

1 **In-depth analysis of *Bacillus anthracis* 16S rRNA genes and**
2 **transcripts reveals intra- and intergenomic diversity and**
3 **facilitates anthrax detection**

4 Peter Braun^a, Fee Zimmermann^a, Mathias C. Walter^a, Sonja Mantel^a, Karin Aistleitner[†], Inga
5 Stürz^a, Gregor Grass^{a##*} and Kilian Stoecker^{a#}

6

7

8 ^aBundeswehr Institute of Microbiology, Munich, Germany

9 [†]Deceased 31 March 2019

10 * corresponding author

11 #shared last

12

13 **Abstract**

14 Analysis of 16S ribosomal RNA (rRNA) genes provides a central means of taxonomic
15 classification of bacterial species. Based on presumed sequence identity among species of
16 the *Bacillus cereus sensu lato* group, the 16S rRNA genes of *B. anthracis* have been
17 considered unsuitable for diagnosis of the anthrax pathogen. With the recent identification
18 of a single nucleotide polymorphism in some 16S rRNA gene copies, specific identification of
19 *B. anthracis* becomes feasible. Here, we designed and evaluated a set of *in situ*-, *in vitro*- and
20 *in silico*-assays to assess the yet unknown 16S-state of *B. anthracis* from different
21 perspectives. Using a combination of digital PCR, fluorescence *in situ* hybridization, long-read
22 genome sequencing and bioinformatics we were able to detect and quantify a unique 16S
23 rRNA gene allele of *B. anthracis* (16S-BA-allele). This allele was found in all available *B.*
24 *anthracis* genomes and may facilitate differentiation of the pathogen from any close relative.
25 Bioinformatics analysis of 959 *B. anthracis* genome data-sets inferred that abundances and
26 genomic arrangements of the 16S-BA-allele and the entire rRNA operon copy-numbers differ
27 considerably between strains. Expression ratios of 16S-BA-alleles were proportional to the
28 respective genomic allele copy-numbers. The findings and experimental tools presented here
29 provide detailed insights into the intra- and intergenomic diversity of 16S rRNA genes and
30 may pave the way for improved identification of *B. anthracis* and other pathogens with
31 diverse rRNA operons.

32

33 Introduction

34 Anthrax, caused by the spore-forming bacterium *Bacillus anthracis*, is a disease of animals
35 but can also affect humans either through contact with infected animals and their products
36 or as a consequence of deliberate acts of bioterrorism^{1,2}. Because of its high pathogenicity,
37 rapid, sensitive and unambiguous identification of the pathogen is vital. However, diagnostic
38 differentiation of *B. anthracis* from its closest relatives of the *Bacillus cereus sensu lato* group
39 is challenging. Phenotypic properties are not species-specific and nearly identical derivatives
40 of the anthrax virulence plasmids can also be found in related bacilli².

41 In spite of earlier work³ rRNA gene sequences have not been deemed discriminatory
42 for unambiguous distinction of *B. anthracis* from its closest relatives due to the lack of
43 specific sequence variations. Recent analysis of 16S rRNA gene alleles of *B. anthracis* and
44 relatives, however, revealed an unexpected SNP (Single Nucleotide Polymorphism) at
45 position 1110 (position 1139 in⁴; 1110 according to *B. anthracis* strain Ames Ancestor,
46 NC_007530) in some of the 16S rRNA gene copies⁴. This SNP has previously been missed,
47 most likely because it is present only in some of the total eleven 16S rRNA gene copies⁴.
48 Despite the high abundance of more than 1,000 publicly available short-read genomic
49 datasets and more than 260 genome assemblies, reliable information about sequence
50 variations within *B. anthracis* rRNA operons is still scarce due to the limitations of short-read
51 whole genome sequencing (WGS) and subsequent reference mapping to detect sequence
52 variations in paralogous, multi-copy genes. Producing high-quality genomes e.g. through
53 hybrid assemblies of long- and short-read approaches would help bridge this gap.

54 In this study, we validated a species-discriminatory SNP within the 16S rRNA genes of *B.*
55 *anthracis* using a set of different *in situ*, *in vitro* and *in silico* approaches on both genomic and
56 transcript levels. Through this work, we established new diagnostic tools for *B. anthracis*
57 including a fluorescence *in situ* hybridization (FISH) assay and a digital PCR (dPCR) test for
58 both genomic and transcript identification and quantification. Finally, we expanded our
59 analysis on all short-read *B. anthracis* data sets available in the NCBI Short Read Archive
60 (SRA) and calculated the rRNA operon copy-numbers and allele frequencies using a coverage-
61 ratio based bioinformatics approach.

