

1     **A global *Corynebacterium diphtheriae* genomic framework sheds light on**  
2                                   **current diphtheria reemergence**

3  
4     **Authors**

5     Melanie Hennart <sup>a,b,c</sup>, Chiara Crestani <sup>a</sup>, Sebastien Bridel <sup>a</sup>, Nathalie Armatys <sup>a,b</sup>, Sylvie  
6     Brémont <sup>a,b</sup>, Annick Carmi-Leroy <sup>a,b</sup>, Annie Landier <sup>a,b</sup>, Virginie Passet <sup>a,b</sup>, Laure Fonteneau <sup>e</sup>,  
7     Sophie Vaux <sup>e</sup>, Julie Toubiana <sup>a,b,d</sup>, Edgar Badell <sup>a,b</sup> and Sylvain Brisse <sup>a,b,\*</sup>

8  
9     **Affiliations**

10    <sup>a</sup> Institut Pasteur, Université Paris Cité, Biodiversity and Epidemiology of Bacterial  
11    Pathogens, F-75015, Paris, France

12    <sup>b</sup> Institut Pasteur, National Reference Center for Corynebacteria of the Diphtheriae Complex,  
13    Paris, France

14    <sup>c</sup> Sorbonne Université, Collège doctoral, F-75005 Paris, France

15    <sup>d</sup> Department of General Pediatrics and Pediatric Infectious Diseases, Hôpital Necker-Enfants  
16    Malades, APHP, Université de Paris, Paris, France

17    <sup>e</sup> Santé publique France, Saint-Maurice, France

18

19    **\*Correspondence:**

20    Sylvain Brisse: Institut Pasteur, Biodiversity and Epidemiology of Bacterial Pathogens, 25-28  
21    rue du Docteur Roux, F-75724, Paris, France; Phone: +33 1 45 68 83 34 ; E-mail:

22    [sylvain.brisse@pasteur.fr](mailto:sylvain.brisse@pasteur.fr)

23

24    **Keywords:** diphtheria, genomic sequencing, antimicrobial resistance, virulence,  
25    epidemiology, transmission, 2022 reemergence, bioinformatics tool

26

27    **Running Title:** Genomic surveillance of diphtheria using DIPHTOSCAN

28

## Abstract

### 29 **Background**

30 Diphtheria, caused by *Corynebacterium diphtheriae*, reemerges in Europe since 2022. Genomic  
31 sequencing can inform on transmission routes and genotypes of concern, but currently, no  
32 standard approach exists to detect clinically important genomic features and to interpret  
33 emergence in the global *C. diphtheriae* population framework.

34

### 35 **Methods**

36 We developed the bioinformatics pipeline DIPHTOSCAN (available at  
37 <https://gitlab.pasteur.fr/BEBP/diphtoscan>) to extract from genomes of *Corynebacteria* of the  
38 *diphtheriae* species complex, medically relevant features including *tox* gene presence and  
39 disruption. We analyzed 101 human *C. diphtheriae* isolates collected in 2022 in metropolitan  
40 and overseas France (France-2022). To define the population background of this emergence,  
41 we sequenced 379 additional isolates (mainly from France, 2018-2021) and collated 870  
42 publicly-available genomes.

43

### 44 **Results**

45 The France-2022 isolates comprised 45 *tox*-positive (44 toxigenic) isolates, mostly imported,  
46 belonging to 10 sublineages (<500 distinct core genes). The global dataset comprised 245  
47 sublineages and 33.9% *tox*-positive genomes, with DIPHTOSCAN predicting non-toxigenicity in  
48 16.0% of these. 12% of the global isolates, and 43.6% of France-2022 ones, were multidrug  
49 resistant. Convergence of toxigenicity with penicillin and erythromycin resistance was  
50 observed in 2 isolates from France-2022. Phylogenetic lineages Gravis and Mitis contrasted  
51 strikingly in their pathogenicity-associated genes.

52

### 53 **Conclusions**

54 This work provides a bioinformatics tool and global population framework to analyze  
55 *C. diphtheriae* genomes, revealing important heterogeneities in virulence and resistance  
56 features. Emerging genotypes combining toxigenicity and first-line antimicrobial resistance  
57 represent novel threats. Genomic epidemiology studies of *C. diphtheriae* should be intensified  
58 globally to improve understanding of reemergence and spatial spread.

59

## Introduction

60 Diphtheria was a leading cause of infant mortality before the implementation of anti-  
61 toxin therapy and mass vaccination programs. Classical diphtheria is a respiratory infection  
62 mainly caused by the *tox* gene-positive strains of the bacterium *Corynebacterium diphtheriae*.  
63 The disease is classically characterized by the presence of a pseudomembranes on the tonsils,  
64 pharynx and larynx. Only some strains of *C. diphtheriae* can produce the diphtheria toxin,  
65 which is encoded by the *tox* gene carried by a prophage integrated into the chromosome of these  
66 strains. The toxigenic strains can induce severe systemic symptoms that include myocarditis  
67 and peripheral neuropathies. Other forms of infection include bacteriemic infections, most often  
68 caused by non-toxigenic strains, and cutaneous infections, which are considered to play an  
69 important role in the transmission of the pathogen.

70 Diphtheria has been virtually eliminated by mass vaccination, but can cause large  
71 outbreaks where vaccination coverage is insufficient (du Plessis *et al.*, 2017; Polonsky *et al.*,  
72 2021; Badell *et al.*, 2021). In France, no case was reported between 1990 and 2001 (Bonmarin  
73 *et al.*, 2009), and in the 2017-2021 period only 6.4 *tox*-positive *C. diphtheriae* were detected  
74 per year by the French surveillance (our unpublished data). In striking contrast, in 2022, 45 *tox*-  
75 positive isolates were detected, including 34 from metropolitan France, mostly associated with  
76 recent arrival from abroad. *C. diphtheriae* also reemerges in several European countries,  
77 strongly associated with non-vaccinated young adults with cutaneous infections with a travel  
78 history from Afghanistan and other countries (Kofler *et al.*, 2022; Badenschier *et al.*, 2022).

79 Whole genome sequencing (WGS) is a powerful approach to understand transmission  
80 and define the pathogenicity-associated characteristics of infectious isolates. *C. diphtheriae* is  
81 a genetically diverse species with multiple phylogenetic sublineages among which a large  
82 heterogeneity of virulence or antimicrobial resistance factors is observed (Sangal and  
83 Hoskisson, 2016; Seth-Smith and Egli, 2019; Hennart *et al.*, 2020; Guglielmini *et al.*, 2021).  
84 One prominent polymorphism in *C. diphtheriae* is the variable presence of the *tox* gene, but the  
85 population dynamics and drivers of *tox* acquisition or loss remain poorly understood. In  
86 addition, non-toxigenic *tox*-bearing (NTTB) *C. diphtheriae* isolates represent 5-20% of *tox*-  
87 positive isolates, but our capacity to predict toxigenicity from genomic sequences is still  
88 limited. Several other experimentally-demonstrated virulence factors have been described in  
89 *C. diphtheriae* (Ott, 2018). Although early 1930s literature suggested a higher virulence of  
90 isolates of biovar Gravis (McLeod, 1943; Barksdale, 1970), it is unknown whether this  
91 historical observation applies to extant diphtheria cases, as recent Gravis isolates are more  
92 rarely *tox*-positive than those of biovar Mitis (Hennart *et al.*, 2020). More generally, the  
93 population variation of virulence factors, and its interactions with clinical outcomes, remain

94 largely to be characterized. Despite being rare, antimicrobial resistance (AMR) in  
95 *C. diphtheriae* is increasingly reported (Mina *et al.*, 2011; Zasada, 2014; Forde *et al.*, 2020;  
96 Hennart *et al.*, 2020), but the mechanisms of resistance that are prevalent across world regions  
97 are not well known, and the evolutionary emergence and dissemination of multi-drug resistant  
98 *C. diphtheriae*, and its possible convergence with toxigenicity in the same strains, should be  
99 carefully monitored.

100 Although WGS of *C. diphtheriae* clinical isolates is increasingly performed for  
101 surveillance purposes, no simple tool currently exists for *C. diphtheriae* genomic feature  
102 extraction and interpretation in clinical, surveillance and research contexts. Besides, analyses  
103 of *C. diphtheriae* genomes remain largely unstandardized, which limits the interpretation of  
104 local genomic epidemiology studies in their global context. Advances towards standardization  
105 include the 7-gene MLST genotyping approach and attached nomenclature of sequence types  
106 (ST) (Bolt *et al.*, 2010), and its core-genome MLST (cgMLST) extension and associated  
107 nomenclature of sublineages and genomic clusters (Guglielmini *et al.*, 2021).

108 Here, we aimed to provide insights into the France 2022 diphtheria emergence by  
109 reporting on its epidemiology and by placing the involved isolates in the global genomic context  
110 of *C. diphtheriae* populations. We introduce DIPHTOSCAN, a genotyping tool designed for rapid  
111 and standardized genomic analyses of Corynebacteria of the *C. diphtheriae* species complex  
112 (CdSC), and illustrate its use by analyzing the 101 *C. diphtheriae* isolates (including *tox*-  
113 negative ones) collected in 2022 in France (henceforth, the France-2022 dataset). We provide  
114 context of this emergence by analyzing 1249 other *C. diphtheriae* genomes of diverse  
115 geographic and temporal origins, including 379 newly sequenced isolates collected by the  
116 French national surveillance laboratory, mostly between 2018 and 2021. We uncovered novel  
117 insights into the global population structure of *C. diphtheriae*, including a striking contrast in  
118 pathogenesis-associated gene clusters between phylogenetic lineages Gravis and Mitis, and  
119 describe high-risk sublineages with convergence of resistance and virulence features.

120

## Results

### 121 **1. The re-emergence of *C. diphtheriae* in France in 2022**

122 In 2022, the French NRC has received 101 human samples of *C. diphtheriae*, from  
123 metropolitan France (n=76) as well as in the Indian Ocean islands of Mayotte (n=10) and La  
124 Reunion (n=6), and in French Guiana (n=9). There were 45 isolates carrying the *tox* gene coding  
125 for diphtheria toxin (*tox*-positive isolates), whereas in the five previous years a total of 32 *tox*-  
126 positive *C. diphtheriae* were detected (**Figure S1A**). *C. diphtheriae* were isolated in  
127 metropolitan France (n=34) and in Mayotte/La Reunion (n=11), while none were found in  
128 French Guiana. The metropolitan France isolates were isolated only in the second part of the  
129 year (**Figure S1B**) and were associated with a recent travel history from Afghanistan (n=24) or  
130 other countries from West Africa, North Africa, Middle East and Southern Asia; These isolates  
131 were predominantly from cutaneous infections, whereas 7 were from respiratory infections  
132 (**Table S1; Figure 1**).

133

### 134 **2. Development of the DIPHTOSCAN pipeline**

135 To provide a tool to extract information from genomes of *C. diphtheriae* and related  
136 potentially toxigenic species, we developed DIPHTOSCAN. The technical characteristics of  
137 DIPHTOSCAN are summarized in **Figure S2-S4** and the methodological details for genotyping  
138 are provided in the Methods section.

139 The DIPHTOSCAN pipeline (**Figure S2**) starts with taxonomic assignment of species.  
140 Recent taxonomic updates have defined, besides the three classical species *C. diphtheriae*,  
141 *C. ulcerans* and *C. pseudotuberculosis*, three novel species of the Corynebacteria of the  
142 *diphtheriae* species complex (CdSC): *C. belfantii* (Dazas *et al.*, 2018), *C. rouxii* (Badell *et al.*,  
143 2020) and *C. silvaticum* (Dangel *et al.*, 2020). If the genome is confirmed to belong to the  
144 CdSC, 7-gene MLST analysis (Bolt *et al.*, 2010) is performed. For *C. diphtheriae*, additional  
145 genotype categorizations can be performed using the BIGSdb-Pasteur database tool: cgST,  
146 genomic cluster and sublineage assignment (Guglielmini *et al.*, 2021). Next, the detection of  
147 antimicrobial resistance determinants (mutations in core genes and horizontally acquired genes)  
148 and virulence factors is performed. DIPHTOSCAN also includes a prediction of the functionality  
149 or disruption of the *tox* gene, the most important virulence factor of CdSC isolates. DIPHTOSCAN  
150 next searches for genomic markers associated with biovars Gravis, Mitis and Belfanti, a  
151 biochemical-based classification that was initiated in the 1930s (Anderson *et al.*, 1931;  
152 McLeod, 1943) and which is still in use for *C. diphtheriae* strain characterization.  
153 IntegronFinder2 (Néron *et al.*, 2022) was included in the pipeline to contextualize resistance  
154 genes. Last, a rapid phylogenetic method based on k-mer distances, JolyTree (Criscuolo, 2020),

155 was integrated to provide quick phylogenetic trees for the genomic assembly datasets under  
156 study. The two latter steps are optional.

