TcdB of R20291 (TcdB_{R20291}), the shared AAI for TcdB_{ST649_C-I}, TcdB_{ST637_C-II} and TcdB_{ST369_C-III} were 94.0%, 90.5% and 89.4%, respectively. This sequence heterogeneity was confirmed through the detection of five distinct *HincII/AccI* digestion profiles of *tcdB* B1 fragments possibly reflecting novel toxinotypes (Supplementary File 1n). TcdB phylogenies identified clade C2 as the most recent common ancestor for TcdB_{ST649_C-I} (Fig. 7). Phylogenetic subtyping analysis of the TcdB receptor-binding domain (RBD) showed the respective sequences in C-I, C-II and C-III clustered with *tcdB* alleles belonging to virulent C2 strains (Supplementary File 1o). Notably, the TcdB-RBD of ST649 (C-I) shared an AAI of 93.5% with TcdB-RBD allele type 8 belonging to hypervirulent STs 1 (RT027)¹³ and 231 (RT251)³⁰. Similarly, the closest match to *tcdB*-RBDs of ST637 (C-II) and ST369 (C-III) was allele type 10 (ST41, RT244)³¹.

Discussion

Through phylogenomic analysis of the largest and most diverse collection of *C. difficile* genomes to date, we identified major incoherence in *C. difficile* taxonomy, provide the first WGS-based phylogeny for the *Peptostreptococcaceae* and provide new insight into intra-species diversity and evolution of pathogenicity in this major One Health pathogen.

Our analysis found high nucleotide identity (ANI > 97%) between *C. difficile* clades C1-4, indicating that strains from these four clades (comprising 560 known STs) belong to the same species. On the other hand, ANI between C5 and C1-4 is on the borderline of the accepted species threshold (95.9-96.2%). This degree of speciation likely reflects the unique ecology of C5 – a lineage comprising 33 known STs which is well established in non-human animal reservoirs worldwide and associated with CDI in the community setting³². Conversely, we identified major taxonomic incoherence among the three cryptic clades and C1-5, evident by ANI values (compared to ST3, C1) far below the species threshold (~91%, C-I; ~94%, C-II; and ~89%, C-III). Similar ANI value differences were seen between the cryptic clades themselves, indicating they are as divergent from each other as they are individually from C1-5. This extraordinary level of discontinuity is substantiated by our core genome and Bayesian analyses. Our study estimated the most recent common ancestor of *C. difficile* clades C1-4 and C1-5 existed between 0.46 to 2.77 mya and between 1.11 to 6.71 mya, respectively, whereas the common ancestors of clades C-I, C-II and C-III were estimated to have existed at least 1.5 to 75 million years before the common ancestor of C1-5. For context, divergence dates for other notable pathogens range from 10 Ma (*Campylobacter coli*

307 and *C. jejuni*)³³, 47 Ma (*Burkholderia pseudomallei* and *B. thailandensis*)³⁴ and 120 Ma
308 (*Escherichia coli* and *Salmonella enterica*)³⁵. Corresponding whole-genome ANI values for these
309 species are 86%, 94% and 82%, respectively (**Supplementary File 1j**).

310 Although BEAST provided wider confidence intervals (and therefore less certainty compared to
311 BactDating), it estimates the time of divergence for all clades within the same order of magnitude
312 and, importantly, provides robust support for the same branching order of clades with clade C-III
313 the most ancestral of lineages, followed by the emergence of C-I, C-II, and C5. After this point,
314 there appears to have been rapid population expansion into the four closely related clades described
315 today, which include many of the most prevalent strains causing healthcare-associated CDI
316 worldwide¹¹.

317 We acknowledge that the dating of ancient taxa is often imprecise and that using a strict
318 clock model for such a diverse set of taxa leads to considerable uncertainty in divergence estimates.
319 However, we tried to mitigate this as much as possible by using two independent tools and
320 evaluated multiple molecular clock estimates (covering almost an order of magnitude), ultimately
321 using the same fixed clock model as Kumar *et al.*²⁷ (2.5×10^{-9} – 1.5×10^{-8}). The branching order of the
322 clades is robust, supported by comprehensive and independent comparative genomic and
323 phylogenomic analyses. Notwithstanding this finding, if variations in the molecular clock happen
324 over time and across lineages, which is likely the case for such a genetically diverse spore-forming
325 pathogen, then the true age ranges for *C. difficile* clade emergence are likely far greater (and
326 therefore less certain) than we report here.