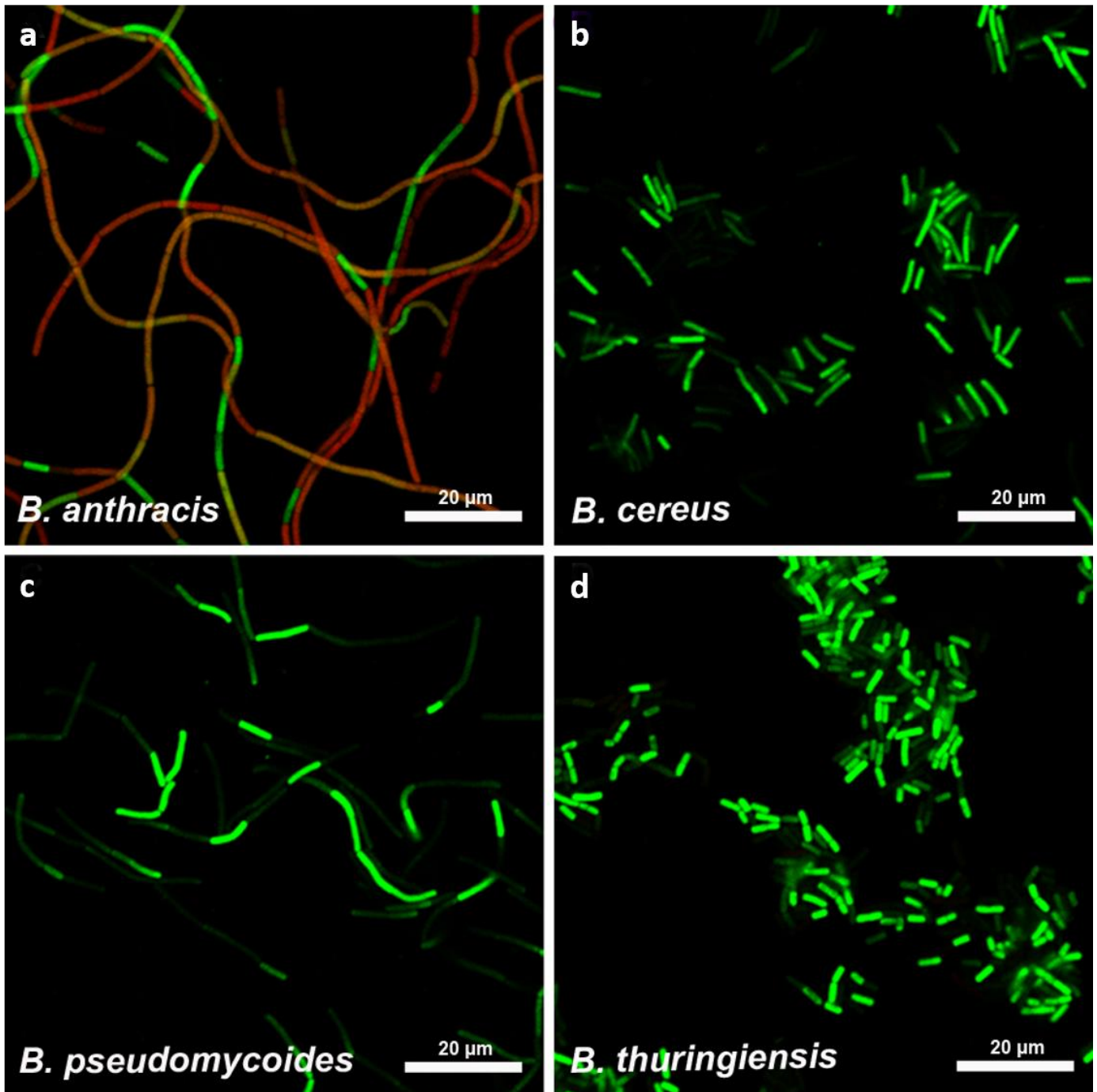
62 **Results**

63 **A SNP in transcripts of 16S rRNA genes enables specific microscopic detection** 64 **of *B. anthracis* by fluorescence *in situ* hybridization.**

65 Triggered by earlier data on a unique SNP position in some copies of the 16S rRNA gene of *B.*
66 *anthracis* (guanine to adenine transition at position 1110)⁴, we aimed at developing a new
67 FISH assay for the identification of *B. anthracis*. Previous work has introduced a probe set for
68 the FISH based identification of *B. anthracis*⁵. Evaluation of the probe sequences revealed,
69 however, that they are unsuitable for unambiguous *B. anthracis* identification due to
70 unspecific probe binding⁶. Thus, we designed a FISH probe for discriminating *B. anthracis*
71 from all of its close relatives targeting this specific SNP in 16S rRNA genes (Probe
72 BA_SNP_Cy3). Additionally, we developed probe BC_SNP_FAM which binds to 16S rRNA
73 sequence found in all *B. cereus s. l.* strains, including *B. anthracis* (Supplementary Table S2).
74 No other bacterial or archaeal 16S rRNA gene in the SILVA database had a full match for both

75 of the newly designed probes (assessed 2021-03-01). In order to increase signal intensity and
76 stringency⁷ we incorporated two locked nucleic acids (LNA) in probe BA_SNP_Cy3 and one
77 LNA in probe BC_SNP_FAM. Optimum formamide concentrations in the hybridization buffer
78 of this FISH assay was titrated and finally set at 30% (v/v) formamide for species
79 differentiation (Supplementary Figure S2).

80 For assay validation, the 16S rRNA probes were tested against a broad panel of
81 *B. cereus s. l.* strains (Supplementary Table S1). The FISH assay allowed differentiation of *B.*
82 *anthracis* from all other *B. cereus s. l.* group members. *B. anthracis* cells displayed red
83 fluorescence Cy3-signals after hybridization of the specific 16S rRNA variation at position
84 1110, and green fluorescence FAM-signals resulting from hybridization to the divergent 16S
85 rRNA featuring no *B. anthracis* specific SNP (Figure 1). No red Cy3-signals were detected in
86 any of the non-*B. anthracis B. cereus s. l.* group strains.



87

88 **Figure 1. FISH-based microscopic differentiation of *B. anthracis* from other *B. cereus s. l.***
89 **group species.** Representative images for *B. anthracis* (a, strain Bangladesh 28/01) *B. cereus*
90 (b, strain ATCC 6464), *B. pseudomycooides* (c, strain WS 3119) and *B. thuringiensis* (d, strain
91 WS 2614) are shown as overlay images of red (probe BA_SNP_Cy3 / 568nm) and green
92 fluorescent channels (probe BC_SNP_FAM / 520nm).

93 While we found Cy3 FISH signals for all *B. anthracis* strains, we discovered broad
94 variations in Cy3 fluorescence signal intensities for different cells of the same and between
95 different *B. anthracis* strains. Even for cells of the same chain, there were individual cells
96 showing almost uniquely either the Cy3 or the FAM signal, resulting in a mosaic-like pattern
97 (Figure 1). Total fluorescence intensities varied between different *B. anthracis* strains from
98 very strong Cy3 signals to the extreme cases of *B. anthracis* strains ATCC 4229 Pasteur, SA20
99 and A3783, for which Cy3 signals were very weak (for signal intensities see Supplementary
100 Table S1). These findings strongly indicate that the 16S rRNA of *B. anthracis* can be used for
101 microscopy-based specific pathogen detection. Notably, variations in fluorescence
102 intensities suggests differences in the rRNA expression level. As these differences might be
103 caused by a gene dose effect we decided to analyze the genomic distribution of the *B.*
104 *anthracis* specific SNP in 16S rRNA genes.

105 **Genomic analysis of *B. anthracis* genomes reveals variations in 16S-BA-allele** 106 **frequencies.**

107 We correlated FISH results with the abundance of 16S rRNA gene copies harboring the *B.*
108 *anthracis* specific SNP within different *B. anthracis* genomes. Despite the significant number
109 of *B. anthracis* genomes published, the vast majority of sequences has been generated using
110 short-read-sequencing with subsequent mapping to the reference genome (Ames Ancestor
111 NC_007530, ⁸). Due to multiple copies of the rRNA operons, conventional short-read-
112 sequencing and mapping approaches do not allow for reliable detection of allele variations.
113 During *de novo* assembly of short reads, near identical regions like rRNA operon are
114 collapsed into one contig representing only a consensus sequence missing any minor allele