157 DIPHTOSCAN was developed using code from Kleborate v2.2.0 (Lam *et al.*, 2021),  
158 AMRfinderPlus (Feldgarden *et al.*, 2021) and BIGSdb (Jolley and Maiden, 2010) with some  
159 modifications (**Figure S3**). A custom code was created for DIPHTOSCAN initiation,  
160 interpretation and for displaying results. The *C. diphtheriae* specific genes (genomic markers,  
161 AMR determinants and virulence factors) for which the genomes are screened by DIPHTOSCAN  
162 (**Figure S4**) are provided in a custom database similar in its structure to the AMRfinderPlus  
163 database ([https://ftp.ncbi.nlm.nih.gov/pathogen/Antimicrobial\\_resistance/](https://ftp.ncbi.nlm.nih.gov/pathogen/Antimicrobial_resistance/)); this database can  
164 be further enriched with novel features in the future. When launching DIPHTOSCAN, the  
165 AMRfinderPlus and custom databases are merged. We used the functions of species  
166 determination, MLST genotyping, and full CDS prediction from Kleborate.

167

### 168 **3. Genetic diversity of *C. diphtheriae* isolates from France, 2022**

169 The *C. diphtheriae* isolates belonging to the France-2022 dataset were sequenced and their  
170 genomic sequences were analyzed using DIPHTOSCAN. Sublineage classification of the isolates  
171 showed that the France-2022 dataset comprised 41 distinct sublineages (defined using the 500  
172 cgMLST mismatch threshold). The nomenclature of these sublineages was established using an  
173 inheritance rule that captures their majority MLST denomination, where possible (Guglielmini  
174 *et al.*, 2021; Hennart *et al.*, 2022), resulting in a strong concordance of sublineage  
175 denominations with the classical MLST identifiers (**Figure 1**). There were 51 different STs, as  
176 9 sublineages comprised two or more closely related STs; in 7 of 9 cases, they only differed by  
177 a single locus. Sublineages thus appeared as useful classifiers for closely related STs.

178 There were four frequently isolated *tox*-positive sublineages: SL824 included 10 isolates  
179 from Mayotte and La Reunion; these all belonged to the same genomic cluster (GC756),  
180 indicating recent transmission. Three other frequent *tox*-positive sublineages were SL377 (n=11  
181 isolates, 10 of which were *tox*-positive), SL698 (n=9) and SL384 (n=7), which were associated  
182 with travel from Afghanistan and countries of the Middle East (**Figure 1**). Whereas SL384 was  
183 genetically homogeneous (GC805), SL377 and SL698 both comprised two genomic clusters  
184 (SL377: GC817 and GC71; SL698: GC795-ST574 and GC804-ST698). SL377-GC71 was not  
185 associated with Afghanistan and one isolate from Senegal was *tox*-negative.

186 Besides the above four frequent sublineages, six additional *tox*-positive sublineages were  
187 isolated: three isolates of sublineage SL486 associated with Senegal and Tunisia; two SL852  
188 isolates associated with Mali; and one SL466 isolate associated with travel from Afghanistan  
189 and one SL464 isolate associated with Thailand. SL91 comprised one non-toxigenic, *tox*-gene  
190 bearing (NTTB) isolate, and SL830 comprised 2 isolates: one *tox*-positive and one *tox*-negative.

191 Besides, there were 31 *tox*-negative sublineages, which were typically isolated once or  
192 twice only; a notable exception was SL297, which comprised six *tox*-negative isolates  
193 associated with travel from Egypt, Senegal, and Mali (**Figure 1**).

194

#### 195 **4. The global phylogenetic framework of *C. diphtheriae***

196 We investigated the global diversity of *C. diphtheriae* to provide context to the France-  
197 2022 emerging genotypes. A dataset of 1,249 comparative *C. diphtheriae* genomes were  
198 sequenced or gathered from previous studies (see Methods). cgMLST grouped these isolates  
199 into 245 sublineages. The 7-gene MLST analysis revealed 364 distinct STs. Almost all (360;  
200 98.6%) STs corresponded one-to-one with the sublineage level, *i.e.*, all isolates of these STs  
201 belonged to the same sublineage. However, 72 sublineages (29.4%) comprised at least two STs.  
202 Of the 123 novel sublineages uncovered here, 114 sublineages were given an identifier inherited  
203 from the 7-gene MLST nomenclature (whereas 9 were attributed an arbitrary number, see  
204 Methods).

205 There were 576 genomic clusters, many of which comprised previously documented  
206 epidemiological clusters of related isolates. For example, GC456 comprised 43 isolates from a  
207 Vancouver inner city outbreak (Chorlton *et al.*, 2019). Whereas 47 GCs had between 5 and 27  
208 isolates (**Table S1; Figure S5A**), the 529 remaining ones had only 1 and 4 isolates. 106 (43.3%)  
209 of the 245 sublineages comprised at least two genomic clusters.

210 To eliminate the population bias introduced by multiple sampling of outbreak strains, we  
211 created a non-redundant subset by randomly selecting one genome per genomic cluster,  
212 isolation year and city (if city was unavailable, the country was used instead) and with the same  
213 resistance genes profile and *tox* status (see column ‘Dataset’ in **Table S1**). These 976  
214 deduplicated genomes (hereafter, the *global dataset*) define the background population of  
215 *C. diphtheriae*.

216 Within the global dataset, 35 sublineages were represented 7 times or more (**Figure 2**). The  
217 two predominant sublineages were SL8 (n=61) and SL5 (n=48); their main 7-gene MLST  
218 sequence types were ST8 and ST5, previously noted to be predominant in the ex-USSR 1990s  
219 outbreak. The most represented *tox*-positive sublineages in the global dataset were SL8, SL453,  
220 SL486, SL377 and SL91, and SL50 was a predominant NTTB sublineage (**Figure 2**).

221 Of the 10 sublineages with *tox*-positive isolates observed in France-2022, 7 were found in  
222 the global dataset; of which 5 were among the 35 frequent global sublineages. Besides, 9 *tox*-  
223 negative sublineages from France-2022 were also frequent in the global dataset (**Figure 2**). Of  
224 the common France-2022 sublineages, SL377, SL384 and SL297 were also common in the  
225 global dataset (**Figure 2**), and their toxigenicity and resistance features matched those observed

226 in the global dataset. In contrast, SL698 (metropolitan France) and SL824 (Indian Ocean) were  
227 uniquely common in the France-2022 dataset (**Figure S5B**).

228 The phylogenetic structure of *C. diphtheriae* revealed a star-like phylogeny with multiple  
229 deeply-branching sublineages as previously reported (Berger *et al.*, 2019; Seth-Smith and Egli,  
230 2019; Hennart *et al.*, 2020; Guglielmini *et al.*, 2021) (**Figure 3**). Sublineages were clustered  
231 according to biovars Gravis (and its *spuA* marker gene) and Mitis as previously noted (Hennart  
232 *et al.*, 2020), and formed two main lineages named Gravis (green branches) and Mitis (purple),  
233 defined by the presence of the *spuA* gene (**Table S1**). cgMLST-defined sublineages were highly  
234 concordant with the phylogeny and often comprised more than one 7-gene ST (**Figure 3**;  
235 **Table S1**). The frequent *tox*-positive sublineages SL377 and SL384 were phylogenetically  
236 related within lineage Gravis (**Figure 3**), suggesting they share ancestrally-acquired genetic  
237 features.

238 We placed within this population background, the France-2022 isolates (**Figure S6**), which  
239 appeared to be dispersed in multiple branches of the global phylogeny. The isolates previously  
240 collected by the French reference laboratory appeared even more diverse and largely dispersed  
241 across the global phylogenetic diversity of *C. diphtheriae* (**Figure S6**), indicating that a large  
242 fraction of the global diversity has been sampled by the French surveillance system.

243 Ribotyping was previously used as a classification and nomenclature system of  
244 *C. diphtheriae* strains (Grimont *et al.*, 2004; Mokrousov, 2009). The 71 ribotype reference  
245 strains sequenced herein or previously (Hennart *et al.*, 2020) were placed in the global  
246 phylogeny (**Figure S7**), showing that these strains are highly diverse. However, this ribotype  
247 subset is biased towards *tox*-positives (40 of 71 strains) and appears to represent unevenly and  
248 incompletely, the currently sampled *C. diphtheriae* diversity.

249

## 250 **5. Population distribution of the diphtheria toxin gene**

251 To evaluate DIPHTOSCAN for its ability to detect the *tox* gene and to predict its toxigenicity,  
252 we used the 855 isolates for which data on *tox* qPCR and Elek test were available. DIPHTOSCAN  
253 detected that *tox* was located at the end of a contig and therefore incomplete in 3 cases (reported  
254 with a '\$' suffix, indicating genomic assembly truncation). Of the 852 remaining isolates, 221  
255 were *tox*-positive and 631 *tox*-negative by the reference qPCR method. DIPHTOSCAN detected  
256 the *tox* gene in 219 (99.1%) of the *tox*-positives, and reported its absence in 2 isolates. Among  
257 the 631 *tox*-negative isolates, DIPHTOSCAN reported the absence of the gene in 625 (99.0)  
258 isolates. Of 198 Elek-positives, 195 (98.5%) were predicted to be toxigenic by DIPHTOSCAN,  
259 whereas 1 was predicted to be non-toxigenic and for two isolates the *tox* gene was not detected.  
260 Of the Elek-negative isolates, 11 (50.0%) were predicted as non-toxigenic by DIPHTOSCAN.  
261 Thus, *tox* detection by DIPHTOSCAN was both sensitive and specific, whereas toxigenicity



262 prediction was highly sensitive but not highly specific, likely due to unexplained non-  
263 toxigenicity in isolates with a full-length toxin gene.

264 In the France 2022 dataset, 45 genomes were detected as *tox*-positive and 44 of these were  
265 predicted as toxigenic, with 100% concordance with the Elek test. In comparison, within the  
266 global dataset, approximately one third of the isolates (331/976; 33.9%) were *tox*-positive, as  
267 defined using DIPHTOSCAN, which detected a truncation and hence predicted non-toxigenicity  
268 in 16.0% of these (52/331).

269 The diversity of *tox*-positive isolates was evident from their distribution in the  
270 *C. diphtheriae* phylogenetic tree, but it was striking that the Gravis branch comprised much less  
271 *tox*-positive sublineages than the Mitis branch (**Figure 3**): in the Gravis lineage, there were only  
272 three main branches of *tox*-positive isolates: (i) an early-branching group of sublineages; (ii) a  
273 branch comprising SL377 and SL384 (two frequent sublineages in France-2022), and (iii) SL8.  
274 NTTB isolates were only observed in the Mitis lineage (with one exception in Gravis-SL384)  
275 and this phenotype was acquired through multiple independent evolutionary events (**Figure 3**).

276 A high diversity of *tox*-negative sublineages was also observed in the global dataset:  
277 whereas 173 of 245 (70.6%) sublineages were entirely *tox*-negative, only 73 (29.8%) of them  
278 had at least 1 *tox*-positive isolate. Of these, 50 sublineages were homogeneous for *tox* status  
279 (*i.e.*, they included uniquely *tox*-positive genomes), whereas 23 sublineages (9.3%) included  
280 both *tox*-positive and *tox*-negative genomes (**Table S1**; **Figure 2**), indicating that the gain or  
281 loss of the *tox* gene is not uncommon within sublineages. When considering the genomic  
282 clusters, almost all were either *tox*-positive or *tox*-negative in the global dataset. Accordingly,  
283 sublineages in the France-2022 dataset were all either *tox* positive or negative, but notably,  
284 SL377-GC71 comprised both types of isolates (**Figure 1**).

285

## 286 6. Antimicrobial resistance

287 DIPHTOSCAN includes a screen of *C. diphtheriae* genomes for the presence of antimicrobial  
288 resistance genes or mutations against 10 classes of antimicrobial agents. DIPHTOSCAN also  
289 computes a resistance score, defined as the number of antimicrobial classes for which at least  
290 one resistance gene or mutation is detected. The resistance score varied from 0 to 8 in the global  
291 dataset; 38.2% non-redundant global isolates had at least one genomic resistance feature, and  
292 118 isolates (12.1%) were multidrug resistant (acquired resistance to  $\geq 3$  drug classes; **Table**  
293 **S1**).

294 Resistance feature frequencies are shown in **Figure 4B** for the global dataset. The highest  
295 frequencies of resistance genes were observed for sulfonamides (exclusively gene *sulI*; rarely  
296 present in two copies; 260 non-redundant isolates; 26.6%) and for tetracycline resistance, where  
297 *tet(O)*, *tet(W)* and *tet(33)* were present in approximately equal proportions (132 isolates; 13.5%

298 in total). The phenicol resistance gene *cmx* was also commonly found. *pbp2m* was present in  
299 34 (3.5%) isolates, and *ermX* [sometimes named *erm(X)*] in 36 (3.7%) isolates, with 14 (1.4%)  
300 isolates carrying both *pbp2m* and *ermX*.

301 Antimicrobial resistance genes were dispersed across the global *C. diphtheriae*  
302 phylogenetic tree (**Figure 3**). The distribution of resistance at the sublineage level showed that  
303 just above half of the sublineages (128; 52.0%) comprised at least one strain with at least one  
304 resistance genomic feature (**Table S1**). The two sublineages with the most resistant strains were  
305 SL8 (the main sublineage involved in the ex-USSR outbreak; 46 strains) and SL377 (17 strains)  
306 (**Figure 2**). 19 sublineages carried at least one multidrug resistant isolate, and SL377 and SL405  
307 were the most frequent of these (**Figure 2**).