327 Comparative ANI analysis of the cryptic clades with >5000 reference genomes across 21
328 phyla failed to provide a better match than *C. difficile* (89–94% ANI). Similarly, our revised ANI-
329 based taxonomy of the *Peptostreptococcaceae* placed clades C-I, C-II and C-III between *C. difficile*
330 and *C. manganotii*. Our analyses of the *Clostridioides* spp. highlights the major discordance
331 between WGS data and 16S rRNA data which has historically been used to classify bacterial
332 species. In 2016, Lawson *et al.*²³ used 16S rRNA data to categorise *C. difficile* and *C. manganotii* as
333 the sole members of the *Clostridioides*. These species have 94.7% similarity in 16S rRNA sequence
334 identity, yet our findings indicate that *C. manganotii* and *C. difficile* share 77% ANI and should not
335 be considered within the same genus. The rate of 16S rRNA divergence in bacteria is estimated to
336 be 1–2% per 50 Ma³⁵. Contradicting our ANI and core genome data, 16S rRNA sequences were
337 highly conserved across all 8 clades. This indicates that in *C. difficile*, 16S rRNA gene similarity
338 correlates poorly with measures of genomic, phenotypic and ecological diversity, as reported in
339 other taxa such as *Streptomyces*, *Bacillus* and *Enterobacteriaceae*^{36, 37}. Another interesting
340 observation is that C5 and the three cryptic clades had a high proportion (>90%) of MLST alleles
341 that were absent in other clades (**Supplementary File 1e**) suggesting minimal exchange of essential
342 housekeeping genes between these clades. Whether this reflects divergence or convergence of two
343 species, as seen in *Campylobacter*³⁸, is unknown. Taken together, these data strongly support the
344 reclassification of *C. difficile* clades C-I, C-II and C-III as novel independent *Clostridioides*
345 genomospecies. There have been similar genome-based reclassifications in *Bacillus*³⁹,
346 *Fusobacterium*⁴⁰ and *Burkholderia*⁴¹. Also, a recent Consensus Statement⁴² argues that the
347 genomics and big data era necessitate easing of nomenclature rules to accommodate genome-based
348 assignment of species status to nonculturable bacteria and those without ‘type material’, as is the
349 case with these genomospecies.

350 We also found that the significant taxonomic incoherence observed in *C. difficile* was also
351 evident in other medically important clostridia, supporting calls for taxonomic revisions^{23, 24}. The
352 entire published collections of *C. perfringens*, *C. sporogenes* and *C. botulinum* all contained
353 sequenced strains with pairwise ANI below the 96% demarcation threshold, with 8% of
354 *C. sporogenes* and 31% of *C. botulinum* sequenced strains below 90% ANI. These findings
355 highlight a significant problem with the current classification of the clostridia and further
356 demonstrate that high-resolution approaches such as whole-genome ANI can be a powerful tool for
357 the re-classification of these bacteria^{23, 24, 42}.

358 The NCBI SRA was dominated by C1 and C2 strains, both in number and diversity. This
 359 apparent bias reflects the research community's efforts to sequence the most prominent strains
 360 causing CDI in regions with the highest-burden, e.g. ST 1 from humans in Europe and North
 361 America. As such, there is a paucity of sequenced strains from diverse environmental sources,
 362 animal reservoirs or regions associated with atypical phenotypes. Cultivation bias - a historical
 363 tendency to culture, preserve and ultimately sequence *C. difficile* isolates that are concordant with
 364 expected phenotypic criteria, comes at the expense of 'outliers' or intermediate phenotypes.
 365 Members of the cryptic clades fit this criterion. They were first identified in 2012 but have been
 366 overlooked due to atypical toxin architecture which may compromise diagnostic assays (discussed
 367 below). Our updated MLST phylogeny shows as many as 55 STs across the three cryptic clades (C-
 368 I, n=25; C-II, n=9; C-III, n=21) (**Fig. 2**). There remains a further dozen 'outliers' that could either
 369 fit within these new taxa or be the first typed representative of additional genomospecies. The
 370 growing popularity of metagenomic sequencing of animal and environmental microbiomes will
 371 certainly identify further diversity within these taxa, including nonculturable strains^{43, 44}.

372 By analysing 260 STs across eight clades, we provide the most comprehensive pangenome
 373 analysis of *C. difficile* to date. Importantly, we also show that the choice of algorithm significantly
 374 affects pangenome estimation. The *C. difficile* pangenome was determined to be open (i.e., an
 375 unlimited gene repertoire) and vast in scale (over 17000 genes), much larger than previous
 376 estimates (~10000 genes) which mainly considered individual clonal lineages^{16, 22}. Conversely,
 377 comprising just 12.8% of its genetic repertoire (2,232 genes), the core genome of *C. difficile* is
 378 remarkably small, consistent with earlier WGS and microarray-based studies describing ultralow
 379 genome conservation in *C. difficile*^{11, 45}. Considering only C1-5, the pangenome reduced in size by
 380 12% (2,082 genes); another 519 genes were lost when considering only C1-4. These findings are
 381 consistent with our taxonomic data, suggesting the cryptic clades, and to a lesser extent C5,
 382 contribute a significant proportion of evolutionarily divergent and unique loci to the gene pool. A
 383 large open pangenome and small core genome are synonymous with a sympatric lifestyle,
 384 characterised by cohabitation with, and extensive gene transfer between, diverse communities of
 385 prokarya and archaea⁴⁶. Indeed, *C. difficile* shows a highly mosaic genome comprising many
 386 phages, plasmids and integrative and conjugative elements¹¹, and has adapted to survival in multiple
 387 niches including the mammalian gastrointestinal tract, water, soil and compost, and invertebrates³².