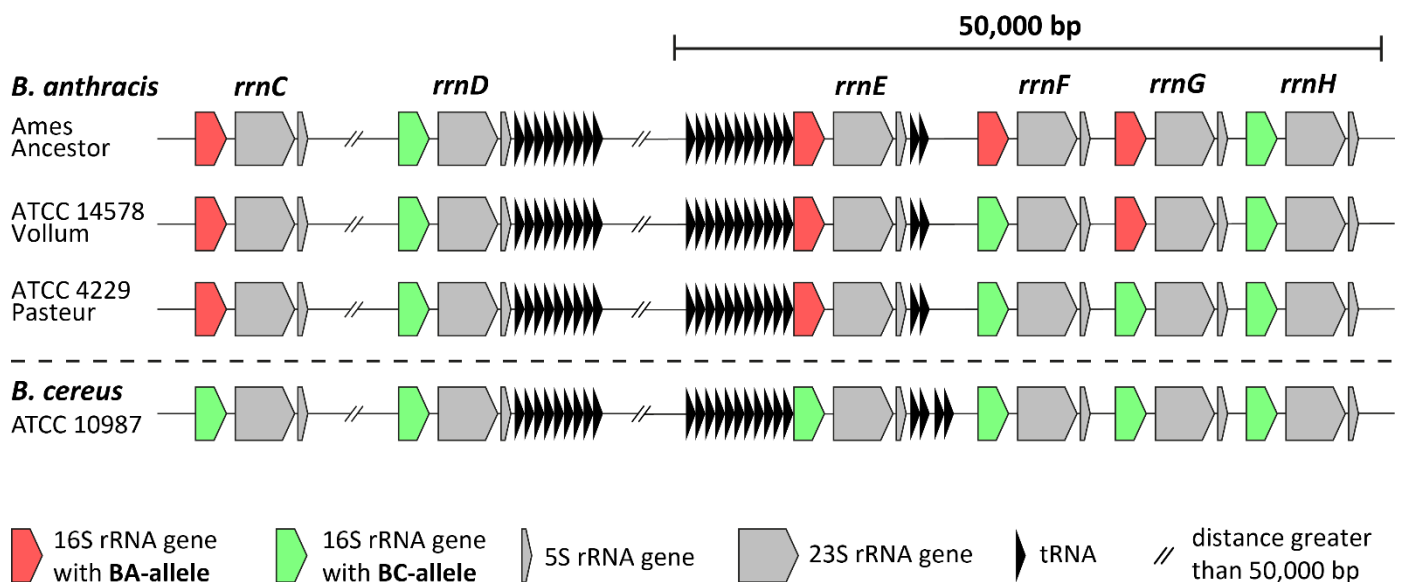
115 variations. Thus, potential differences in allele frequencies can easily be missed. Because of
116 mapping to the reference genome, consensus sequences always feature identical 16S rRNA
117 allele distribution as the reference. Hence, there is a need for high-quality genomes
118 generated by hybrid assemblies using long- and short-read-sequences for obtaining insights
119 into the real distribution and diversity of 16S rRNA alleles in *B. anthracis* genomes.

120 To start meeting this need, we analyzed and compared the 16S gene sequences and
121 locations in all available high-quality genomes of *B. anthracis* (assessed at the end of 2020)
122 that are based on long-read-sequencing and *de novo* assembly. Figure 2 shows a schematic
123 illustration of the genomic organization of rRNA operons including 16S, 23S and 5S ribosomal
124 subunits as well as tRNA genes from operons *rrnC* to *H* (outlying operons A, B, I, J and K are
125 not shown) of representative strains for different 16S rRNA genotypes (Ames Ancestor -
126 NC_007530,⁸; ATCC 14578 Vollum (in-house sequenced; this work; Supplementary Table
127 S3); ATCC 4229 Pasteur - NZ_CP009476,⁹), and closely related *B. cereus* strain ATCC 10987 -
128 NC_003909,¹⁰).

129 We found that all 16S rRNA gene copies featuring the *B. anthracis* specific SNP to
130 have 100% sequence identity representing a distinct allele. For simplification, copies
131 featuring this guanine to adenine transition at position 1110 were termed 16S-BA-(*B.*
132 *anthracis*)-alleles, while all other variants lacking this transition were designated 16S-BC
133 (*B. cereus s. l.*)-alleles.

134 The three *B. anthracis* strains Ames Ancestor, ATCC 14578 Vollum and ATCC 4229
135 Pasteur analyzed above, harbored different 16S-BA/BC-allele frequencies with 4/7, 3/8 and
136 2/9 copies, respectively (Figure 2). No 16S-BA-alleles were found in *B. cereus* ATCC 10987 or
137 any other non-*B. anthracis* strain. In all three *B. anthracis* strains, rRNA operons *rrnA*, *B*, *D*, *H*,

138 *I, J* and *K* carried 16S-BC-alleles, while for *rrnC* and *rrnE* exclusively the 16S-BA-allele was
 139 identified. Only two rRNA operons, *rrnF* and *rrnG*, were found to be variable, with strain
 140 Ames Ancestor harboring two 16S-BA-alleles and strain ATCC 4229 Pasteur only the BC-
 141 alleles for *rrnF* and *rrnG*. Strain ATCC 14578 Vollum exhibited an intermediate state with a
 142 16S-BA-allele in *rrnG* and a BC-allele in *rrnF* (Figure 2). It is thus possible that these
 143 differences in 16S rRNA allele distributions may have caused the observed variations in *B.*
 144 *anthracis* specific FISH signals (Figure 1) by gene-dosage-mediated differences in rRNA
 145 transcription levels.
 146



147 **Figure 2. Schematic illustration of the genomic organization of rRNA operons and**
 148 **distribution of 16S alleles in *B. anthracis*.** Depicted are the 16S, 23S, 5S ribosomal subunit,
 149 and tRNA genes from operons *rrnC* to *H* in strains Ames Ancestor, ATCC 14578 Vollum, ATCC
 150 4229 Pasteur and *B. cereus* ATCC 10987. The 16S rRNA genes are either displayed in red for
 151 16S-BA-alleles or in green for 16S-BC-alleles. Not shown are operons *rrnA, B, I, J* and *K*
 152 exclusively carrying the 16S-BC-allele in any strain. Distances are not to scale.

153 **A tetraplex dPCR assay enables the absolute quantification of species-specific**

154 **16S rRNA gene allele numbers in *B. anthracis*.**

155 To verify this finding and to quantify the ratios of each allele in a diverse panel of
156 *B. anthracis* strains, we designed and tested a hydrolysis-probe-based digital PCR (dPCR)
157 assay (Figure 3a). This assay utilized HEX (green) and FAM (blue) fluorescent dyes labelled
158 allele-specific probes for the 16S-BC-allele and BA-allele, respectively, with both probes
159 targeting the 1110 SNP of the 16S rRNA genes (Supplementary Table S2). In parallel, a
160 previously published second hydrolysis-probe-based PCR assay using HEX dye was adopted
161 for dPCR. This assay targets the *B. anthracis* specific, chromosomal *PL3* gene¹¹. Finally, a pan-
162 *B. cereus s. l.* hydrolysis-probe-based PCR assay on the *gyrA* (gyrase gene) marker using FAM
163 dye was designed, facilitating the detection and quantification of *B. cereus s. l.* species
164 (including *B. anthracis*) chromosomes. In these dPCR assays, the *PL3* and *gyrA* dPCR-tests
165 served as internal controls (for *B. anthracis* and *B. cereus s.l.*, respectively): each positive for
166 *B. anthracis* genomic DNA vs. negative for *PL3* and positive for *gyrA* using genomic DNA of
167 other members of the *B. cereus s. l.* group.

168 These four assays were combined into a single tetraplex dPCR assay. To achieve the
169 required signal separation of the four individual dPCR reactions (on our dPCR-analysis
170 instrument featuring only two channels, FAM and HEX), we deliberately altered the signal
171 output levels by titrating concentrations of probes labeled with the same dye (Figure 3a).
172 Thus, the *PL3* marker assay was tuned to produce high HEX signals vs. low HEX signals
173 coming from 16S-BC-alleles. Likewise, the *gyrA* marker assay was set to produce high FAM
174 signals vs. low FAM signals originating from 16S-BA-alleles. Since both *PL3* and *gyrA* are

175 single-copy genes located on the chromosome of *B. anthracis*, these markers should result in
176 very similar quantitative outputs when individual *B. anthracis* DNA samples are analyzed.
177 Therefore, these markers served as internal quantification controls in this work.