308 Against this background, the France-2022 isolates appeared to carry resistance features  
309 much more frequently, including *pbp2m*, *ermX* and quinolone-resistance determining mutations  
310 (**Figures 1 and 4**). 61 (60.4%) isolates presented at least one resistance feature (**Table S1**;  
311 **Figure 1**), and 44 (43.6%) were multidrug resistant.

312 First-line treatments of diphtheria are penicillin or amoxicillin and macrolides in case  
313 of allergy to beta-lactams. The *pbp2m* gene confers decreased susceptibility to penicillin and  
314 other beta-lactams (Forde *et al.*, 2020; Hennart *et al.*, 2020), whereas *ermX* (and rarely *ermC*)  
315 are associated with erythromycin resistance in *C. diphtheriae* (Tauch *et al.*, 1995, 2003). In the  
316 global dataset, 34 isolates (**Table S1**; including strain BQ11 with three copies consistent with  
317 Forde *et al.* 2020) carried *pbp2m* and 35 carried *ermX*; 14 (1.4%) isolates carried both genes.  
318 Sublineages SL297 and SL484 were the most common carriers of these genes, whereas the  
319 frequent multidrug resistant sublineages SL377, SL384 and SL301 did not carry *ermX* and  
320 *pbp2m* (**Figure S8**). In France-2022, 8 (7.9%) isolates carried both *pbp2m* and *ermX*. These  
321 were observed in patients with travel history from Mali (SL395, SL542, SL852) and Egypt  
322 (SL297-GC820).

323 Antimicrobial susceptibility phenotypes were determined for the France-2022 dataset,  
324 and were highly concordant with the presence of resistance features (**Table S4**). Resistance to  
325 penicillin and macrolides was associated with *pbp2m* and *ermX*, respectively, although some  
326 *ermX*-carrying isolates remained susceptible to erythromycin (**Table S4**).

327 We included in DIPHTOSCAN a search for integrons, which may harbor multiple resistance  
328 genes in *C. diphtheriae* (Barraud *et al.*, 2011; Arcari *et al.*, 2023). In the global dataset, we  
329 identified 45 (4.6%) isolates carrying integrons (including integrase-less ones, *i.e.*, CALINs)  
330 (**Table S1**), which were highly dispersed in the phylogeny (not shown). In France-2022, we  
331 found the presence of complete integrons in 9 isolates and integrase-less integrons in 9  
332 additional isolates (18; 17.8%). These structures were strongly associated with antimicrobial  
333 resistance, particularly to trimethoprim and sulfonamides (**Figure 1**; **Table S1**).

334

## 335 **7. Dual risk isolates: convergence of diphtheria toxin and multidrug resistance, including** 336 **to first-line treatments**

337 The presence within the same isolates of multidrug resistance and toxigenicity could  
338 cause particularly threatening infections. We therefore explored the co-occurrence of these two  
339 genotypes (**Figure 2**). In the global dataset, 57 (5.8%) isolates were both multidrug resistant  
340 and *tox*-positive. The majority of these isolates belonged to a few sublineages (**Figure 2**),  
341 including SL377, which comprised 9 *tox*-positive multidrug resistant isolates mostly from India  
342 (and also observed in France-2022). Eight convergent isolates of SL301 were also observed  
343 from India, Austria and Syria. SL453 had three *tox*-positive multidrug resistant isolates, which  
344 were isolated in Spain and France with links to Afghanistan (Arcari *et al.*, 2023). In  
345 metropolitan France, there were 22 *tox*-positive isolates that were multidrug resistant (21.8%),  
346 with SL377 and SL696 being predominant among these (**Table S1, Figure 1**).

347 Regarding resistance genes to first-line treatments, there was not a single isolate  
348 carrying at the same time *tox*, *pbp2m* and *ermX* in the global dataset (**Table S1**). However, in  
349 France-2022, SL852 isolates (from two patients with travel history from Mali) were *tox*-positive  
350 and carried *pbp2m* and *ermX*. Furthermore, they carried other resistance genes including *cmx*,  
351 *sull*, *dfrA1*, and in addition *tet33* and *aadA15* for isolate FRC1688. This latter isolate only  
352 lacked resistance features to quinolones and rifampicin. No other isolate of this particularly  
353 concerning sublineage (SL852) was found in the global dataset.

354

## 355 **8. Lineages Gravis and Mitis differ in the presence of pathogenicity-associated genes**

356 Biovars represent an early attempt to discriminate among *C. diphtheriae* strains (Anderson  
357 *et al.*, 1931) and are still commonly reported. We found that lineages Mitis and Gravis, defined  
358 genetically based on the presence of the *spuA* gene probably involved in starch utilization,  
359 correspond to two distinct parts of the phylogenetic tree (**Figure 3**) as previously reported  
360 (Hennart *et al.*, 2020; Guglielmini *et al.*, 2021). Note that the match between lineage and *spuA*  
361 or biovar phenotype is not absolute, as a few isolates within the Gravis branch were *spuA*-  
362 negative (in particular SL625, SL130, SL102, and SL377) and 42 (5.1%) isolates of the Mitis  
363 lineage were *spuA*-positive. Among the France-2022 isolates, for which biovars were in  
364 addition determined phenotypically, the two biovars were also phylogenetically distinct  
365 (**Figure 1**). Nearly four in five (n=78) of the France-2022 isolates had a Mitis biotype (including  
366 37 *tox*-positives), with 23 Gravis strains (8 *tox*-positive).

367 To provide a population-level view of pathogenesis features in *C. diphtheriae*, we included  
368 in the DIPHTOSCAN database of searched genes, in addition to the *tox* gene, all virulence genes  
369 previously demonstrated or strongly suspected to be involved in diphtheria pathogenesis (see

370 **Table S2** for pathogenesis involvement evidence). These include genes involved in iron and  
371 heme acquisition, fimbriae biosynthesis and assembly, and other adhesins (Ott *et al.*, 2022).

372 Screening for these genes in the global dataset revealed highly heterogeneous patterns of  
373 presence and phylogenetic distribution (**Table S1; Figure S9**). We found that a number of  
374 virulence factors are highly conserved within *C. diphtheriae*; for example, DIP1546 was  
375 present in all genomes except in 28\_DSM43988, and DIP0733, DIP1281, DIP1621, and  
376 DIP1880 were fully conserved (**Table S1**). The corynebactin transport (*ciuA-D*) gene cluster  
377 was present in all genomes, with one exception, whereas the corynebactin synthesis (*ciuEFG*)  
378 locus was absent or incomplete in only 5.4% of genomes (n=29 Mitis, n=25 Gravis); of these,  
379 33 lacked the *ciuE* gene, which is essential for siderophore synthesis. One of the genomes  
380 lacking *ciuE* corresponds to the vaccine strain PW8, which is defective for corynebactin  
381 synthesis (Russell and Holmes, 1985). The heme-acquisition genes *hmuTUV* were also largely  
382 conserved (921 genomes; 94.4%).

383 In contrast, some genes were infrequent: DIP2014, a gene encoding for a BigA-like adhesin,  
384 was detected in only a few sublineages of the Gravis branch (133 isolates), and the DIP0543  
385 (also known as *nanH*, coding for a sialidase) was present in only a few sublineages distributed  
386 across the phylogeny (not shown).

387 Remarkably, we uncovered a sharp divide between lineages Gravis and Mitis in terms of  
388 iron metabolism-associated genes, fimbriae gene clusters and other genes (**Figure S9**). The  
389 putative siderophore synthesis and transport operon *irp2ABCDEFGHI-irp2JKLMN* was strongly  
390 associated with the Mitis lineage: 513 out of 567 (90.5%) Mitis isolates were *irp2*-positive,  
391 whereas only 1 of 406 Gravis isolates was *irp2*-positive. The iron transport cluster *irp1ABCD*  
392 was also mainly present in the Mitis lineage. Differently, the *htaA* gene, which is part of the  
393 same gene cluster as *hmuTUV* and codes for a membrane protein that binds hemoglobin, was  
394 absent or truncated in most genomes from the Mitis branch (92.1%), whereas it was largely  
395 conserved in the Gravis branch (99.8% *htaA*-positive). Similar to *htaA*, genes *chtA* and *chtB*,  
396 which have sequence and functional similarity to *htaA* and *htaB*, were also strongly associated  
397 with the Gravis lineage: 304 of 406 Gravis isolates were *chtAB*-positive (74.9%), whereas only  
398 7 of 567 Mitis isolates were *chtAB*-positive (1.2%). In sharp contrast, the *htaC* gene, which is  
399 suspected to be involved in hemin transport, and which is also in genetic linkage with the  
400 *hmuTUV* gene cluster, was entirely absent from the Gravis branch, but was detected in 68.6%  
401 of Mitis genomes.

402 Three main fimbriae gene clusters, encoding fimbrial proteins, SpaA, SpaD and SpaH, have  
403 been described in *C. diphtheriae* (Rogers *et al.*, 2011; Reardon-Robinson and Ton-That, 2014;  
404 Sangal and Hoskisson, 2016). We found that these were more commonly found in the Gravis  
405 branch compared to the Mitis branch (**Figure S9**). The SpaH gene cluster (*spaGHI-srtDE*) was

406 present in its entirety in 254 genomes and as a cluster with one missing gene in 29 isolates, all  
407 of which belonged to the Gravis lineage. The other two systems showed some variability in the  
408 distribution of their genes. The sortase-mediated assembly genes of the SpaA type pili, *spaABC*,  
409 were found in biovar Gravis in similar proportions (87.2% *spaA*, 86.2% *spaB* and 86.0% *spaC*-  
410 positive), whereas in Mitis *spaB* was present in about half of the genomes (49.0%) and *spaA*  
411 and *spaC* in one third (17.5%, and 18.2%, respectively). The distribution of the SpaA pilin-  
412 specific sortase gene *srtA* was similar to that of *spaB* (98.8% in Gravis, 49.9% in Mitis), and  
413 the complete SpaA gene cluster *spaABC-srtA* was found in only 299 genomes (30.6%), the  
414 majority of which were of Gravis lineage (n=256). Last, genes of the SpaD cluster were less  
415 frequent (*spaD* 8.7%, *spaE* 14.9%, *spaF* 9.3%, *srtB* 33.2%, *srtC* 33.7%) compared to the other  
416 pili types, and the complete gene cluster (*spaDEF-srtBC*) was found only in 11 genomes, all of  
417 which belonged to lineage Gravis. Interestingly, the presence of SpaD and SpaH complemented  
418 each other in the Gravis branch (**Figure S9**).

419 We further found that the collagen-binding protein DIP2093 (Peixoto et al.,2017) is  
420 strongly associated with the Gravis lineage: 118 of 406 (29.1%) Gravis isolates were DIP2093-  
421 positive, whereas only 3 of 567 (0.5%) Mitis isolates were.

422 The complement of virulence genes of the France-2022 isolates was in full agreement  
423 with their Gravis/Mitis placement and the above observations. For example, the *irp2A-I* and  
424 *irp2J-N* gene clusters were present uniquely in sublineages belonging to the Mitis branch, and  
425 the *htaC* gene was present only in 64.2% of the Mitis genomes (**Table S1**); *chtA* and *chtB* were  
426 completely absent in Mitis and the collagen-binding protein DIP2093 uniquely in Gravis  
427 isolates (n=16, 47.1%). None of the France-2022 isolates carried a complete SpaD fimbriae  
428 cluster; in particular, they all lacked at least the *spaD* gene; and only 8 Gravis genomes carried  
429 the complete SpaH cluster. The latter were dispersed among various lineages (SL32, SL374,  
430 SL502, SL542, SL130).

431

## Discussion

432 In recent years, large epidemics of diphtheria have been observed, *e.g.*, in South Africa,  
433 Bangladesh and Yemen (du Plessis *et al.*, 2017; Polonsky *et al.*, 2021; Badell *et al.*, 2021),  
434 while a progressive increase of diphtheria cases has been noted in multiple countries (Bernard  
435 *et al.*, 2019; Truelove *et al.*, 2020). However, so far, our understanding of diphtheria  
436 reemergence has been hindered by a lack of background knowledge on the population diversity  
437 of *C. diphtheriae*, its sublineages of concern and the epidemiology of their local or global  
438 dissemination. Here, we report on a sharp increase in *tox*-positive *C. diphtheriae* in France in  
439 2022, and developed a bioinformatics pipeline, DIPHTOSCAN, which enables to harmonize the  
440 way genomic diversity and genetic features of medical concern are detected, reported and  
441 interpreted. We illustrate how this novel tool provides clinically-relevant genomic profiling and  
442 evolutionary understanding of emergence, by placing the 2022 *C. diphtheriae* from France in  
443 the context of 1,249 global *C. diphtheriae* genomes.