388 Through a robust Pan-GWAS approach we identified loci that are enriched or unique in the
 389 genomospecies. C-I strains were associated with the presence of transporter AbgB and absence of a
 390 mannose-type phosphotransferase (PTS) system. In *E. coli*, AbgAB proteins allow it to survive on
 391 exogenous sources of folate⁴⁷. In many enteric species, the mannose-type PTS system is essential
 392 for catabolism of fructosamines such as glucoselysine and fructoselysine, abundant components of
 393 rotting fruit and vegetable matter⁴⁸. C-II strains contained Zn transporter loci *znuA* and *yeiR*, in
 394 addition to Zn transporter ZupT which is highly conserved across all eight *C. difficile* clades.
 395 *S. enterica* and *E. coli* harbour both *znuA/yeiR* and ZupT loci, enabling survival in Zn-depleted
 396 environments⁴⁹. C-III strains were associated with major gene clusters encoding systems for
 397 ethanolamine catabolism, heavy metal transport and spermidine uptake. The C-III *eut* gene cluster
 398 encoded six additional kinases, transporters and transcription regulators absent from the highly
 399 conserved *eut* operon found in other clades. Ethanolamine is a valuable source of carbon and/or
 400 nitrogen for many bacteria, and *eut* gene mutations (in C1/C2) impact toxin production *in vivo*⁵⁰.
 401 The C-III metal transport gene cluster encoded a chelator of heavy metal ions and a multi-
 402 component transport system with specificity for iron, nickel, and glutathione. The conserved
 403 spermidine operon found in all *C. difficile* clades is thought to play an important role in various
 404 stress responses including during iron limitation⁵¹. The additional, divergent spermidine
 405 transporters found in C-III were similar to regions in closely related genera *Romboutsia* and
 406 *Paenibacillus* (data not shown). Together, these data provide preliminary insights into the
 407 biology and ecology of the genomospecies. Most differential loci identified were responsible for
 408 extra or alternate metabolic processes, some not previously reported in *C. difficile*. It is therefore

tempting to speculate that the evolution of alternate biosynthesis pathways in these species reflects distinct ancestries and metabolic responses to evolving within markedly different ecological niches.

This work demonstrates the presence of toxin genes on PaLoc and CdtLoc structures in all three genomospecies, confirming their clinical relevance. Monotoxin PaLocs were characterised by the presence of *tcdR*, *tcdB* and *tcdE*, the absence of *tcdA* and *tcdC*, and flanking by transposases and recombinases which mediate LGT^{20, 21, 52}. These findings support the notion that the classical biotoxin PaLoc common to clades C1-5 was derived by multiple independent acquisitions and stable fusion of monotoxin PaLocs from ancestral Clostridia⁵². Moreover, the presence of syntenic PaLoc and CdtLoc (in ST369, C-I), the latter featuring two copies of *cdtA* and *cdtR*, and a recombinase (*xerC*), further support this PaLoc fusion hypothesis⁵².

Bacteriophage holin and endolysin enzymes coordinate host cell lysis, phage release and toxin secretion⁵³. Monotoxin PaLocs comprising phage-derived holin (*tcdE*) and endolysin (*cwlH*) genes were first described in C-I strains⁵². We have expanded this previous knowledge by demonstrating that syntenic *tcdE* and *cwlH* are present within monotoxin PaLocs across all three genomospecies. Moreover, since some strains contained *cwlH* but lacked toxin genes, this gene seems to be implicated in toxin acquisition. These data, along with the detection of a complete and functional²⁹ CdtLoc contained within ΦSemix9P1 in ST343 (C-III), further substantiate the role of phages in the evolution of toxin loci in *C. difficile* and related Clostridia⁵³.

The CdtR and TcdR sequences of the new genomospecies are unique and further work is needed to determine if these regulators display different mechanisms or efficiencies of toxin expression¹². The presence of dual copies of CdtR in ST369 (C-I) is intriguing, as analogous duplications in PaLoc regulators have not been documented. One of these CdtR had a mutation at a key phosphorylation site (Asp61→Asn61) and possibly shows either reduced wild-type activity or non-functionality, as seen in ST11⁵⁴. This might explain the presence of a second CdtR copy.

TcdB alone can induce host innate immune and inflammatory responses leading to intestinal and systemic organ damage⁵⁵. Our phylogenetic analysis shows TcdB sequences from the three genomospecies are related to TcdB in Clade 2 members, specifically ST1 and ST41, both virulent lineages associated with international CDI outbreaks^{13, 31}, and causing classical or variant (*C. sordellii*-like) cytopathic effects, respectively⁵⁶. It would be relevant to explore whether the divergent PaLoc and CdtLoc regions confer differences in biological activity, as these may present challenges for the development of effective broad-spectrum diagnostic assays, and vaccines. We have previously demonstrated that common laboratory diagnostic assays may be challenged by changes in the PaLoc of C-I strains²¹. The same might be true for monoclonal antibody-based treatments for CDI such as bezlotoxumab, known to have distinct neutralizing activities against different TcdB subtypes⁵⁷.

Our findings highlight major incongruence in *C. difficile* taxonomy, identify differential patterns of diversity among major clades and advance understanding of the evolution of the PaLoc and CdtLoc. While our analysis is limited solely to the genomic differences between *C. difficile* clades, our data provide a robust genetic foundation for future studies to focus on the phenotypic, ecological and epidemiological features of these interesting groups of strains, including defining the biological consequences of clade-specific genes and pathogenic differences *in vitro* and *in vivo*. Our findings reinforce that the epidemiology of this important One Health pathogen is not fully understood. Enhanced surveillance of CDI and WGS of new and emerging strains to better inform the design of diagnostic tests and vaccines are key steps in combating the ongoing threat posed by *C. difficile*. Last, besides *C. difficile*, we also demonstrate that a similar approach can be applied to other clostridia making a useful tool for the reclassification of these bacteria.