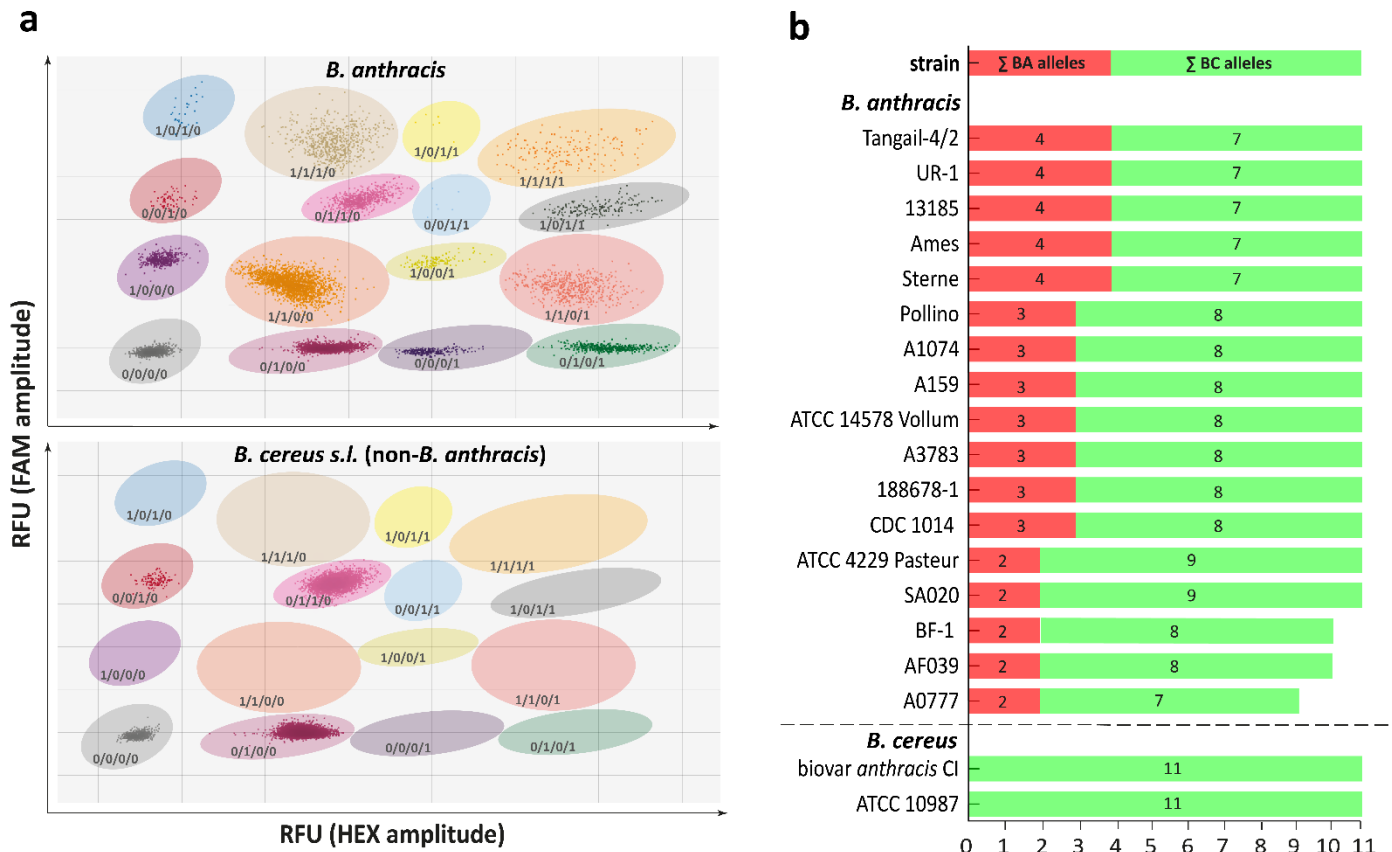
178 A typical analysis output of this tetraplex dPCR assay is exemplified in Figure 3a. In a
179 two-dimensional plot (FAM signal amplitude on the y-axis and HEX signal amplitude on the x-
180 axis) of such tetraplex dPCR data, one can discriminate a specific fluorescence patterns after
181 dPCR representing 16 clusters (when *B. anthracis* DNA was used as a template). Each of the
182 droplets within a cluster contained a certain target combination of *gyrA*, *PL3*, 16S-BA-allele
183 and/or 16S-BC-allele (for example $gyrA^+/PL3^+/16S-BC-allele^+/16S-BA-allele^+$ or $gyrA^-/PL3^-$
184 $/16S-BC-allele^-/16S-BA-allele^-$). Using template DNA originating from a non-*B. anthracis*
185 member of the *B. cereus s. l.* group (i.e., not harboring any 16S-BA-allele), resulted in the
186 expected formation of only four droplet clusters i.e., lacking all signals of *B. anthracis*-
187 specific clusters containing combinations of the *PL3* marker or the 16S-BA-allele (Figure 3a).

188 Testing the assay on the reference strains Ames, ATCC 14578 Vollum and ATCC 4229
189 Pasteur we found four, three and two 16S-BA-alleles, respectively, and eleven 16S rRNA total
190 copies per cell in all three strains. This agreed with the values determined by genomic
191 analysis and, therefore, validated the dPCR assay being able to accurately quantify 16S rRNA
192 alleles in *B. anthracis*.

193 Using the validated tetraplex dPCR assay we analyzed the same strain panel as tested
194 by FISH (Supplementary Table S1). Similar to FISH, there was no signal for 16S-BA-alleles in
195 the 32 non-*B. anthracis* strains of the *B. cereus s. l.* group. However, all of the 17 *B. anthracis*
196 strains harbored at least two (up to four) copies of the 16S-BA-allele per cell (Figure 3b). The
197 majority of *B. anthracis* strains exhibited either the genotypes 4/7 or 3/8 (16S-BA/BC-alleles;

198 six and seven strains, respectively). These predominant genotypes, together with genotype
199 2/9 (strain ATCC 4229 Pasteur and strain SA020) were all found to harbor eleven rRNA
200 operons in total, which agrees with previously determined numbers of rRNA operons in
201 these strains. Conversely, strains A182 and BF-1 harbored only ten 16S gene copies in total
202 (genotype 2/8). Notably, strain A0777 exhibited just nine rRNA copies, two of which
203 contained the *B. anthracis* specific SNP (genotype 2/7)