444 Our results provide an updated overview of the population diversity of *C. diphtheriae*  
445 based on currently available genomic sequences. As previously reported (Berger *et al.*, 2019;  
446 Seth-Smith and Egli, 2019; Hennart *et al.*, 2020; Guglielmini *et al.*, 2021), *C. diphtheriae* is  
447 made up of multiple sublineages that are related through a star-like phylogeny. We here  
448 uncovered 123 novel sublineages, for a total of 253 described ones. We observed that, compared  
449 to previous datasets, there was no sublineage fusion upon adding novel genomes, which  
450 indicated an excellent stability of *C. diphtheriae* sublineage classification. The latter provides  
451 a broad classification of isolates that correlates strongly with classical MLST, and which  
452 facilitates a deep-level approach to *C. diphtheriae* diversity and evolution. The naming of  
453 sublineages by inheritance of ST numbers will facilitate continuity with classical MLST.  
454 Besides, sublineage classification is more congruent with phylogenetic relationships: whereas  
455 most (140/146; 95.8%) non-singleton sublineages were monophyletic, only 134 of 167 (79.8%)  
456 non-singleton STs were (data not shown). We therefore strongly recommend transitioning from  
457 MLST to the cgMLST-based nomenclature, which is available on the BIGSdb-Pasteur  
458 platform. Our phylogenetic analysis of reference strains of the historical ribotype nomenclature  
459 provides a first overview of their relationships, to our knowledge, and allows revisiting  
460 genealogical inferences that were made among ribotypes based on CRISPR spacer variation  
461 (Mokrousov, 2009).

462 Genomic clusters represent a much narrower genetic classification of *C. diphtheriae*  
463 isolates, compatible with recent transmission (Guglielmini *et al.*, 2021). Therefore, genomic  
464 clusters appear more relevant than sublineages for epidemiological investigation purposes, as  
465 illustrated for example within SL377: whereas GC817 was associated with Afghanistan, GC71

466 was associated with Senegal and these two genomic clusters of sublineage SL377 were clearly  
467 distinct phylogenetically (**Figure 1**).

468 The diagnostic and surveillance of diphtheria is largely based on the detection of the *tox*  
469 gene and its expression (WHO, 2018). We found that the determination of the *tox* gene presence  
470 by DIPHTOSCAN was highly concordant with the experimental reference qPCR. We also found  
471 that DIPHTOSCAN can predict a large proportion of non-toxigenic *tox* gene-bearing (NTTB)  
472 isolates. Still, some NTTB isolates were not identified by DIPHTOSCAN. These cases may be  
473 attributable to (i) a lack of detection by the Elek test due to a low level of expression of the  
474 toxin gene in some strains, or (ii) yet unknown genetic mechanisms that abort *tox* gene  
475 expression entirely (unexplained true NTTB). Future work is needed to define the genotype-  
476 phenotype links underlying toxigenicity and to improve our predictive capacity of toxigenicity  
477 from genomic sequences. In the non-redundant global dataset, 16.0% of *tox*-positive isolates  
478 were predicted as NTTB, which provides a quantitative view of the relevance of differentiating  
479 mere *tox* gene presence from actual toxigenicity. The capacity to predict toxigenicity from  
480 sequences opens interesting perspectives as to the diagnostic of diphtheria based on rapid  
481 genomic sequencing. Our phylogenetic analysis showed that gain or loss of the *tox* gene is a  
482 rare event at the timescale of genomic cluster diversification. The phenomenon of *tox* status  
483 switch by phage acquisition or loss during infection or transmission was suspected  
484 previously (Pappenheimer and Murphy, 1983) and deserves further study given its importance  
485 for public health and clinical management.

486 Up until now, antimicrobial resistance has been considered of moderate clinical  
487 concern in *C. diphtheriae* (WHO, 2018; Zasada, 2014). Although resistant strains have been  
488 described, clinical susceptibility breakpoints have lacked standardization and the prevalence,  
489 origin and dissemination of resistance genetic features are largely unknown. Here, we identified  
490 in the France-2022 isolates as well as in the global *C. diphtheriae*, multidrug resistant isolates  
491 and/or isolates resistant to first-line treatments. We provide an overview of the prevalence and  
492 distribution of resistance genes or mutations in *C. diphtheriae*, and identify sublineages that  
493 carry multiple resistance genes. Because antimicrobial resistance phenotypes are typically  
494 unattached to publicly available genomic sequences, it is not possible to link these genomic  
495 features complements to resistance phenotypes systematically. However, this (**Table S4**) and  
496 previous works clearly showed that most resistance genetic features identified here may impact  
497 resistance phenotypes (Tauch *et al.*, 1995, 2003; Hennart *et al.*, 2020; Forde *et al.*, 2020). Of  
498 particular concern, *tox*-positive isolates that are resistant to multiple drugs and/or first-line  
499 treatments were identified herein, with the convergence of *tox*, *pbp2m* and *ermX* in two 2022  
500 cases with a travel history from Mali, which were resistant to 9 and 11 out of 23 tested

501 antimicrobials, respectively. Such isolates may pose serious clinical management difficulties,  
502 and multidrug resistant *C. diphtheriae* should therefore be closely monitored.

503 The combined analysis of the France-2022 and global datasets using a unique pipeline  
504 provides context to the reemergence of diphtheria (**Figure S6**). Here, we found that some  
505 sublineages contributing to the reemergence were previously observed, whereas others are  
506 described for the first time. For example, SL377, one of the major toxigenic and resistant  
507 sublineages observed in France-2022, had been circulating in India during 2016 and was  
508 reported in Europe (Spain and France) since 2015 (**Table S1**). In contrast, SL698 was absent  
509 from the global dataset. Of the 10 *tox*-positive France-2022 sublineages, five were associated  
510 with travel from Afghanistan, and were recently described in other European countries too  
511 (Kofler *et al.*, 2022; Badenschier *et al.*, 2022).

512 The DIPHTOSCAN tool will facilitate the harmonized characterization of *C. diphtheriae*  
513 sublineages of concern. Several virulence-associated genes were largely conserved in the entire  
514 *C. diphtheriae* population analyzed; these genomic features may therefore be central for  
515 *C. diphtheriae* colonization and transmission among humans, as there appears to be a strong  
516 selective pressure to maintain them. The distribution of other, more variably present, virulence-  
517 associated genes uncovers a very striking dichotomy between the Gravis and Mitis lineages, as  
518 heme and iron-acquisition systems and Spa-encoded fimbriae gene clusters were either  
519 associated with the Mitis or the Gravis lineages, in a largely mutually exclusive way. Based on  
520 these observations, the Gravis lineage may preferentially capture iron from hemin, whereas the  
521 Mitis one could be associated with the ability to synthesize and use siderophores. There might  
522 be important implications for the regulation and expression level of the *tox* gene, which is  
523 controlled by the iron-dependent DtxR repressor. Importantly, the toxin gene and its NTTB-  
524 leading disruptions were also unequally distributed between Gravis and Mitis lineages. It was  
525 noted early that toxin production is less inhibited by infection-relevant iron concentrations in  
526 Gravis strains (Mueller, 1941; McLeod, 1943), and our results shed a new light and provides  
527 experimentally testable hypotheses on this critical difference in the biology of infection of the  
528 Gravis and Mitis lineages.

529 Another striking feature we uncovered is the distribution of gene clusters coding for  
530 fimbriae. Previous work reported SpaA as being largely conserved in *C. diphtheriae*, with SpaD  
531 and SpaH being more variably present (Reardon-Robinson and Ton-That, 2014; Ott, 2018;  
532 Sangal and Hoskisson, 2016). We found that SpaA was largely present in our dataset, however,  
533 the complete gene cluster *spaABC-srtA* was mostly found in the Gravis branch. SpaD was also  
534 more common among Gravis genomes, although the complete cluster (*spaDEF-srtBC*) was  
535 only detected in a minority of genomes. None of the Mitis isolates were positive for SpaH.  
536 These three Spa systems were experimentally shown to be involved in adhesion to different



537 human tissues: pharyngeal (SpaA), laryngeal (SpaD) and pulmonary (SpaH) epithelial cells  
538 (Mandlik *et al.*, 2007; Reardon-Robinson and Ton-That, 2014). The Gravis/Mitis dichotomy in  
539 Spa-type fimbriae may have important implications regarding a possible differential ecology,  
540 transmission, tissue tropism and pathophysiology of these two major *C. diphtheriae* lineages.

541 In conclusion, we developed and applied to a large dataset, the bioinformatics tool  
542 DIPHTOSCAN. Its public availability and ease of use will enable to conveniently extract and  
543 interpret genomic features that are relevant to the clinical and public health management of  
544 diphtheria cases, and to future research on the genotype-clinical phenotype links in  
545 *C. diphtheriae*. This dedicated tool is also applicable to the other members of the *C. diphtheriae*  
546 complex, such as *C. ulcerans* (data not shown). Harmonization of genomic studies in this group  
547 of pathogens, which have been largely forgotten but currently undergo re-emergence in Europe  
548 and elsewhere, will support genomic surveillance of diphtheria, will contribute to enhance our  
549 understanding of the pathogenesis of modern diphtheria, and opens interesting hypotheses as to  
550 the underlying mechanisms of variation in clinical severity and forms of diphtheria.

551

## Material & Methods

### 552 **Clinical isolates inclusion and global genomic sequence dataset**

553 To investigate the epidemiology of diphtheria in France, we included all cases of  
554 *C. diphtheriae* infections detected by the French surveillance in 2022. Among 144 isolates  
555 received by the National Reference Center, there were 101 deduplicated isolates when retaining  
556 only one from each patient. These were isolated in metropolitan France as well as in Mayotte,  
557 La Reunion and French Guiana (**France-2022 dataset, Table S1**). Note that metropolitan  
558 France comprises mainland France and Corsica, as well as nearby islands in the Atlantic Ocean,  
559 the English Channel (French: la Manche), and the Mediterranean Sea. All isolates collected in  
560 2022 from metropolitan France were from mainland France. Overseas France is the collective  
561 name for all the French territories outside Europe.

562 In addition, a total of 1,249 comparative genomes were included (**Table S1**). First, we  
563 sequenced for the present study 373 additional isolates, including 320 collected prospectively  
564 between 2008 and 2021 by the French National Reference Center (NRC), 34 historical clinical  
565 isolates mostly from metropolitan France and 19 isolates from Algeria (Benamrouche *et al.*,  
566 2016). These new genomes were sequenced to complement the 226 previous genomes from  
567 *C. diphtheriae* from the French diphtheria surveillance system (Hennart *et al.*, 2020;  
568 Guglielmini *et al.*, 2021), including 43 isolates from Yemen (Badell *et al.*, 2021). Together,  
569 these represent 599 produced by the NRC for Corynebacteria of the *diphtheriae* complex (**non-**  
570 **2022 French NRC dataset, Table S1**). Nearly four-fifths (532; 88.7%) of these isolates were  
571 prospectively collected between 2008 and 2021 from French metropolitan and overseas  
572 territories, 54 isolates (9.0%) were collected between 1990 and 2007 from France and Algeria  
573 and 14 (2.3%) isolates collected between 1951 and 1987 from metropolitan France.

574 Second, we included publicly-available genomes from NCBI, mostly previously  
575 published and isolated in South Africa (du Plessis *et al.*, 2017), Germany-Switzerland (Meinel  
576 *et al.*, 2016), Germany (Dangel *et al.*, 2018; Berger *et al.*, 2019), Canada (Chorlton *et al.*, 2019)  
577 Austria (Schaeffer *et al.*, 2020), the USA (Xiaoli *et al.*, 2020; Williams *et al.*, 2020), Spain  
578 (Hoefler *et al.*, 2020), India (Will *et al.*, 2021) and Australia (Timms *et al.*, 2018). Altogether,  
579 this represents a dataset of 579 genomes (**non-French public dataset, Table S1**).

580 Further, we sequenced 6 ribotype reference strains (Grimont *et al.*, 2004). Together with  
581 65 previously sequenced (Hennart *et al.*, 2020), this represents a dataset of 71 genomes of  
582 ribotype reference strains (**Table S1**).

583 From the global set of 1,249 genomes (**non-2022 French NRC + non-French public**  
584 **dataset + ribotype datasets**), we created a non-redundant subset of genomes by randomly  
585 selecting one genome per genomic cluster (threshold: 25 cgMLST mismatches; see below),

586 isolation year and city (if city was unavailable, the country was used instead); this deduplicated  
587 subset comprised 976 genomes (hereafter, the *global dataset*).

588

### 589 **Microbiological characterization of isolates at the French National Reference Laboratory**

590 *C. diphtheriae* isolates were grown and purified on Tinsdale agar. Strains were  
591 characterized biochemically for pyrazinamidase, urease, and nitrate reductase and for  
592 utilization of maltose and trehalose using API Coryne strips (BioMérieux, Marcy l’Etoile,  
593 France) and the Rosco Diagnostica reagents (Eurobio, Les Ulis, France). The Hiss serum water  
594 test was used for glycogen fermentation. The biovar of isolates was determined based on the  
595 combination of nitrate reductase (positive in Mitis and Gravis, negative in Belfanti) and  
596 glycogen fermentation (positive in Gravis only). Antimicrobial susceptibility was determined  
597 by disc diffusion (BioRad, Marnes-la-Coquette, France). Zone diameter interpretation  
598 breakpoints are given in **Table S3**.

599 The presence of the diphtheria toxin *tox* gene was determined by real-time PCR assay  
600 (Badell *et al.*, 2019), whereas the production of the toxin was assessed using the modified Elek  
601 test (Engler *et al.*, 1997).