Materials and Methods

Key Resources Table				
Reagent type or resource	Designation	Source or reference	Identifiers	Additional information

software, algorithm	ABRicate	https://github.com/tseemann/abricate	RRID:SCR_021093	
software, algorithm	ACT: Artemis Comparison Tool	http://www.sanger.ac.uk/resources/software/act/	RRID:SCR_004507	
software, algorithm	BactDating	https://github.com/xavierdidelot/BactDating	RRID:SCR_021092	
software, algorithm	BEAST	http://beast.bio.ed.ac.uk/	RRID:SCR_010228	
software, algorithm	Clustal Omega	http://www.ebi.ac.uk/Tools/msa/clustalo/	RRID:SCR_001591	
software, algorithm	Easyfig	http://easyfig.sourceforge.net/	RRID:SCR_013169	
software, algorithm	FastANI	https://github.com/ParBLISS/FastANI	RRID:SCR_021091	
software, algorithm	Geneious	http://www.geneious.com/	RRID:SCR_010519	
software, algorithm	Gubbins	https://sanger-pathogens.github.io/gubbins/	RRID:SCR_016131	
software, algorithm	iToL	https://itol.embl.de/	RRID:SCR_018174	
other	KEGG	http://www.kegg.jp/	RRID:SCR_012773	online database
software, algorithm	Kraken2	http://www.ebi.ac.uk/research/enright/software/kraken	RRID:SCR_005484	
software, algorithm	MAFFT	http://mafft.cbrc.jp/alignment/server/	RRID:SCR_011811	
software, algorithm	MEGA	http://megasoftware.net/	RRID:SCR_000667	
software, algorithm	MUSCLE	http://www.ebi.ac.uk/Tools/msa/muscle/	RRID:SCR_011812	
other	NCBI RefSeq	http://130.14.29.110/BLAST/	RRID:SCR_008420	online database
other	NCBI Sequence Read Archive	http://www.ncbi.nlm.nih.gov/sra	RRID:SCR_004891	online database
software, algorithm	Panaroo	https://github.com/gtonkinhill/panaroo	RRID:SCR_021090	
software, algorithm	PanGP	https://pangp.zhaopage.com/	RRID:SCR_021089	
software, algorithm	Phandango	http://phandango.net/	RRID:SCR_015243	

software, algorithm	Prokka	http://www.vicbioinformatics.com/software.prokka.shtml	RRID:SCR_014732	
other	PubMLST	http://pubmlst.org/	RRID:SCR_012955	online database
software, algorithm	pyani	https://pypi.org/project/pyani/	RRID:SCR_021088	
software, algorithm	QUAST	http://bioinf.spbau.ru/quast	RRID:SCR_001228	
software, algorithm	RAxML	https://github.com/stamatak/standard-RAxML	RRID:SCR_006086	
software, algorithm	Roary	https://sanger-pathogens.github.io/Roary/	RRID:SCR_018172	
software, algorithm	Scoary	https://github.com/AdmiralEnOla/Scoary	RRID:SCR_021087	
software, algorithm	SPAdes	http://bioinf.spbau.ru/spades/	RRID:SCR_000131	
software, algorithm	SPSS	https://www.ibm.com/products/spss-statistics	RRID:SCR_019096	
software, algorithm	SRST2	https://github.com/katholt/srst2	RRID:SCR_015870	
software, algorithm	TrimGalore	http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/	RRID:SCR_011847	

456 **Genome collection.** We retrieved the entire collection of *C. difficile* genomes (taxid ID 1496) held
457 at the NCBI Sequence Read Archive [<https://www.ncbi.nlm.nih.gov/sra/>]. The raw dataset (as of 1st
458 January 2020) comprised 12,621 genomes. After filtering for redundancy and Illumina paired-end
459 data (all platforms and read lengths), 12,304 genomes (97.5%) were available for analysis.

460 **Multi-locus sequence typing.** Sequence reads were interrogated for multi-locus sequence type (ST)
461 using SRST2 v0.1.8⁵⁸. New alleles, STs and clade assignments were verified by submission of
462 assembled contigs to PubMLST [<https://pubmlst.org/cdifficile/>]. A species-wide phylogeny was
463 generated from 659 ST alleles sourced from PubMLST (dated 01-Jan-2020). Alleles were
464 concatenated in frame and aligned with MAFFT v7.304. A final neighbour-joining tree was
465 generated in MEGA v10⁵⁹ and annotated using iTOL v4 [<https://itol.embl.de/>].

466 **Genome assembly and quality control.** Genomes were assembled, annotated and evaluated using
467 a pipeline comprising TrimGalore v0.6.5, SPAdes v3.6.043, Prokka v1.14.5, and QUAST v2.344¹⁶.
468 Next, Kraken2 v2.0.8-beta⁶⁰ was used to screen for contamination and assign taxonomic labels to
469 reads and draft assemblies. Based on metadata, read depth and assembly quality, a final dataset of
470 260 representative genomes of each ST present in the ENA were used for all subsequent
471 bioinformatics analyses (C1, n=149; C2, n=35; C3, n=7; C4, n=34; C5, n=18; C-I, n=12; C-II, n=3,
472 C-III, n=2). The list of representative genomes is available in **Table S2** in **Supplementary File 1b**.