204



205 **Figure 3. Detection and quantification of 16S rRNA gene alleles in *B. anthracis* and**
 206 ***B. cereus s. l.* strains.** (a) typical results of a tetraplex dPCR assay using *B. anthracis* template
 207 DNA (upper panel) and DNA of a non-*B. anthracis* member of the *B. cereus s. l.* group (lower
 208 panel). With each dot representing a droplet plotted according to its FAM signal-amplitude
 209 (RFU: Relative Fluorescence Units) on the y-axis and HEX signal-amplitude on the x-axis, a
 210 total of 16 (for *B. anthracis*, upper panel) or 4 (non-*B. anthracis* members of the *B. cereus*
 211 *sensu lato* group, lower panel) clusters (defined by shaded areas) can be assigned to a
 212 certain dPCR marker combination of *gyrA* (FAM high signal), *PL3* (HEX high signal), 16S-BA-
 213 (FAM low signal), and 16S-BC-allele (HEX low signal). Since both the *PL3* gene and the 16S-
 214 BA-allele are exclusively found in *B. anthracis*, the 16S-BA- and 16S-BC-allele copy numbers
 215 can be calculated from the positive droplets of single copy genes (*PL3* and *gyrA*) and multi-
 216 copy 16S rRNA genes. All dPCR patterns lacking either (or both) the *PL3* gene and the 16S-
 217 BA-allele clusters represent DNA of a non-*B. anthracis* member of the *B. cereus s. l.* group.
 218 (b) Copy-numbers for 16S-BA- and 16S-BC-alleles for all *B. anthracis* strains (and *B. cereus*
 219 biovar *anthracis* CI) tested.

220 **16S-BA-allele frequencies and total rRNA operon copy numbers vary between**
221 **different *B. anthracis* strains**

222 In order to further confirm dPCR results and to exclude underestimation by dPCR as a
223 possible cause of the unexpected low number of total rRNA operons in strains A182, BF-1
224 and A0777, we conducted a combination of long- and short-read sequencing on these and 32
225 additional *B. anthracis* strains (Supplementary Table S1). A mean read-length of about 15 kb
226 generated by Nanopore sequencing combined with Illumina 2x300 bp paired-end sequencing
227 allowed for the precise assembly of complete genomes including correct positioning of rRNA
228 operons on chromosome. Coverage values of more than 200-fold enabled the accurate
229 quantification of SNPs and therefore, genotypes based on 16S-BA/BC-allele distribution could
230 be reliably determined. The results matched those obtained from dPCR, confirming the
231 accuracy and reliability of the tetraplex assay. We found that strain A0777 lacked rRNA
232 operons *rrnG* and *rrnH*. rRNA operon *rrnG* was not present in strains AF039, SA020 and BF-1.
233 The genome regions downstream of the missing rRNA operons and upstream of the next
234 rRNA operon were also absent.

235 In order to extend our analysis of 16S rRNA allelic states to more *B. anthracis* strains,
236 we expanded our investigation on all publicly available short-read sequence data for *B.*
237 *anthracis* generated using Illumina sequencing technology. Starting from our newly-
238 generated high-quality hybrid assemblies, we developed a *k*-mer and coverage-ratio based
239 tool to calculate the rRNA operon copy-numbers and allele frequencies from all SRA
240 datasets, published until the end of 2020. These numbers of rRNA operons and 16S-BA-
241 alleles (from short-read datasets) were identical to the long-read data of the same genomes

242 (Supplementary Table S4). After this method validation, we analyzed 986 SRA Illumina
243 sequenced datasets for 16S rRNA operon and BA/BC-allele distribution. After assembly and
244 filtering, 959 genomes remained for a detailed comparison. The majority (n=735, 76.64%)
245 contained 11 rRNA operons, 189 genomes (19.71%) harbored 10 rRNA operons and only 35
246 genomes (3.65%) contained 9 rRNA operons (Table 1). This ratio is comparable to that found
247 in our initial strain-set tested with FISH and dPCR (11 copies: 82.35%, 10 copies: 11.76% and
248 9 copies: 5.88%). Of these 959 genomes the 16S-BA-allele distributions were: 23.04% had 2,
249 58.39% had 3 and 17.10% had 4 copies (Table 1), respectively. As with the rRNA operon
250 copy-numbers, this distribution correlated with the 16S-BA-allele distribution in our strain-
251 set analyzed by dPCR and WGS (2: 29.41%, 3: 41.17%, 4: 29.41%). Notably, a few strains
252 were calculated to possess 1 (0.31%) or 5 (1.15%) 16S-BA-alleles. The overall diversity of 16S
253 rRNA genotypes (BA alleles/BC alleles) was higher than in our initial strain-set (genotypes
254 4/7, 3/8 and 2/9). Additional major genotypes (frequency >5) obtained from SRA analysis
255 comprised 16S-BA-/BC-allele-ratios of 2/8 and 2/7, minor genotypes were 5/6, 4/6, 4/5, 3/7,
256 3/6, 1/9 and 1/8, each with frequencies < 5.

257 Interestingly, ten of the genomes which were calculated to possess five BA-alleles are
258 from the same originating lab and were sequenced with 100 bp single-end technique only
259 (Supplementary Table S4). Thus, without genomic context it is hardly possible to validate the
260 presence of a fifth 16S-BA-allele from single-end short reads. The same applies to the only
261 other strain (BC038/2000031523) sequenced with 2x100 bp paired-end reads and a mean
262 insert size of 520 bp. Along with three strains putatively containing a single 16S-BA-allele
263 only, strains with five BA-alleles should be re-sequenced using long-read technology for
264 validation.

265 **Table 1: 16S rRNA genotypes obtained from *k*-mer based SRA analysis.** Numbers of 16S-BA-
 266 alleles, overall rRNA operon numbers, and 16S rRNA genotypes resulting from these values
 267 are listed with their respective frequencies.

16S-BA-alleles	# of strains	# of rRNA operon copies per genome	16S rRNA genotype [16S-BA-alleles / BC-alleles]	# of strains
1	3 (0.31%)	9	1 / 8	1 (0.10%)
		10	1 / 9	2 (0.21%)
2	221 (23.04%)	9	2 / 7	28 (2.92%)
		10	2 / 8	116 (12.10%)
		11	2 / 9	77 (8.03%)
3	560 (58.39%)	9	3 / 6	3 (0.31%)
		10	3 / 7	39 (4.10%)
		11	3 / 8	518 (54.01%)
4	164 (17.10%)	9	4 / 5	3 (0.31%)
		10	4 / 6	32 (3.34%)
		11	4 / 7	129 (13.45%)
5	11 (1.15%)	11	5 / 6	11 (1.15%)

268

269 Finally, we tested to which degree 16S rRNA genotypes fit the phylogenetic placement of
 270 strains. For this we correlated established phylogeny of *B. anthracis* based on a number of
 271 canonical SNPs¹² with the distribution of 16S-BA-alleles within ten major canonical SNP
 272 groups of the three branches A, B and C of *B. anthracis*. Figure S2 shows that there is limited
 273 correlation. Notably, B-branch featured a small set of genotypes besides the major 2/8 type.
 274 The few C-branch strains all had the 3/7 genotype. A-branch (comprising the majority of
 275 isolates) was the most diverse, dominantly showing the 2/9 genotype (with the exception of
 276 canSNP group Ames: 4/7). Although the 16S rRNA genotypes did not follow the established