602 For genomic sequencing, isolates were retrieved from  $-80^{\circ}\text{C}$  storage and plated on  
603 tryptose-casein soy agar for 24 to 48 h. A small amount of bacterial colony biomass was  
604 resuspended in a lysis solution (20 mM Tris-HCl [pH 8], 2 mM EDTA, 1.2% Triton X-100, and  
605 lysozyme [20 mg/ml]) and incubated at  $37^{\circ}\text{C}$  for 1 h DNA was extracted with the DNeasy  
606 Blood&Tissue kit (Qiagen, Courtaboeuf, France) according to the manufacturer’s instructions.  
607 Genomic sequencing was performed using a NextSeq500 instrument (Illumina, San Diego, CA)  
608 with a  $2 \times 150$ -nucleotide (nt) paired-end protocol following Nextera XT library preparation  
609 (Hennart *et al.*, 2020).

610 For de novo assembly, paired-end reads were clipped and trimmed using AlienTrimmer v0.4.0  
611 (Criscuolo and Brisse, 2013), corrected using Musket v1.1 (Liu *et al.*, 2013), and merged (if  
612 needed) using FLASH v1.2.11 (Magoč and Salzberg, 2011). For each sample, the remaining  
613 processed reads were assembled and scaffolded using SPAdes v3.12.0 (Bankevich *et al.*, 2012).

614

### 615 **Merging of the Oxford and Pasteur MLST databases**

616 Two *C. diphtheriae* databases using the BIGSdb framework were originally designed  
617 separately for distinct purposes: while Oxford’s PubMLST database mainly offered 7-gene  
618 MLST (Bolt *et al.*, 2010), the Pasteur database was used for the *Corynebacterium* cgMLST  
619 typing (Guglielmini *et al.*, 2021). To facilitate the use of these resources and avoid redundancy  
620 in the curation of the two independent genomic libraries, a merging of the databases was  
621 decided in agreement with PubMLST administrators. In order to merge the data available in the

622 two databases, we proceeded as per BIGSdb dual design: isolates genomes and provenance data  
623 were imported into the “isolates” database, whereas allelic definitions of MLST were imported  
624 into the “seqdef” database.

625 Regarding the isolates database, we first downloaded Oxford’s PubMLST  
626 *C. diphtheriae* database. To avoid isolate entries duplication, we identified common isolates  
627 between the two databases, and filtered duplicate isolates before import into the Pasteur  
628 database. In total, 684 out of 934 (73%) isolates from the Oxford database were imported. To  
629 facilitate the tracing of isolates and their possible previous existence in Oxford’s database,  
630 isolates identification numbers (BIGSdb-Pasteur ID number) of isolates from the Oxford  
631 database were numbered from 1,520 to 2,003. We also collated them into a public project  
632 collection called “Oxford” (project ID 13).

633 Regarding the sequence and profiles definition database, we imported MLST alleles and  
634 profiles into an initially void MLST scheme container within the BIGSdb-Pasteur database.  
635 MLST analysis was performed on all isolates of the BIGSdb-Pasteur database, including the  
636 ones imported from Oxford, which were therefore assigned the same MLST genotype as  
637 previously in the Oxford database.

638 At the end of the merging process, all isolates and MLST data from PubMLST’s  
639 *C. diphtheriae* database were available into the BIGSdb-Pasteur *C. diphtheriae* species  
640 complex database (<https://bigsdB.pasteur.fr/diphtheria/>), and Oxford’s PubMLST  
641 *C. diphtheriae* database was shut down. As of September 22<sup>nd</sup>, 2022, the database resulting  
642 from the merged datasets comprised 1,478 public isolates records with 794 associated genomes,  
643 and 2,392 isolates in total when considering private entries. The number of entries varied across  
644 species: *C. diphtheriae* (n = 1,291; 87.4%) and *C. ulcerans* (n = 131; 8.9%), *C. belfantii* (n =  
645 45; 3.0%) and *C. rouxii* (n = 10; 0.7%). The MLST scheme comprised 854 registered STs.

646

#### 647 **cgMLST and nomenclature of sublineages**

648 The MLST and cgMLST genotypes (cgST) were defined using the Institut Pasteur  
649 *C. diphtheriae* species complex database at <https://bigsdB.pasteur.fr/diphtheria/>.

650 A core genome MLST (cgMLST) scheme comprising 1,305 loci (Guglielmini *et al.*,  
651 2021) was employed to define the alleles and cgST of the 1,249 genomic sequences using  
652 BIGSdb (<https://bigsdB.pasteur.fr/diphtheria/>). Using the 1,249-genomes dataset, the mean  
653 number of missing alleles per profile was 12 (0.9%) and almost all (n=1,242; 99.4%) genomes  
654 had a cgMLST profile with fewer than 65 (5%) missing alleles. A cgST number was defined  
655 for all but one cgMLST profiles (one genome had 219 missing alleles, whereas the admissible  
656 threshold is 10%, i.e., 130 missing alleles).

657 Genomes were classified using the single-linkage cluster-profile.pl function of BIGSdb  
658 into genomic clusters (25 mismatch threshold) and sublineages (500 mismatches). Sublineages  
659 were attributed numbers by using an ST inheritance rule (Hennart *et al.*, 2022), which was  
660 applied from SL1 to SL744, after which the numbers are attributed consecutively with no  
661 reference to MLST identifiers, starting at 10,000 (see column ‘SL’ in **Table S1**).

662

### 663 **Phylogenetic analysis based on a core genome**

664 Panaroo v1.2.3 was used to generate from the assembled genomic sequences, a core  
665 genome used to construct a multiple sequence alignment (cg-MSA). The genome sequences  
666 were first annotated using prokka v1.14.5 with default parameters, resulting in GFF files.  
667 Protein-coding gene clusters were defined with a threshold of 70% amino acid identity, and  
668 core genes were concatenated into a cg-MSA when present in 95% of genomes. IQtree version  
669 2 was used to build a phylogenetic tree based on the cg-MSA, with the best fitting model  
670 TVM+F+R5. The tree was constructed from 1,948 core genome loci, for a total alignment  
671 length of 1,986,172 bp (79.8% of NCTC13129 genome length, of 2,488,635 bp), was rooted  
672 using *C. belfantii* strain FRC0043<sup>T</sup>, and is available at:  
673 <https://itol.embl.de/tree/1579917435471751662784292>.

674

### 675 **Development of the DIPHTOSCAN pipeline**

676 To develop DIPHTOSCAN, we combined code from Kleborate (Lam *et al.*, 2021), NCBI  
677 database of AMR genes (<https://www.ncbi.nlm.nih.gov/pathogens/refgene/#>), and  
678 AMRfinderPlus (Feldgarden *et al.*, 2021). The structures of DIPHTOSCAN and its custom  
679 database are presented in **Figure S3** and **Figure S4**. The functionalities are presented in  
680 **Figure S2**. To facilitate readability and downstream analyses, the output of DIPHTOSCAN is  
681 generated in a tab-delimited format. The execution time of DIPHTOSCAN increases linearly with  
682 the number of input genomes. Roughly, 40 seconds are needed to scan a single genome with 1  
683 cpu. DIPHTOSCAN computations can be parallelized, as AMRFinderPlus and JolyTree use  
684 parallelization.

685

### 686 **Assignment of species, MLST and Sequence Types (ST)**

687 To perform rapid and accurate species identification, DIPHTOSCAN uses the k-mer-derived  
688 Mash distances (Ondov *et al.*, 2016). DIPHTOSCAN calculates Mash distances (Mash v2.2)  
689 between the query genomes and a collection of reference assemblies of the *CdSC*, and reports  
690 the species with the smallest distance. *C. diphtheriae* genomes were confirmed as  
691 *C. diphtheriae* based on a Mash distance smaller than 0.05 with either the *C. diphtheriae* type

692 strain NCTC11397<sup>T</sup> (= C7S), the reference genome strain NCTC13129, or the vaccine strain  
693 PW8 (Park-Williams 8).

694 Mash distance  $\leq 0.05$  is reported as a strong match,  $\leq 0.1$  as weak. We have used and adapted  
695 the structure of the Kleborate tool for this function. This approach was validated by comparing  
696 DIPHTOSCAN species assignments with those obtained by average nucleotide identity (ANI;  
697 Konstantinidis and Tiedje, 2005) using FastANI (Jain *et al.*, 2018) using the global dataset;  
698 100% concordance was achieved.

699 MLST profiles and sequence types (ST) were defined using the international MLST  
700 scheme for *C. diphtheriae* and *C. ulcerans*. DIPHTOSCAN defines these genotypes for genomic  
701 sequences using the analogous script from Kleborate. In order to use an up-to-date version of  
702 the MLST nomenclature, which is regularly updated, the MLST profiles and alleles are  
703 downloaded at the start of the pipeline before genotyping the genomes. The  
704 download\_alleles.py script from BIGSdb is used for this purpose  
705 ([https://github.com/kjolley/BIGSdb/tree/develop/scripts/rest\\_examples](https://github.com/kjolley/BIGSdb/tree/develop/scripts/rest_examples)).

706

#### 707 **Biovar-associated markers detection**

708 The three main biovars of *C. diphtheriae* can be distinguished based on isolate abilities to  
709 reduce nitrate and to metabolize glycogen. Previously, a strong concordance was found between  
710 the biovar and the presence in the genome of several genomic markers including *spuA*, which  
711 codes for a putative alpha-1,6-glycosidase, and the *narKGHJI* operon for nitrate reductase  
712 (Sangal *et al.*, 2014; Santos *et al.*, 2018; Hennart *et al.*, 2020). We therefore included in the  
713 custom DIPHTOSCAN query database the *spuA* marker and its adjacent genes (DIP0351;  
714 DIP0353; DIP0354; DIP0357=*spuA*), which are strongly associated with biovar Gravis, and the  
715 *narIJHGK* cluster, which is typically absent or partly disrupted, mainly due to mutations in the  
716 *narG* (Hennart *et al.*, 2020) or *narI* (Sangal *et al.*, 2014) in isolates of biovar Belfanti. In the  
717 future, markers of the two biovars of *C. pseudotuberculosis* may be added.

718

#### 719 **Detection of antibiotic resistance genes**

720 Antibiotic resistant genes were identified using AMRfinderPlus, with the database  
721 found at: [https://ftp.ncbi.nlm.nih.gov/pathogen/Antimicrobial\\_resistance/](https://ftp.ncbi.nlm.nih.gov/pathogen/Antimicrobial_resistance/). Features are  
722 detected by using the BLAST family of tools, with identity and coverage defined for each  
723 family of antibiotics (fam.tab). A few genes particularly relevant for the *CdSC* were added to  
724 this database: *pbp2m* (Forde *et al.*, 2020; Hennart *et al.*, 2020) and mutation points of *rpoB*  
725 (WP\_004566675.1) and *gyrA* (WP\_010933942.1). AMRfinderPlus v3.11.2 is used within  
726 DIPHTOSCAN with no modifications.

727

## 728 **Detection of virulence genes from the *C. diphtheriae* species complex**

729 A custom database of virulence features of *C. diphtheriae* and related species was compiled  
730 from literature for the purposes of this work. We included in the custom query database, a panel  
731 of genetic features for which published experimental evidence of their clinical relevance exists  
732 in *C. diphtheriae* or closely related species (*i.e.*, increased virulence in animal models, or  
733 decreased antimicrobial susceptibility *in vitro*) (**Table S2**). These target genes are the  
734 following: *tox*, SpaA-, SpaD-, and SpaH-type pili gene clusters, DIP0733 (67-72p), the genes  
735 DIP1281 and DIP1621 that code for proteins of the NlpC/P60 family, DIP0543 (*nanH*),  
736 DIP1546 and DIP2093 (Ott, 2018) and *pld* (phospholipase). A second panel of genetic features  
737 with no experimental evidence but with strong suspicion for a role in virulence, based on  
738 homology with genes from other pathogens, was also included for broader screening of  
739 virulence features (**Table S2**).

740 For the main virulence factor, the *tox* gene, we used a reference sequence of this gene  
741 from each of *C. diphtheriae*, *C. ulcerans* and *C. pseudotuberculosis* (WP\_003850266.1,  
742 WP\_014835773.1 and WP\_014654963.1, respectively), as the toxin differs between these  
743 species (Dangel *et al.*, 2019).

744 The *tox* gene may be disrupted in some strains by the occurrence of stop codons or other  
745 genetic events, leading to non-toxigenic, *tox*-gene bearing (NTTB) isolates (Zakikhany *et al.*,  
746 2014; Melnikov *et al.*, 2022). DIPHTOSCAN provides information on the putative toxicity of a  
747 strain from the *tox* gene sequence using a categorization into four possible outputs, following  
748 the convention proposed in Kleborate (Lam *et al.*, 2021): (i) if the sequence in the analyzed  
749 genome is identical to the reference *tox* sequence from NTCT13129 strain, the output provides  
750 the name of the sequence with the denomination of the species (*e.g.*, *tox\_diphtheriae*); (ii) If  
751 the sequence in the analyzed genome has a coverage length identical to the reference, but an  
752 identity different from 100%, then an asterisk (\*) is added (*e.g.*, *tox\_diphtheriae\**); (iii) If the  
753 hit coverage length is smaller than the reference length, the tag '-NTTB?- xx%' is added, where  
754 xx is the percentage of the missing sequence length compared to the reference length); (iv)  
755 Finally, if the truncated *tox* sequence is located at the end of a contig, the symbol '\$' is added,  
756 to highlight that the prediction is uncertain.