473 **Taxonomic analyses.** Species-wide genetic similarity was determined by computation of whole-
474 genome ANI for 260 STs. Both alignment-free and conventional alignment-based ANI approaches
475 were taken, implemented in FastANI⁵ v1.3 and the Python module pyani⁶¹ v0.2.9, respectively.
476 FastANI calculates ANI using a unique *k*-mer based alignment-free sequence mapping engine,

whilst pyani utilises two different classical alignment ANI algorithms based on BLAST+ (ANiB) and MUMmer (ANIm). A 96% ANI cut-off was used to define species boundaries⁴. For taxonomic placement, ANI was determined for divergent *C. difficile* genomes against two datasets comprising (i) members of the *Peptostreptococcaceae* (n=25)²³, and (ii) the complete NCBI RefSeq database (n=5895 genomes, <https://www.ncbi.nlm.nih.gov/refseq/>, accessed 14th Jan 2020). Finally, comparative identity analysis of consensus 16S rRNA sequences for *C. mangenotii* type strain DSM1289T²³ (accession FR733662.1) and representatives of each *C. difficile* clade was performed using Clustal Omega <https://www.ebi.ac.uk/Tools/msa/clustalo/>.

Estimates of clade and species divergence. BactDating v1.0.1⁶² was applied to the recombination-corrected phylogeny produced by Gubbins (471,708 core-genome sites) with Markov chain Monte Carlo (MCMC) chains of 10^7 iterations sampled every 10^4 iterations with a 50% burn-in. A strict clock model was used with a rate of 2.5×10^{-9} to 1.5×10^{-8} substitutions per site per year, as previously defined by He *et al.*¹³ and Kumar *et al.*²⁷. The effective sample sizes (ESS) were >200 for all estimated parameters, and traces were inspected manually to ensure convergence. To provide an independent estimate from BactDating, BEAST v1.10.4⁶³ was run on a recombination-filtered gap-free alignment of 10,466 sites with MCMC chains of 5×10^8 iterations, with a 9×10^{-7} burn-in, that were sampled every 10^4 iterations. The strict clock model described above was used in combination with the discrete GTR gamma model of heterogeneity among sites and skyline population model. MCMC convergence was verified with Tracer v1.7.1 and ESS for all estimated parameters were >150. For ease of comparison, clade dating from both approaches were transposed onto a single MLST phylogeny. Tree files are available as **Supplementary Files 2-4** at <http://doi.org/10.6084/m9.figshare.12471461>.

Pangenome analysis. The 260 ST dataset was used for pangenome analysis with Panaroo v1.1.0⁶⁴ and Roary v3.6.0⁶⁵. Panaroo was run with default thresholds for core assignment (98%) and blastP identity (95%). Roary was run with a default threshold for core assignment (99%) and two different thresholds for BlastP identity (95%, 90%). Sequence alignment of the final set of core genes (Panaroo; n=2,232 genes, 2,606,142 bp) was performed using MAFFT v7.304 and recombinative sites were filtered using Gubbins v7.304⁶⁶. A recombinant adjusted alignment of 471,708 polymorphic sites was used to create a core genome phylogeny with RAxML v8.2.12 (GTR gamma model of among-site rate-heterogeneity), which was visualized alongside pangenome data in Phandango⁶⁷. Pangenome dynamics were investigated with PanGP v1.0.1¹⁶.

Scoary⁶⁸ v1.6.16 was used to identify genetic loci that were statistically associated with each clade via a Pangenome-Wide Association Study (pan-GWAS). The Panaroo-derived pangenome (n=17,470) was used as input for Scoary with the evolutionary clade of each genome depicted as a discrete binary trait. Scoary was run with 1,000 permutation replicates and genes were reported as significantly associated with a trait if they attained *p*-values (empirical, naïve, and Benjamini-Hochberg-corrected) of ≤ 0.05 , a sensitivity and specificity of > 99% and 97.5%, respectively, and were not annotated as “hypothetical proteins”. All significantly associated genes were reannotated using prokka and BlastP and functional classification (KEGG orthology) was performed using the Koala suite of web-based annotation tools⁶⁹.

Comparative analysis of toxin gene architecture. The 260 ST genome dataset was screened for the presence of *tcdA*, *tcdB*, *cdtA* and *cdtB* using the Virulence Factors Database (VFDB) compiled within ABRicate v1.0 [<https://github.com/tseemann/abricate>]. Results were corroborated by screening raw reads against the VFDB using SRST2 v0.1.8⁵⁸. Both approaches employed minimum coverage and identity thresholds of 90 and 75%, respectively. Comparative analysis of PaLoc and CdtLoc architecture was performed by mapping of reads with Bowtie2 v2.4.1 to cognate regions in reference strain R20291 (ST1, FN545816). All PaLoc and CdtLoc loci investigated showed sufficient coverage for accurate annotation and structural inference. Genome comparisons were visualized using ACT and figures prepared with Easyfig²¹. MUSCLE-aligned TcdB sequences were visualized in Geneious v2020.1.2 and used to create trees in iTOL v4.