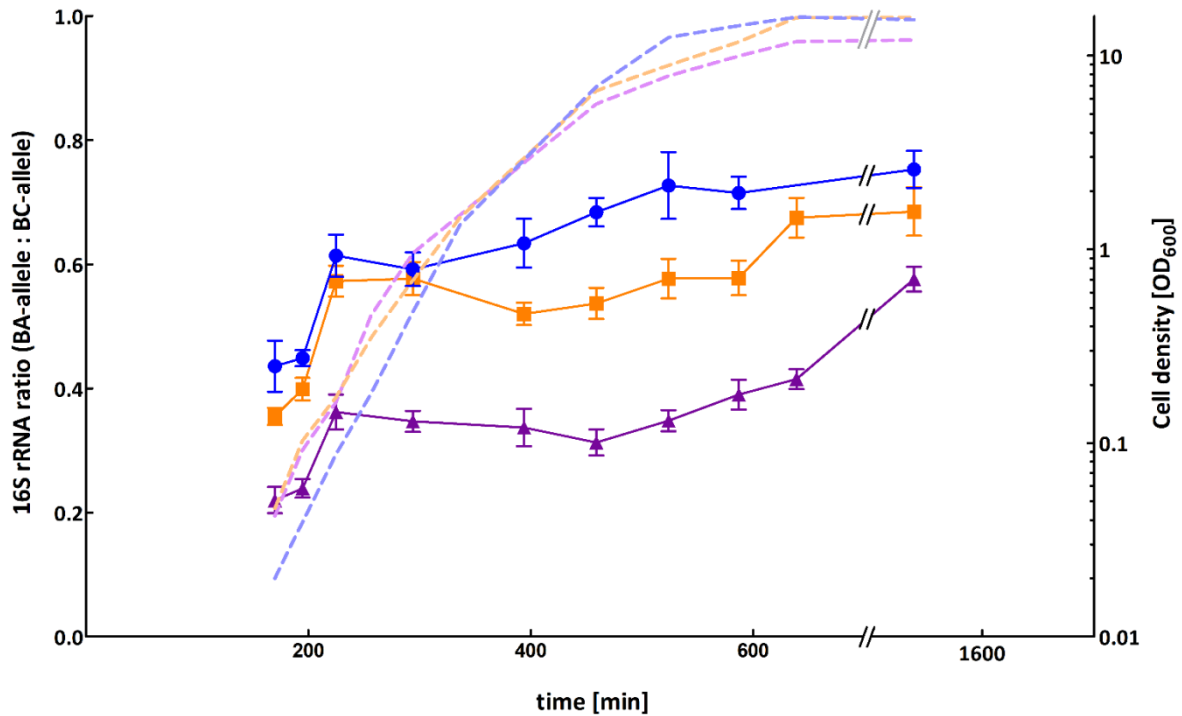
277 phylogeny of *B. anthracis*, the newly developed tools (tetraplex dPCR and *k*-mer based SRA
278 analysis) might still be harnessed as an alternative typing system for *B. anthracis* strains.

279 **Expression of 16S-BA-alleles is proportional to gene copy-number.**

280 Varying ratios of 16S-BA/BC-alleles constitute possible explanations for differences in FISH
281 signals of cells of diverse *B. anthracis* strains (compare Figure 1). Indeed, we found a
282 significant correlation between 16S-BA/BC-allele ratios in sequenced genomes and mean
283 intensities of the Cy3 FISH signals targeting the 16S-BA-allele (tested with the `cor.test`
284 function in R, Pearson's $r=0.61$, $p\text{-value}=0.009$), confirming this assumption.

285 In order to investigate whether the 16S-BA-alleles are differentially expressed
286 throughout different growth phases of *B. anthracis* we quantified 16S rRNA from growth
287 experiments (Figure 4). For this, culture samples of *B. anthracis* strains Sterne, CDC 1014 and
288 Pasteur ATCC 4229 representing three major 16S-BA/BC-allele genotypes 4/7, 3/8 and 2/9,
289 respectively, were taken for total RNA extraction at several time points during lag, log and
290 stationary growth phase. To compare rRNA levels with FISH signals, we also took parallel
291 samples from six of these time points for FISH analysis. By a one-step reverse transcription
292 duplex dPCR, the two 16S allele targets were interrogated for the expression ratios of the
293 16S-BA- *vis-à-vis* the BC-alleles. *B. anthracis* RNA yielded four clusters of droplets in 2D
294 analysis plots, namely 16S-BC-allele⁻/16S-BA-allele⁻, 16S-BC-allele⁺/16S-BA-allele⁻, 16S-BC-
295 allele⁻/16S-BA-allele⁺ and 16S-BC-allele⁺/16S-BA-allele⁺ (Figure 4). RNA of other *B. cereus s. l.*
296 strains produced only two cluster types lacking 16S-BC-allele⁻/16S-BA-allele⁺ and 16S-BC-
297 allele⁺/16S-BA-allele⁺. Absolute quantification of the two initial target concentrations of 16S-
298 BA-alleles/BC-alleles in samples from growth cultures made it possible to determine their