757 Virulence genes were identified using the method of AMRfinderPlus but based on our  
758 custom database of virulence features. The virulence genes are detected by BLASTn with  
759 thresholds of minimum 80% identity and 50% coverage. Based on the output of  
760 AMRfinderPlus, the gene completion and allele similarity is reported as described above for  
761 the *tox* gene following the Kleborate convention.

762

763

764 **Code availability**

765 The DIPHTOSCAN code is available at <https://gitlab.pasteur.fr/BEBP/diphtoscan>.

766



## 767 **Acknowledgements**

768 We thank Martin Maiden and Keith Jolley (Oxford University) for maintaining the previous  
769 MLST data from Oxford's PubMLST database and for providing the data for import into the  
770 BIGSdb-Pasteur *C. diphtheriae* species complex database.

771

## 772 **Funding**

773 MH was supported financially by the PhD grant "Codes4strains" from the European Joint  
774 Programme One Health, which has received funding from the European Union's Horizon 2020  
775 Research and Innovation Programme under Grant Agreement No. 773830. This work used the  
776 computational and storage services provided by the IT department at Institut Pasteur. The  
777 National Reference Center for Corynebacteria of the Diphtheriae Complex is supported  
778 financially by the Ministry of Health (Public Health France) and Institut Pasteur.

779

## 780 **Declaration of interest statement**

781 The authors declare no conflict of interest.

782

## 783 **Ethical approval statement**

784 Diphtheria is a notifiable disease in France. Phenotypic and genotypic analyses of bacterial  
785 isolates were carried out within the framework of the mandate given to the National Reference  
786 Center for Corynebacteria of the Diphtheriae Complex by the Ministry of Health (Public Health  
787 France). All French bacteriological samples and data were collected in the frame of the French  
788 national diphtheria surveillance and are collected, coded, shipped, managed and analyzed  
789 according to the French National Reference Center protocols. Other strains were obtained from  
790 culture collections.

791

## 792 **Author contributions**

793 S. Brisse (S.B.) conceived, designed, and coordinated the study. Melanie Hennart (M.H.)  
794 developed the DIPHTOSCAN tool with input from SB. M.H. and S.B. analyzed the genomic data.  
795 M.H. created the figures and tables. S.B. and M.H. created the first draft of the manuscript,  
796 worked together to improve it and reviewed the final version. Chiara Crestani analyzed the iron  
797 metabolism and fimbriae genes distribution and wrote the first version of the corresponding  
798 sections. Sebastien Bridel performed the merger of the Oxford PubMLST and BIGSdb-Pasteur  
799 databases. Annick Carmi-Leroy, Sylvie Brémont, Annie Landier, Nathalie Armatys and  
800 Virginie Passet provided technical assistance with the microbiological characterization and  
801 sequencing of the *C. diphtheriae* isolates. Edgar Badell and Julie Toubiana contributed to the

802 NRC operations coordination. Laure Fonteneau and Sophie Vaux coordinated diphtheria  
803 epidemiological surveillance in France. All authors reviewed and approved the final contents  
804 of the manuscript.

805

### 806 **Authors license statement**

807 This research was funded, in whole or in part, by Institut Pasteur and by European Union's  
808 Horizon 2020 research and innovation programme. For the purpose of open access, the authors  
809 have applied a CC-BY public copyright license to any Author Manuscript version arising from  
810 this submission.

811

## Figure legends

### 812 **Figure 1. Phylogenetic tree of *Corynebacterium diphtheriae* from France, 2022**

813 The tree was obtained by maximum likelihood based on a multiple sequence alignment of the  
814 core genome. The scale bar represents the number of nucleotide substitutions per site. The first  
815 column that follows the isolates identifiers indicates the geographic origin (place of isolation;  
816 see key). Travel history provides the most distant geographic region of reported travel (see key);  
817 note that Afghanistan was included in Near and Middle East; and Egypt was included in North  
818 Africa. The stars represent the presence (red star), presence but disruption (NTTB, orange) or  
819 absence (white star) of the diphtheria toxin *tox* gene. Biovars are represent in colored squares,  
820 and *spuA* gene presence by a dark green circle. MLST STs, sublineage (SL) and genomic  
821 clusters are provided with an alternation of colored strips. Identifiers of the main STs are  
822 indicated (note the strong concordance between ST and cgMLST sublineages). The 10 next  
823 colored columns correspond to the presence of at least one gene or mutation (for quinolone and  
824 rifamycin classes) involved in resistance to the indicated class of antimicrobial agents. Last, the  
825 presence of integron-related structures (Cury *et al.*, 2016) is indicated: In0 (integron integrase  
826 and no *attC* sites), CALIN (clusters of *attC* sites lacking integron-integrases) and complete  
827 integrons (integrase and at least one *attC* site). The simultaneous presence of In0 and CALIN  
828 may denote their presence in different contigs even though the integron might be complete.

829

### 830 **Figure 2. Sublineage distribution of *tox* gene and resistance score**

831 (Top) Bar length correspond to the number of isolates per sublineage (deduplicated global  
832 dataset, 976 isolates). Upper part: isolates with non-disrupted *tox* are colored in red, with  
833 disrupted *tox* (NTTB) in orange, and not carrying the *tox* gene in white. Lower part: bar sectors  
834 are colored by resistance score (including beta-lactams and macrolides; see key).

835 (Bottom) Bar length correspond to the number of isolates per sublineage (France, 2022 dataset,  
836 101 isolates). Bar sectors are colored as in the top panel.

837

### 838 **Figure 3. Phylogenetic tree of *Corynebacterium diphtheriae***

839 The tree was obtained by maximum likelihood based on a multiple sequence alignment of the  
840 core genome, and was rooted with *C. belfantii* (not shown). The scale bar gives the number of  
841 nucleotide substitutions per site. The main lineages Mitis and Gravis are drawn using purple  
842 and green branches, respectively. The two inner circles indicate MLST and sublineage  
843 alternation, respectively; main sublineages are labeled within the sectors. first ten colored

844 circles around the tree correspond to the different classes of antibiotics. The following circle  
845 indicates the presence, disruption or absence of the diphtheria toxin *tox* gene (see key). The  
846 beta-lactam resistance circle indicates the presence of the *pbp2m* gene, while the macrolide  
847 circle corresponds to the presence of *ermX* or *ermC* (darker color: presence of the genomic  
848 determinant). The most external circle indicates the non-beta-lactam, non-macrolide (NBNM)  
849 resistance score (number of classes with at least one resistance feature), as a blue gradient (see  
850 key). Four reference strains are indicated: strain NCTC13129, which is used as genomic  
851 sequence reference; strain NCTC10648, which is used as the *tox*-positive and toxinogenic  
852 reference strain in PCR and Elek tests, respectively; strain NCTC11397<sup>T</sup>, which is the  
853 taxonomic type strain of the *C. diphtheriae* species; and the vaccine production strain PW8.

854

#### 855 **Figure 4. Observed frequencies of resistance genes or mutations**

856 The number of genomes with a genetic feature associated with resistance, per antimicrobial  
857 class. Left: Isolates from France, 2022 (n=101 genomes); Right: global deduplicated dataset  
858 (n=976 genomes). The bars are ordered vertically by decreasing frequency in the right panel  
859 and the bar sectors are colored according to the presence of resistance features (see keys).