527 **Statistical analyses.** All statistical analyses were performed using SPSS v26.0 (IBM, NY, USA).
 528 For pangenome analyses, a Chi-squared test with Yate's correction was used to compare the
 529 proportion of core genes and a One-tailed Mann-Whitney U test was used to demonstrate the
 530 reduction of gene content per genome, with a p-value ≤ 0.05 considered statistically significant.

531 **References**

- 532 1. Doolittle WF, Papke RT. Genomics and the bacterial species problem. *Genome Biol* **7**, 116
 533 (2006).
- 534 2. Konstantinidis KT, Ramette A, Tiedje JM. The bacterial species definition in the genomic era.
 535 *Philos Trans R Soc Lond B Biol Sci* **361**, 1929-1940 (2006).
- 536 3. Wayne LG, *et al.* Report of the ad hoc committee on reconciliation of approaches to bacterial
 537 systematics. *Int J Syst Evol Microbiol* **37**, 463-464 (1987).
- 538 4. Ciufo S, *et al.* Using average nucleotide identity to improve taxonomic assignments in
 539 prokaryotic genomes at the NCBI. *Int J Syst Evol Microbiol* **68**, 2386-2392 (2018).
- 540 5. Jain C, Rodriguez RL, Phillippy AM, Konstantinidis KT, Aluru S. High throughput ANI
 541 analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun* **9**, 5114
 542 (2018).
- 543 6. Richter M, Rossello-Mora R. Shifting the genomic gold standard for the prokaryotic species
 544 definition. *Proc Natl Acad Sci U S A* **106**, 19126-19131 (2009).
- 545 7. Guh AY, *et al.* Trends in US burden of *Clostridioides difficile* infection and outcomes. *N Engl*
 546 *J Med* **382**, 1320-1330 (2020).
- 547 8. CDC. Antibiotic resistance threats in the United States, 2013. Centers for Disease Control and
 548 Prevention. Web citation: <http://www.cdc.gov/drugresistance/threat-report-2013/>. (2013).
- 549 9. CDC. Antibiotic resistance threats in the United States, 2019. Centers for Disease Control and
 550 Prevention. Web citation: <https://www.cdc.gov/drugresistance/biggest-threats.html>. (2019).
- 551 10. Lim S, Knight D, Riley T. *Clostridium difficile* and One Health. *Clinical Microbiology and*
 552 *Infection*, (2019).
- 553 11. Knight DR, Elliott B, Chang BJ, Perkins TT, Riley TV. Diversity and evolution in the genome
 554 of *Clostridium difficile*. *Clin Microbiol Rev* **28**, 721-741 (2015).
- 555 12. Chandrasekaran R, Lacy DB. The role of toxins in *Clostridium difficile* infection. *FEMS*
 556 *Microbiol Rev* **41**, 723-750 (2017).
- 557 13. He M, *et al.* Emergence and global spread of epidemic healthcare-associated *Clostridium*
 558 *difficile*. *Nat Genet* **45**, 109-113 (2013).
- 559 14. Shaw HA, *et al.* The recent emergence of a highly related virulent *Clostridium difficile* clade
 560 with unique characteristics. *Clin Microbiol Infect* **26**, 492-498 (2020).
- 561 15. Imwattana K, *et al.* *Clostridium difficile* ribotype 017 - characterization, evolution and
 562 epidemiology of the dominant strain in Asia. *Emerg Microb Infect* **8**, 796-807 (2019).

16. Knight DR, *et al.* Evolutionary and genomic insights into *Clostridioides difficile* sequence type 11: a diverse, zoonotic and antimicrobial resistant lineage of global One Health importance. *MBio* **10**, e00446-00419 (2019).
17. Dingle KE, *et al.* Evolutionary history of the *Clostridium difficile* pathogenicity locus. *Genome Biol Evol* **6**, 36-52 (2014).
18. Didelot X, *et al.* Microevolutionary analysis of *Clostridium difficile* genomes to investigate transmission. *Genome Biol* **13**, R118 (2012).
19. Janezic S, Potocnik M, Zidaric V, Rupnik M. Highly divergent *Clostridium difficile* strains isolated from the environment. *PLoS One* **11**, e0167101 (2016).
20. Ramirez-Vargas G, Rodriguez C. Putative conjugative plasmids with *tcdB* and *cdtAB* genes in *Clostridioides difficile*. *Clin Infect Dis* **26**, 2287-2290 (2020).
21. Ramirez-Vargas G, *et al.* Novel Clade CI *Clostridium difficile* strains escape diagnostic tests, differ in pathogenicity potential and carry toxins on extrachromosomal elements. *Sci Rep* **8**, 1-11 (2018).
22. Knight DR, Squire MM, Collins DA, Riley TV. Genome analysis of *Clostridium difficile* PCR ribotype 014 lineage in Australian pigs and humans reveals a diverse genetic repertoire and signatures of long-range interspecies transmission. *Front Microbiol* **7**, 2138 (2017).
23. Lawson PA, Citron DM, Tyrrell KL, Finegold SM. Reclassification of *Clostridium difficile* as *Clostridioides difficile* (Hall and O'Toole 1935) Prevot 1938. *Anaerobe* **40**, 95-99 (2016).
24. Oren A, Rupnik M. *Clostridium difficile* and *Clostridioides difficile*: Two validly published and correct names. *Anaerobe* **52**, 125-126 (2018).
25. Knetsch CW, *et al.* Comparative analysis of an expanded *Clostridium difficile* reference strain collection reveals genetic diversity and evolution through six lineages. *Infect Genet Evol* **12**, 1577-1585 (2012).
26. He M, *et al.* Evolutionary dynamics of *Clostridium difficile* over short and long time scales. *Proc Natl Acad Sci U S A* **107**, 7527-7532 (2010).
27. Kumar N, *et al.* Adaptation of host transmission cycle during *Clostridium difficile* speciation. *Nat Genet* **51**, 1315-1320 (2019).
28. Tettelin H, *et al.* Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the microbial "pan-genome". *Proc Natl Acad Sci U S A* **102**, 13950-13955 (2005).
29. Riedel T, *et al.* A *Clostridioides difficile* bacteriophage genome encodes functional binary toxin-associated genes. *J Biotechnol* **250**, 23-28 (2017).
30. Hong S, Knight DR, Chang B, Carman RJ, Riley TV. Phenotypic characterisation of *Clostridium difficile* PCR ribotype 251, an emerging multi-locus sequence type clade 2 strain in Australia. *Anaerobe* **60**, 102066 (2019).