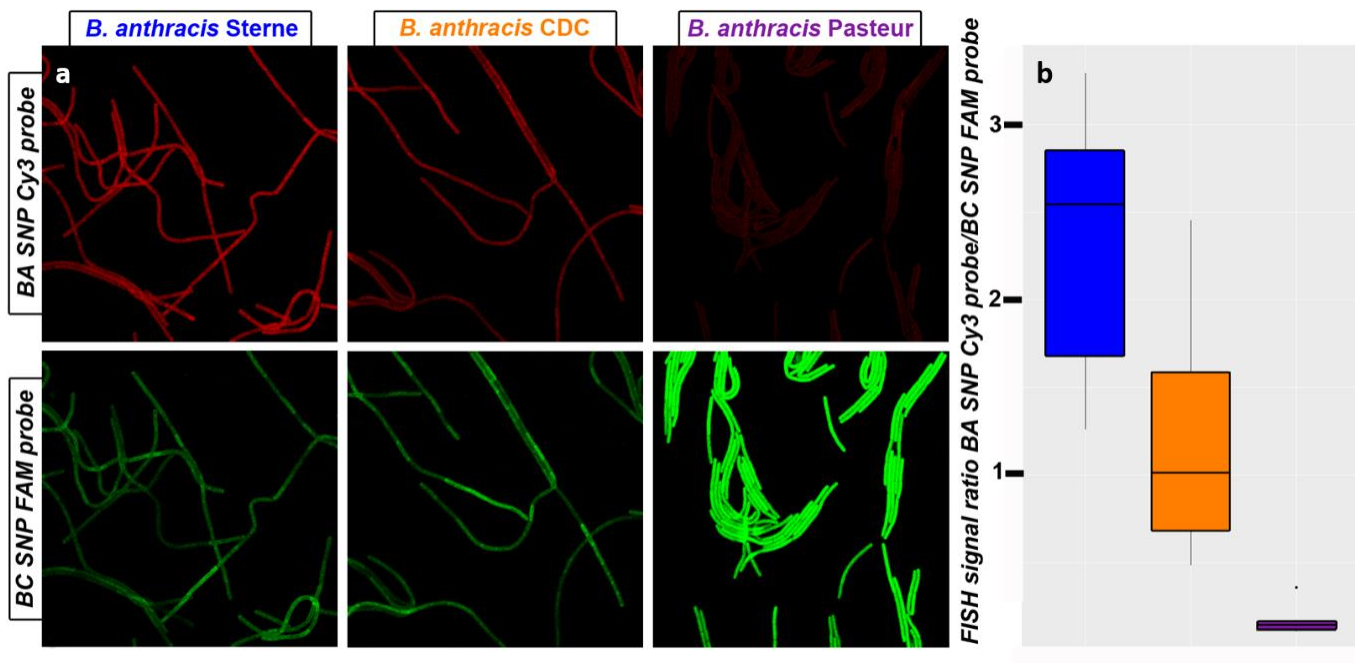
299 ratios representing the expression levels of the 16S-BA-alleles relative to those of 16S-BC-
300 alleles (Figure 4). Notably, 16S-BA/BC-allele rRNA ratios varied during growth and showed
301 similar expression patterns in all three tested *B. anthracis* strains. Starting from a relatively
302 low 16S-BA/BC-allele ratio in early log phase, the fraction of 16S-BA-allele expression
303 increased in early log-phase and decreased in mid log-phase with a final increase towards
304 the stationary phase. While shifts in 16S-BA/BC-allele expression patterns in these strains
305 were similar, differences were observed in numerical expression ratios. *B. anthracis* Sterne
306 showed the highest 16S-BA/BC-allele expression ratio ranging from 0.44 (early exponential
307 phase) up to 0.75 (stationary phase), compared to CDC 1014 with 0.36 to 0.69 and Pasteur
308 ATCC 4229 with 0.22 to 0.58, which was found to have the lowest 16S-BA-allele expression in
309 all growth phases. The largest differences in expression levels between all strains were
310 observed in late log phase (Figure 4). The observed diverging levels of 16S-BA-allele
311 expression in the three tested strains can easily be explained by the different numbers of
312 16S-BA-allele copies per genome (2, 3 or 4). Nevertheless, the proportion of 16S-BA-allele
313 rRNA in late-exponential *B. anthracis* cells is quite disproportionate. If all rRNA operons were
314 transcribed at a constant and equal rate, one would expect a ratio of 0.22 (Pasteur 2/9), 0.38
315 (CDC 3/8), and 0.57 (Sterne 4/7). Instead, we measured ratios, which correlate to a 1.57-
316 (Pasteur), 1.46-(CDC), and 1.09-(Sterne) fold 16S-BA-allele over-representation on average
317 throughout all growth phases and up to 2.59 - (Pasteur), 1.83-(CDC), and 1.32- (Sterne) fold
318 in stationary phase.



319

320 **Figure 4. Expression ratios of 16S-BA- and –BC-alleles in three different *B. anthracis* strains**
321 **at different growth phases.** Expression level ratios of 16S-BA-alleles relative to 16S-BC-alleles
322 were calculated from absolute target concentrations obtained by RT-dPCR. Values were
323 plotted against time-points of each sample taken during growth from early exponential to
324 stationary phase for *B. anthracis* Sterne (blue), CDC 1014 (orange) and Pasteur ATCC 4229
325 (purple) representing three major 16S rRNA genotypes (BA/BC) 4/7, 3/8 and 2/9,
326 respectively. Error bars indicate the Poisson 95 % confidence intervals for each copy-number
327 ratio. Dotted lines depict cell densities over time.

328 The shift towards elevated expression of the 16S-BA-allele genes over time was not
329 significantly reflected in FISH signal intensities, possibly due to the general decrease of FISH
330 signals over time. However, if cells were sampled and fixed at identical time points, 16S-
331 BA/BC-allele ratios were always highest for *B. anthracis* Sterne and lowest for *B. anthracis*
332 Pasteur, which reflects their 16S-BA/BC-allele ratios on the genomic and transcript levels
333 (Figure 5). Also, sampled across all time points, 16S-BA/BC-allele FISH signal ratios correlated
334 well with allele distributions in the three different strains (ANOVA in R, $p=0.0002$, Figure 5).



335 **Figure 5. FISH of *B. anthracis* strains harboring different numbers of 16S-BA-alleles.**
336 (a) Representative FISH images showing signal intensities of *B. anthracis* strains with
337 diverging genomic 16S (BA/BC) allele profiles (Sterne 4/7), CDC 1014 (3/8) and Pasteur ATCC
338 4229 (2/9). Samples were taken and processed after 460 min of continuous growth.
339 (b) Boxplot of BA_SNP_Cy3 and BC_SNP_FAM FISH signal ratios across all sampled time
340 points for *B. anthracis* Sterne (blue), CDC 1014 (orange) and Pasteur ATCC 4229 (purple).

341 Discussion

342 Using a combination of newly developed *in situ*, *in vitro* and *in silico* approaches, we
343 unraveled the elusive heterogeneity of 16S rRNA genes in the biothreat agent *B. anthracis*.
344 Results consistently delineate the organism's intragenomic diversity of 16S rRNA genes, their
345 differential expression across growth phases and their intergenomic heterogeneity in publicly
346 available and newly sequenced genomes. Intragenomic micro-diversity within 16S rRNA
347 genes has long been known from other species^{13,14} and was found to increase with higher
348 copy-numbers of rRNA operons¹⁵. Thus, the species-wide intra- and inter-genomic micro-
349 diversity related to SNP 1110 in the 9 to 11 copies of the 16S rRNA gene of *B. anthracis* is not
350 totally unsurprising^{3,4}. Whereas some of such polymorphic sites are associated with a
351 distinct phenotypic trait (e.g. stress resistance)^{16,17} the functional assignment for the majority
352 of these sequence variations (including those in *B. anthracis* 16S rRNA genes) remains
353 elusive.

354 Though discovered before using Sanger sequencing⁴, the specific SNP in the 16S
355 rRNA genes of *B. anthracis*, was disregarded despite the availability of numerous published
356 genomes. Generally, sequence variations in multi-copy genes such as 16S rRNA genes can
357 hardly be detected when relying on conventional short-read WGS and subsequent reference
358 mapping¹⁸, which was used to generate the majority of publicly available *B. anthracis* whole
359 genomes sequences. SNP calling in different rRNA operons or other paralogous genes gives
360 ambiguous results since assemblers tend to interpret low frequent sequence variations as
361 sequencing errors and correct them prior assembly¹⁹. Even if detected, distances of the SNP
362 to unique flanking regions up- and downstream of the multi-copy gene may be >1000 bases

363 and thus, larger than typical library fragment sizes of 500–800 bases. In such cases,
364 chromosomal locations of SNPs cannot be reconstructed. Instead, all rRNA gene related
365 reads are assembled into one contig with diverse fringes²⁰. The average read length of
366 Nanopore sequencing is typically larger than 5 kb and can therefore cover complete rRNA
367 operons. Thus, any unique SNP occurring in a single or a few rRNA gene alleles can be
368 precisely allocated to a specific chromosome position, especially when combined with short-
369 read sequencing and hybrid assembly as used here. Therefore, the challenges described
370 above will become rather minor for future genomic analysis of *B. anthracis*. Such work is
371 facilitated by the additional 33 complete high-quality genomes we have contributed here.
372 These genomes cover all three major phylogenetic lineages (canSNP groups), all *bona fide*
373 16S-BA-allele frequencies (2, 3, and 4) as well as all known rRNA operon copy numbers.