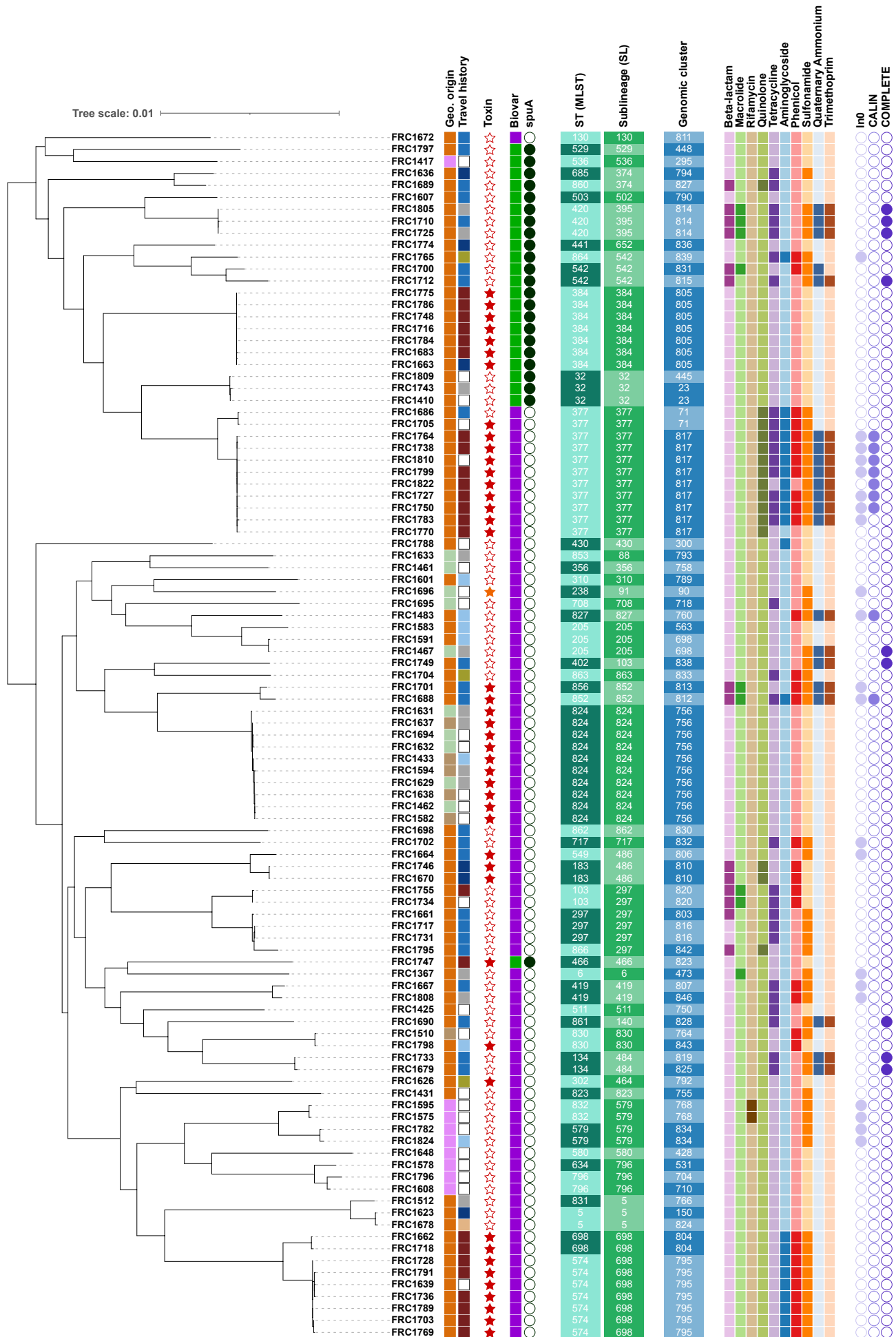
860

## References

- 861 Anderson,J.S. *et al.* (1931) On the existence of two forms of diphtheria bacillus—B. Diphtheriae gravis  
862 and B. Diphtheriae mitis—and a new medium for their differentiation and for the bacteriological  
863 diagnosis of diphtheria. *The Journal of Pathology and Bacteriology*, **34**, 667–681.
- 864 Arcari,G. *et al.* (2023) Multidrug-resistant toxigenic *Corynebacterium diphtheriae* sublineage 453 with  
865 two novel resistance genomic islands. *Microbial Genomics*, **9**.
- 866 Badell,E. *et al.* (2020) *Corynebacterium rouxii* sp. nov., a novel member of the diphtheriae species  
867 complex. *Res. Microbiol.*
- 868 Badell,E. *et al.* (2019) Improved quadruplex real-time PCR assay for the diagnosis of diphtheria. *J. Med.*  
869 *Microbiol.*, **68**, 1455–1465.
- 870 Badell,E. *et al.* (2021) Ongoing diphtheria outbreak in Yemen: a cross-sectional and genomic  
871 epidemiology study. *The Lancet Microbe*, **2**, e386–e396.
- 872 Badenschier,F. *et al.* (2022) Outbreak of imported diphtheria with *Corynebacterium diphtheriae* among  
873 migrants arriving in Germany, 2022. *Euro Surveill*, **27**, 2200849.
- 874 Bankevich,A. *et al.* (2012) SPAdes: a new genome assembly algorithm and its applications to single-  
875 cell sequencing. *J. Comput. Biol.*, **19**, 455–477.
- 876 Barksdale,L. (1970) *Corynebacterium diphtheriae* and its relatives. *Bacteriol Rev*, **34**, 378–422.
- 877 Barraud,O. *et al.* (2011) Antimicrobial drug resistance in *Corynebacterium diphtheriae mitis*. *Emerging*  
878 *Infect. Dis.*, **17**, 2078–2080.
- 879 Benamrouche,N. *et al.* (2016) Microbiological and molecular characterization of *Corynebacterium*  
880 *diphtheriae* isolated in Algeria between 1992 and 2015. *Clin. Microbiol. Infect.*, **22**, 1005.e1-  
881 1005.e7.
- 882 Berger,A. *et al.* (2019) Whole genome sequencing suggests transmission of *Corynebacterium*  
883 *diphtheriae*-caused cutaneous diphtheria in two siblings, Germany, 2018. *Euro Surveill.*, **24**.
- 884 Bernard,K.A. *et al.* (2019) Increase in detection of *Corynebacterium diphtheriae* in Canada: 2006-2019.  
885 *Can Commun Dis Rep*, **45**, 296–301.
- 886 Bolt,F. *et al.* (2010) Multilocus sequence typing identifies evidence for recombination and two distinct  
887 lineages of *Corynebacterium diphtheriae*. *J Clin Microbiol*, **48**, 4177–85.
- 888 Bonmarin,I. *et al.* (2009) Diphtheria: A zoonotic disease in France? *Vaccine*, **27**, 4196–4200.
- 889 Chorlton,S.D. *et al.* (2019) Whole-genome sequencing of *Corynebacterium diphtheriae* isolates  
890 recovered from an inner-city population demonstrates the predominance of a single molecular  
891 strain. *J. Clin. Microbiol.*
- 892 Criscuolo,A. (2020) On the transformation of MinHash-based uncorrected distances into proper  
893 evolutionary distances for phylogenetic inference. *F1000Res*, **9**, 1309.
- 894 Criscuolo,A. and Brisse,S. (2013) AlienTrimmer: A tool to quickly and accurately trim off multiple  
895 short contaminant sequences from high-throughput sequencing reads. *Genomics*,  
896 10.1016/j.ygeno.2013.07.011.
- 897 Cury,J. *et al.* (2016) Identification and analysis of integrons and cassette arrays in bacterial genomes.  
898 *Nucleic Acids Res*, **44**, 4539–4550.
- 899 Dangel,A. *et al.* (2020) *Corynebacterium silvaticum* sp. nov., a unique group of NTTB corynebacteria  
900 in wild boar and roe deer. *Int J Syst Evol Microbiol*, **70**, 3614–3624.
- 901 Dangel,A. *et al.* (2018) Geographically Diverse Clusters of Nontoxigenic *Corynebacterium diphtheriae*  
902 Infection, Germany, 2016-2017. *Emerging Infect. Dis.*, **24**, 1239–1245.
- 903 Dangel,A. *et al.* (2019) NGS-based phylogeny of diphtheria-related pathogenicity factors in different  
904 *Corynebacterium* spp. implies species-specific virulence transmission. *BMC Microbiol.*, **19**, 28.
- 905 Dazas,M. *et al.* (2018) Taxonomic status of *Corynebacterium diphtheriae* biovar Belfanti and proposal  
906 of *Corynebacterium belfantii* sp. nov. *Int. J. Syst. Evol. Microbiol.*, **68**, 3826–3831.
- 907 Engler,K.H. *et al.* (1997) A modified Elek test for detection of toxigenic corynebacteria in the diagnostic  
908 laboratory. *J. Clin. Microbiol.*, **35**, 495–498.
- 909 Feldgarden,M. *et al.* (2021) AMRFinderPlus and the Reference Gene Catalog facilitate examination of  
910 the genomic links among antimicrobial resistance, stress response, and virulence. *Sci Rep*, **11**,  
911 12728.

- 912 Forde,B.M. *et al.* (2020) Fatal respiratory diphtheria caused by  $\beta$ -lactam-resistant *Corynebacterium*  
913 *diphtheriae*. *Clin. Infect. Dis.*
- 914 Grimont,P.A.D. *et al.* (2004) International nomenclature for *Corynebacterium diphtheriae* ribotypes.  
915 *Res Microbiol*, **155**, 162–166.
- 916 Guglielmini,J. *et al.* (2021) Genomic Epidemiology and Strain Taxonomy of *Corynebacterium*  
917 *diphtheriae*. *J Clin Microbiol*, **59**, e0158121.
- 918 Hennart,M. *et al.* (2022) A Dual Barcoding Approach to Bacterial Strain Nomenclature: Genomic  
919 Taxonomy of *Klebsiella pneumoniae* Strains. *Molecular Biology and Evolution*, **39**, msac135.
- 920 Hennart,M. *et al.* (2020) Population genomics and antimicrobial resistance in *Corynebacterium*  
921 *diphtheriae*. *Genome Med*, **12**, 107.
- 922 Hoefer,A. *et al.* (2020) Molecular and epidemiological characterisation of toxigenic and non-toxigenic  
923 *C. diphtheriae*, *C. belfantii* and *C. ulcerans* isolates identified in Spain from 2014 to 2019. *J Clin*  
924 *Microbiol.*
- 925 Jain,C. *et al.* (2018) High throughput ANI analysis of 90K prokaryotic genomes reveals clear species  
926 boundaries. *Nat Commun*, **9**, 5114.
- 927 Jolley,K.A. and Maiden,M.C. (2010) BIGSdb: Scalable analysis of bacterial genome variation at the  
928 population level. *BMC Bioinformatics*, **11**, 595.
- 929 Kofler,J. *et al.* (2022) Ongoing toxin-positive diphtheria outbreaks in a federal asylum centre in  
930 Switzerland, analysis July to September 2022. *Euro Surveill*, **27**, 2200811.
- 931 Konstantinidis,K.T. and Tiedje,J.M. (2005) Genomic insights that advance the species definition for  
932 prokaryotes. *Proceedings of the National Academy of Sciences*, **102**, 2567–2572.
- 933 Lam,M.M.C. *et al.* (2021) A genomic surveillance framework and genotyping tool for *Klebsiella*  
934 *pneumoniae* and its related species complex. *Nat Commun*, **12**, 4188.
- 935 Liu,Y. *et al.* (2013) Musket: a multistage k-mer spectrum-based error corrector for Illumina sequence  
936 data. *Bioinformatics*, **29**, 308–315.
- 937 Magoč,T. and Salzberg,S.L. (2011) FLASH: fast length adjustment of short reads to improve genome  
938 assemblies. *Bioinformatics*, **27**, 2957–2963.
- 939 Mandlik,A. *et al.* (2007) *Corynebacterium diphtheriae* employs specific minor pilins to target human  
940 pharyngeal epithelial cells. *Mol Microbiol*, **64**, 111–124.
- 941 McLeod,J.W. (1943) THE TYPES MITIS, INTERMEDIUS AND GRAVIS OF  
942 *CORYNEBACTERIUM DIPHTHERIAE*: A Review of Observations during the Past Ten  
943 Years. *Bacteriol Rev*, **7**, 1–41.
- 944 Meinel,D.M. *et al.* (2016) Outbreak investigation for toxigenic *Corynebacterium diphtheriae* wound  
945 infections in refugees from Northeast Africa and Syria in Switzerland and Germany by whole  
946 genome sequencing. *Clin. Microbiol. Infect.*, **22**, 1003.e1-1003.e8.
- 947 Melnikov,V.G. *et al.* (2022) Detection of diphtheria toxin production by toxigenic corynebacteria using  
948 an optimized Elek test. *Infection*, **50**, 1591–1595.
- 949 Mina,N.V. *et al.* (2011) Canada’s first case of a multidrug-resistant *Corynebacterium diphtheriae* strain,  
950 isolated from a skin abscess. *J. Clin. Microbiol.*, **49**, 4003–4005.
- 951 Mokrousov,I. (2009) *Corynebacterium diphtheriae*: genome diversity, population structure and  
952 genotyping perspectives. *Infect Genet Evol*, **9**, 1–15.
- 953 Mueller,J.H. (1941) Toxin-production as related to the clinical severity of diphtheria. *J Immunol*, **42**,  
954 353–360.
- 955 Néron,B. *et al.* (2022) IntegronFinder 2.0: Identification and Analysis of Integrons across Bacteria, with  
956 a Focus on Antibiotic Resistance in *Klebsiella*. *Microorganisms*, **10**, 700.
- 957 Ondov,B.D. *et al.* (2016) Mash: fast genome and metagenome distance estimation using MinHash.  
958 *Genome Biol.*, **17**, 132.
- 959 Ott,L. (2018) Adhesion properties of toxigenic corynebacteria. *AIMS Microbiol*, **4**, 85–103.
- 960 Ott,L. *et al.* (2022) Interactions between the Re-Emerging Pathogen *Corynebacterium diphtheriae* and  
961 Host Cells. *International Journal of Molecular Sciences*, **23**, 3298.
- 962 Pappenheimer,A.M. and Murphy,J.R. (1983) Studies on the molecular epidemiology of diphtheria.  
963 *Lancet*, **2**, 923–926.
- 964 Peixoto,R.S. *et al.* Functional characterization of the collagen-binding protein DIP2093 and its influence  
965 on host–pathogen interaction and arthritogenic potential of *Corynebacterium diphtheriae*.  
966 *Microbiology*, **163**, 692–701.

- 967 du Plessis, M. *et al.* (2017) Molecular Characterization of *Corynebacterium diphtheriae* Outbreak  
968 Isolates, South Africa, March-June 2015. *Emerging Infect. Dis.*, **23**, 1308–1315.
- 969 Polonsky, J.A. *et al.* (2021) Epidemiological, clinical, and public health response characteristics of a  
970 large outbreak of diphtheria among the Rohingya population in Cox’s Bazar, Bangladesh, 2017  
971 to 2019: A retrospective study. *PLoS Med*, **18**, e1003587.
- 972 Reardon-Robinson, M.E. and Ton-That, H. (2014) Assembly and function of *Corynebacterium*  
973 *diphtheriae* pili. In, *Corynebacterium diphtheriae and Related Toxigenic Species*. Springer,  
974 Heidelberg, pp. 123–141.
- 975 Rogers, E.A. *et al.* (2011) Adhesion by pathogenic corynebacteria. *Adv Exp Med Biol*, **715**, 91–103.
- 976 Russell, L.M. and Holmes, R.K. (1985) Highly toxinogenic but avirulent Park-Williams 8 strain of  
977 *Corynebacterium diphtheriae* does not produce siderophore. *Infect Immun*, **47**, 575–578.
- 978 Sangal, V. *et al.* (2014) A lack of genetic basis for biovar differentiation in clinically important  
979 *Corynebacterium diphtheriae* from whole genome sequencing. *Infect. Genet. Evol.*, **21**, 54–57.
- 980 Sangal, V. and Hoskisson, P.A. (2016) Evolution, epidemiology and diversity of *Corynebacterium*  
981 *diphtheriae*: New perspectives on an old foe. *Infect. Genet. Evol.*, **43**, 364–370.
- 982 Santos, A.S. *et al.* (2018) Searching whole genome sequences for biochemical identification features of  
983 emerging and reemerging pathogenic *Corynebacterium* species. *Funct. Integr. Genomics*, **18**,  
984 593–610.
- 985 Schaeffer, J. *et al.* (2020) Assessing the Genetic Diversity of Austrian *Corynebacterium diphtheriae*  
986 Clinical Isolates, 2011–2019. *J Clin Microbiol*.
- 987 Seth-Smith, H.M.B. and Egli, A. (2019) Whole Genome Sequencing for Surveillance of Diphtheria in  
988 Low Incidence Settings. *Front Public Health*, **7**, 235.
- 989 Tauch, A. *et al.* (2003) Insights into the genetic organization of the *Corynebacterium diphtheriae*  
990 erythromycin resistance plasmid pNG2 deduced from its complete nucleotide sequence.  
991 *Plasmid*, **49**, 63–74.
- 992 Tauch, A. *et al.* (1995) The *Corynebacterium xerosis* composite transposon Tn5432 consists of two  
993 identical insertion sequences, designated IS1249, flanking the erythromycin resistance gene  
994 *ermCX*. *Plasmid*, **34**, 119–131.
- 995 Timms, V.J. *et al.* (2018) Genome-wide comparison of *Corynebacterium diphtheriae* isolates from  
996 Australia identifies differences in the Pan-genomes between respiratory and cutaneous strains.  
997 *BMC Genomics*, **19**, 869.
- 998 Truelove, S.A. *et al.* (2020) Clinical and Epidemiological Aspects of Diphtheria: A Systematic Review  
999 and Pooled Analysis. *Clin Infect Dis*, **71**, 89–97.
- 1000 WHO (2018) Diphtheria: Vaccine Preventable Diseases Surveillance Standards.  
1001 [https://www.who.int/publications/m/item/vaccine-preventable-diseases-surveillance-](https://www.who.int/publications/m/item/vaccine-preventable-diseases-surveillance-standards-diphtheria)  
1002 [standards-diphtheria](https://www.who.int/publications/m/item/vaccine-preventable-diseases-surveillance-standards-diphtheria).
- 1003 Will, R.C. *et al.* (2021) Spatiotemporal persistence of multiple, diverse clades and toxins of  
1004 *Corynebacterium diphtheriae*. *Nat Commun*, **12**, 1500.
- 1005 Williams, M.M. *et al.* (2020) Detection and Characterization of Diphtheria Toxin Gene-Bearing  
1006 *Corynebacterium* Species through a New Real-Time PCR Assay. *J Clin Microbiol*, **58**.
- 1007 Xiaoli, L. *et al.* (2020) Genomic epidemiology of nontoxigenic *Corynebacterium diphtheriae* from King  
1008 County, Washington State, USA between July 2018 and May 2019. *Microb Genom*, **6**.
- 1009 Zakikhany, K. *et al.* (2014) Emergence and molecular characterisation of non-toxigenic tox gene-bearing  
1010 *Corynebacterium diphtheriae* biovar *mitis* in the United Kingdom, 2003–2012. *Euro Surveill*,  
1011 **19**.
- 1012 Zasada, A.A. (2014) Antimicrobial Susceptibility and Treatment. In, Burkovski, A. (ed),  
1013 *Corynebacterium diphtheriae and Related Toxigenic Species: Genomics, Pathogenicity and*  
1014 *Applications*. Springer Netherlands, Dordrecht, pp. 239–246.
- 1015