31. Eyre DW, *et al.* Emergence and spread of predominantly community-onset *Clostridium difficile* PCR ribotype 244 infection in Australia, 2010 to 2012. *Euro Surveill* **20**, 21059 (2015).
32. Knight DR, Riley TV. Genomic delineation of zoonotic origins of *Clostridium difficile*. *Front Pub Health* **7**, 164 (2019).
33. Sheppard SK, Maiden MC. The evolution of *Campylobacter jejuni* and *Campylobacter coli*. *Cold Spring Harb Perspect Biol* **7**, a018119 (2015).
34. Yu Y, *et al.* Genomic patterns of pathogen evolution revealed by comparison of *Burkholderia pseudomallei*, the causative agent of melioidosis, to avirulent *Burkholderia thailandensis*. *BMC Microbiol* **6**, 46 (2006).
35. Ochman H, Elwyn S, Moran NA. Calibrating bacterial evolution. *Proc Natl Acad Sci U S A* **96**, 12638-12643 (1999).
36. Janda JM, Abbott SL. 16S rRNA gene sequencing for bacterial identification in the diagnostic laboratory: pluses, perils, and pitfalls. *J Clin Microbiol* **45**, 2761-2764 (2007).
37. Chevrette MG, Carlos-Shanley C, Louie KB, Bowen BP, Northen TR, Currie CR. Taxonomic and metabolic incongruence in the ancient genus *Streptomyces*. *Front Microbiol* **10**, 2170 (2019).
38. Sheppard SK, McCarthy ND, Falush D, Maiden MC. Convergence of *Campylobacter* species: implications for bacterial evolution. *Science* **320**, 237-239 (2008).
39. Liu Y, Lai QL, Shao ZZ. Genome analysis-based reclassification of *Bacillus weihenstephanensis* as a later heterotypic synonym of *Bacillus mycoides*. *Int J Syst Evol Microbiol* **68**, 106-112 (2018).
40. Kook JK, *et al.* Genome-based reclassification of *Fusobacterium nucleatum* subspecies at the species level. *Curr Microbiol* **74**, 1137-1147 (2017).
41. Loveridge EJ, *et al.* Reclassification of the specialized metabolite producer *Pseudomonas mesoacidophila* ATCC 31433 as a member of the *Burkholderia cepacia* complex. *J Bacteriol* **199**, e00125-00117 (2017).
42. Murray AE, *et al.* Roadmap for naming uncultivated Archaea and Bacteria. *Nat Microbiol* **5**, 987-994 (2020).
43. Stewart RD, *et al.* Assembly of 913 microbial genomes from metagenomic sequencing of the cow rumen. *Nat Commun* **9**, 1-11 (2018).
44. Lu X, *et al.* Bacterial pathogens and community composition in advanced sewage treatment systems revealed by metagenomics analysis based on high-throughput sequencing. *PLoS One* **10**, e0125549 (2015).
45. Scaria J, Ponnala L, Janvilisri T, Yan W, Mueller LA, Chang YF. Analysis of ultra low genome conservation in *Clostridium difficile*. *PLoS One* **5**, e15147 (2010).
46. Medini D, Donati C, Tettelin H, Massignani V, Rappuoli R. The microbial pan-genome. *Curr Opin Genet Dev* **15**, 589-594 (2005).