**Country**

- Metropolitan France
- Mayotte
- La Reunion
- French Guiana

**Travel history**

- Near and Middle East
- North Africa
- West Africa
- Southern Africa
- Southern and South-Eastern Asia
- North America
- None
- Not documented

**tox gene and Elek's test**

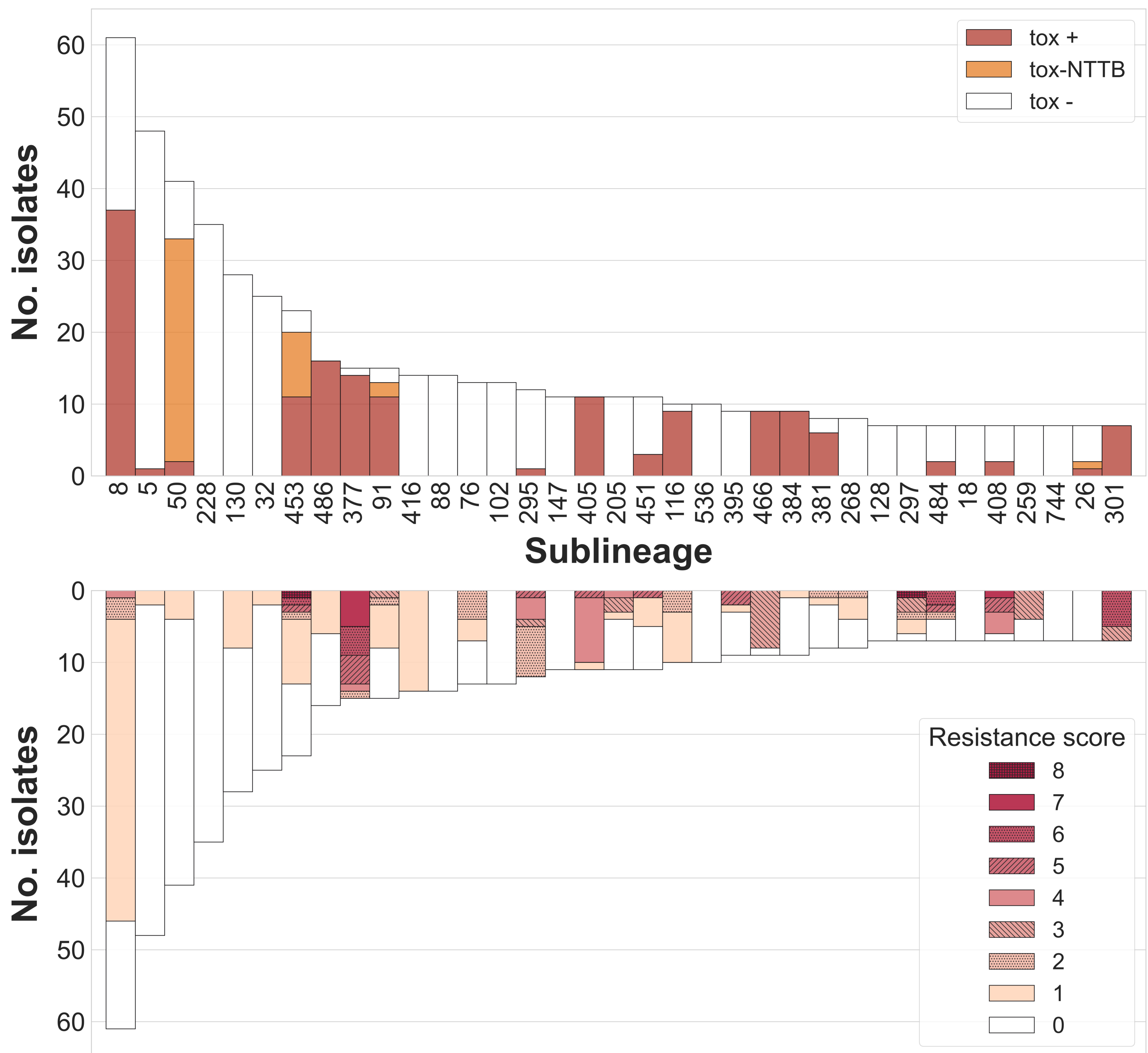
- tox + and Elek +
- NTTB
- tox - and Elek -

**Biovar**

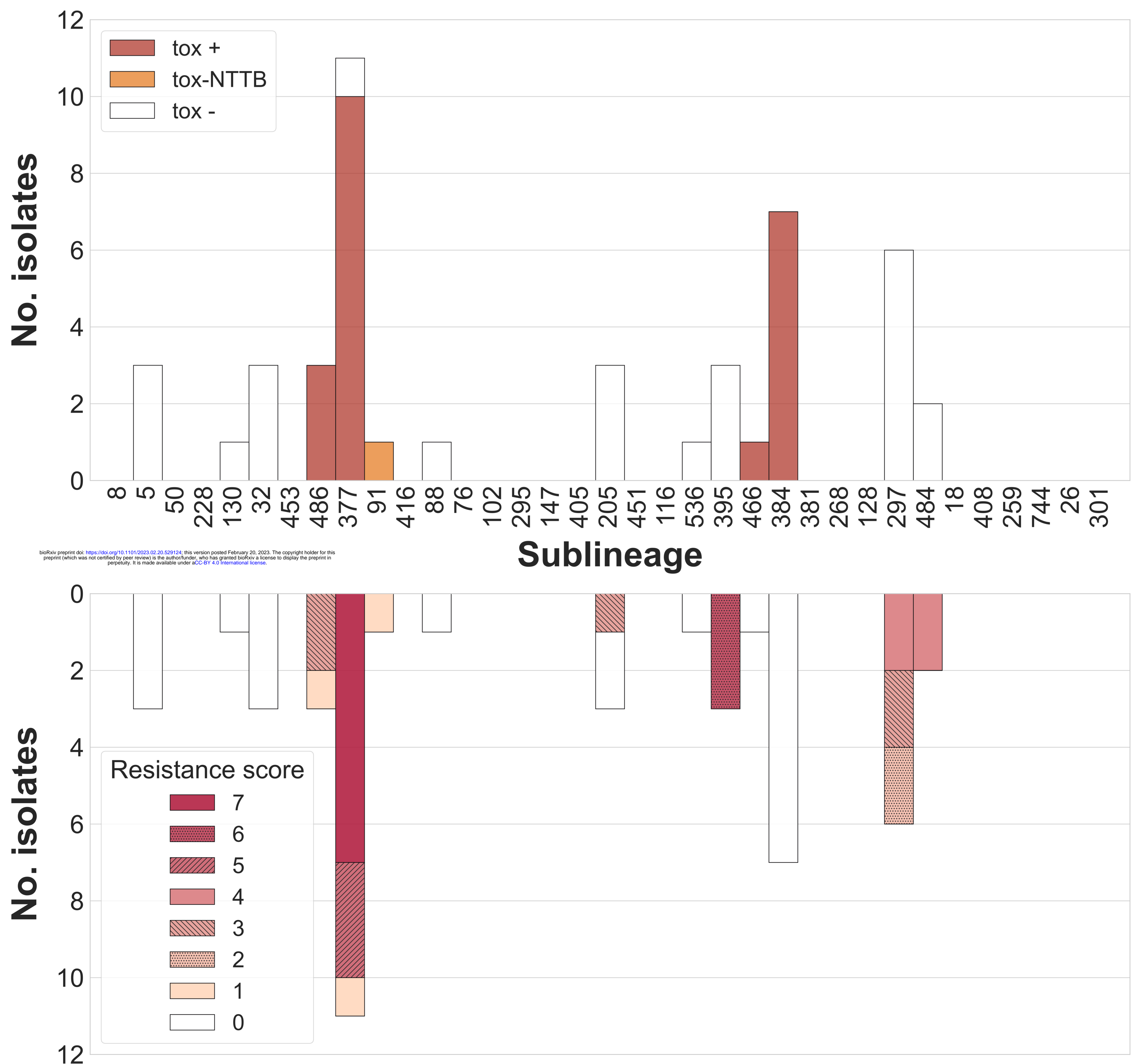
- Mitis
- Gravis



# A. Global dataset

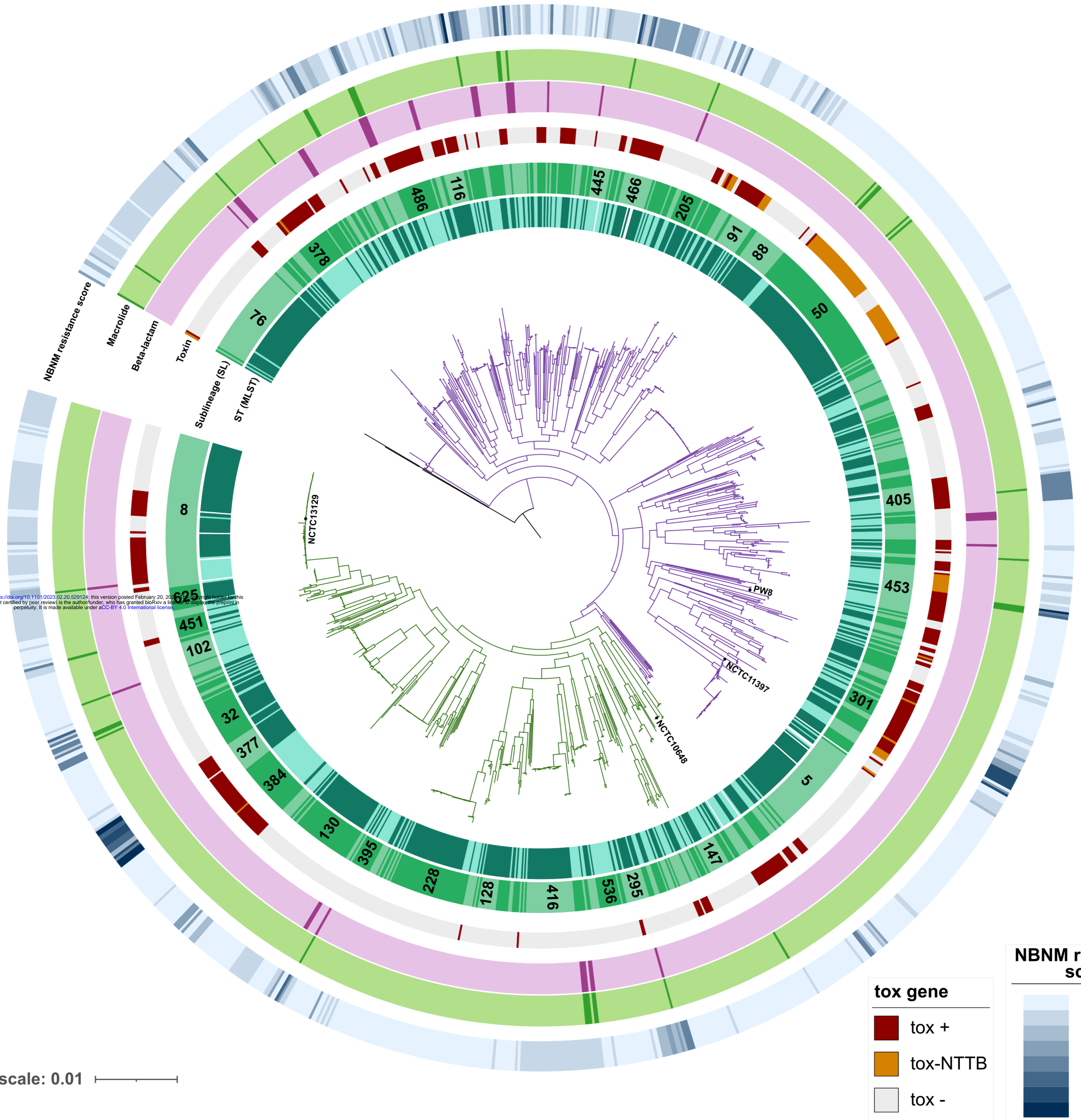


# B. France, 2022



bioRxiv preprint doi: <https://doi.org/10.1101/2023.02.20.529124>; this version posted February 20, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

Tree scale: 0.01



# A. France, 2022

# B. Global dataset