47. Carter EL, Jager L, Gardner L, Hall CC, Willis S, Green JM. *Escherichia coli* abg genes enable uptake and cleavage of the folate catabolite p-aminobenzoyl-glutamate. *J Bacteriol* **189**, 3329-3334 (2007).
48. Miller KA, Phillips RS, Kilgore PB, Smith GL, Hoover TR. A mannose family phosphotransferase system permease and associated enzymes are required for utilization of fructoselysine and glucoselysine in *Salmonella enterica* serovar Typhimurium. *J Bacteriol* **197**, 2831-2839 (2015).
49. Sabri M, Houle S, Dozois CM. Roles of the extraintestinal pathogenic *Escherichia coli* ZnuACB and ZupT zinc transporters during urinary tract infection. *Infect Immun* **77**, 1155-1164 (2009).
50. Nawrocki KL, Wetzel D, Jones JB, Woods EC, McBride SM. Ethanolamine is a valuable nutrient source that impacts *Clostridium difficile* pathogenesis. *Environ Microbiol* **20**, 1419-1435 (2018).
51. Berges M, *et al.* Iron regulation in *Clostridioides difficile*. *Front Microbiol* **9**, 3183 (2018).
52. Monot M, *et al.* *Clostridium difficile*: new insights into the evolution of the pathogenicity locus. *Sci Rep* **5**, 15023 (2015).
53. Fortier LC. Bacteriophages contribute to shaping *Clostridioides (Clostridium) difficile* species. *Front Microbiol* **9**, 2033 (2018).
54. Bilverstone TW, Minton NP, Kuehne SA. Phosphorylation and functionality of CdtR in *Clostridium difficile*. *Anaerobe* **58**, 103-109 (2019).
55. Carter GP, *et al.* Defining the roles of TcdA and TcdB in localized gastrointestinal disease, systemic organ damage, and the host response during *Clostridium difficile* infections. *MBio* **6**, e00551 (2015).
56. Lanis JM, Barua S, Ballard JD. Variations in TcdB activity and the hypervirulence of emerging strains of *Clostridium difficile*. *PLoS Pathog* **6**, e1001061 (2010).
57. Shen E, *et al.* Subtyping analysis reveals new variants and accelerated evolution of *Clostridioides difficile* toxin B. *Commun Biol* **3**, 1-8 (2020).
58. Inouye M, *et al.* SRST2: rapid genomic surveillance for public health and hospital microbiology labs. *Genome Med* **6**, 90 (2014).
59. Kumar S, Stecher G, Li M, Knyaz C, Tamura K. MEGA X: molecular evolutionary genetics analysis across computing platforms. *Mol Biol Evol* **35**, 1547-1549 (2018).
60. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. *Genome Biol* **20**, 257 (2019).
61. Pritchard L, Glover RH, Humphris S, Elphinstone JG, Toth IK. Genomics and taxonomy in diagnostics for food security: soft-rotting enterobacterial plant pathogens. *Anal Methods* **8**, 12-24 (2016).

- 733 62. Didelot X, Croucher NJ, Bentley SD, Harris SR, Wilson DJ. Bayesian inference of ancestral
734 dates on bacterial phylogenetic trees. *Nucleic Acids Res* **46**, e134-e134 (2018).
735
- 736 63. Drummond AJ, Rambaut A. BEAST: Bayesian evolutionary analysis by sampling trees. *BMC*
737 *Evol Biol* **7**, 214 (2007).
738
- 739 64. Tonkin-Hill G, *et al.* Producing polished prokaryotic pangenomes with the panaroo pipeline.
740 *Genome Biol* **21**, 180 (2020).
741
- 742 65. Page AJ, *et al.* Roary: rapid large-scale prokaryote pangenome analysis. *Bioinformatics*
743 **31**, 3691-3693 (2015).
744
- 745 66. Croucher NJ, *et al.* Rapid phylogenetic analysis of large samples of recombinant bacterial
746 whole genome sequences using Gubbins. *Nucleic Acids Res* **43**, e15 (2015).
747
- 748 67. Hadfield J, Croucher NJ, Goater RJ, Abudahab K, Aanensen DM, Harris SR. Phandango: an
749 interactive viewer for bacterial population genomics. *Bioinformatics* **34**, 292-293 (2018).
750
- 751 68. Brynildsrud O, Bohlin J, Scheffer L, Eldholm V. Rapid scoring of genes in microbial pan-
752 genome-wide association studies with Scoary. *Genome Biol* **17**, 238 (2016).
753
- 754 69. Kanehisa M, Sato Y, Morishima K. BlastKOALA and GhostKOALA: KEGG tools for
755 functional characterization of genome and metagenome sequences. *J Mol Biol* **428**, 726-731
756 (2016).

757 **Acknowledgements**

758 This work was supported, in part, by funding from The Raine Medical Research Foundation
759 (RPG002-19) and a Fellowship from the National Health and Medical Research Council
760 (APP1138257) awarded to D.R.K. K.I. is a recipient of the Mahidol Scholarship from Mahidol
761 University, Thailand. This work was also supported by EULac project ‘Genomic Epidemiology of
762 *Clostridium difficile* in Latin America (T020076)’ and by the Millennium Science Initiative of the
763 Ministry of Economy, Development and Tourism of Chile, grant ‘Nucleus in the Biology of
764 Intestinal Microbiota’ to D.P.S. This research used the facilities and services of the Pawsey
765 Supercomputing Centre [Perth, Western Australia] and the Australian Genome Research Facility
766 [Melbourne, Victoria].

767 **Competing Interests**

768 DWE declares lecture fees from Gilead, outside the submitted work. No other author has a conflict
769 of interest to declare.

770 **Additional information**

771 Supplementary Files are available at <http://doi.org/10.6084/m9.figshare.12471461>