374 On the *B. anthracis* chromosome, the 16S rRNA operons *rrnE*, *F*, *G* and *H*, are located
375 in close proximity to each other with only 15.8, 8.5 and 5.2 kb in-between, respectively
376 (forming a genomic region with a high density of four 16S rRNA operons within less than
377 50 kb). Conversely, the other 16S rRNA genes are rather dispersed with distances greater
378 than 50 kb in-between. The 16S-BA-allele is present in operons *rrnC* and *rrnE* in all strains
379 analyzed with long-read WGS while *rrnF* and *rrnG* seem to be variable. Since the four
380 operons *rrnE*, *F*, *G* and *H* are relatively close to each other in the *B. anthracis* chromosome,
381 homologous recombination and gene duplications might be the reason for this allelic
382 variation. Also, compared to all other 16S rRNA alleles on the *B. anthracis* genome, the 16S
383 rRNA copies in this region (*rrnE* – *rrnH*) seem to differ from each other only in SNP position
384 1110. This finding promotes the explanation that the 16S rRNA copies in this region of high
385 rRNA operon density are subjected to an increased recombination-rate between alleles with

386 and without *B. anthracis* specific SNP 1110. This notion is also supported by the fact that
387 only operons *rrnG* and *H* seem to be affected by deletion events in all strains analyzed by
388 long-read WGS. The alternative explanation, horizontal gene transfer of a divergent allele,
389 seems unlikely. We were unable to identify any 16S rRNA gene in public databases matching
390 the 16S-BA-allele outside *B. anthracis*.

391 Recombination and deletion events in 16S rRNA operons of *B. anthracis* do occur as
392 evidenced by a study on bacitracin resistance. Two deletion events, DelFG and DelGH, were
393 described which caused elimination of gene-clusters between rRNA operons *rrnF* and *G* and
394 *G* and *H*, respectively¹⁸. These DelFG- and DelGH-events describe a possible origin of *B.*
395 *anthracis* strains with ten 16S rRNA gene copies, i.e. 21% of all strains (Table 1). Random
396 gene duplication and gene elimination by recombination might also explain another
397 observation: the newly defined 16S rRNA genotypes did not convincingly reflect the
398 established *B. anthracis* phylogeny (Supplementary Figure S1). Instead, some 16S rRNA
399 genotypes seem to be dominant yet not exclusive in separate branches, e.g. 2/8 copies in B-
400 branch or 3/8 in A-branch (Supplementary Figure S1).

401 The recognition of intra- and intergenomic 16S rRNA allele diversity in *B. anthracis*
402 opens possibilities to harness unique SNPs in 16S rRNA gene alleles and their transcripts. This
403 finding strongly highlights the great potential of such genomic variations for both
404 identification of *B. anthracis* and for diagnostics of anthrax disease. This approach is probably
405 also applicable to other pathogens which are otherwise difficult to discriminate from their
406 less notorious relatives.

407 **Materials and Methods**

408 **Cultivation of bacteria**

409 The cultivation of the virulent *B. anthracis* strains was performed in a biosafety level 3
410 laboratory (BSL3). All *Bacillus* strains were cultivated overnight on Columbia blood agar
411 plates (containing 5 % sheep blood, Becton Dickinson, Heidelberg, Germany) at 37°C.
412 For isolation of DNA, a 1 µl loop of colonies was transferred to a 2 mL screw cap microfuge
413 tube, inactivated with 2 % Terralin PAA (Schülke&Mayr GmbH, Norderstedt, Germany) for 30
414 min and washed three times with phosphate-buffered saline (PBS) as described previously ²¹.

415 For FISH, 50 ml centrifuge tubes containing 5 ml of Tryptic Soy Broth (TSB, Merck
416 KGaA, Darmstadt, Germany) were inoculated with one colony from an overnight culture (see
417 above) and incubated at 37° C with shaking at 150 rpm. After four hours of growth, bacteria
418 were pelleted by centrifugation at 5,000 x g for 10 min, washed with PBS, and fixated with 3
419 ml 4% (v/v) formaldehyde for one hour at ambient temperature. After fixation, cells were
420 washed three times with PBS, resuspended in a 1:1 mixture of absolute ethanol and PBS and
421 stored at -20 °C until further use. To ensure sterility 1/10 of the inactivated material was
422 incubated in thioglycolate-medium (Merck KGaA, Darmstadt, Germany) for seven days
423 without growth before material was taken out of the BSL3 laboratory.

424 For growth phase analysis, 1 ml of overnight cultures of attenuated *B. anthracis*
425 (Sterne, CDC 1014 and ATCC 4229 Pasteur) in TSB was used to inoculate 100 ml of fresh TSB
426 in 1 L baffled flasks and incubated at 37°C with shaking at 100 rpm. Every 30 min, turbidity
427 was measured as OD₆₀₀ and 1 ml samples were taken for FISH and RNA isolation,
428 respectively. After pelleting by centrifugation samples for RNA isolation were resuspended

429 and inactivated using 2% Terralin PAA for 30 min and washed three times with PBS. FISH
430 samples were treated as described above.

431 **Design of Primers and Probes**

432 Primers and probes were designed using Geneious 10.1.3 (Biomatters, Auckland, New
433 Zealand) and numerous probe variations were tested to identify the best combination and
434 number of locked nucleic acids for differentiation of *B. anthracis* and the other *B. cereus s. l.*
435 group species based on the SNP (pos. 1110) detected previously⁴. The final probes for FISH
436 included two and one locked nucleic acid while dPCR probes contained 5 and 6 for the *B.*
437 *anthracis* (BA) and the *B. cereus s-l-* (BC) probe, respectively (Supplementary Table S2). For
438 sequences of positive (EUB338²²) and negative (nonEUB,²³) control probes for FISH, see
439 Supplementary Table S2. Primers as well as probes labeled with 6-carboxyfluorescein (6-
440 FAM), hexachlorofluorescein (HEX), indocarbocyanine (Cy3) or indodicarbocyanine (Cy5)
441 were purchased commercially (TIB Molbiol, Berlin, Germany).

442 To determine the ideal formamide concentration for the FISH hybridization buffer, the
443 fluorescence signals of probe BA_SNP_Cy3 and probe BC_SNP_FAM were assessed with
444 *B. anthracis* Sterne and *B. cereus* ATCC 10987 at different formamide concentrations (0, 10,
445 20, 25, 30, 35, 40, 45, 50% FA concentration in the hybridization buffer) as described
446 elsewhere²⁴. Hybridization at 30% formamide was determined to be ideal for differentiation
447 of *B. anthracis* and *B. cereus s. l.* group (Supplementary Figure S3).

448 **Fluorescence in situ Hybridization and Image Processing**

449 FISH was carried out as described elsewhere²⁴. A positive-control probe targeting eubacteria
450 (EUB338,²²) and a nonsense probe targeting no known bacterial species (nonEUB,²³) as a
451 control for unspecific probe binding were included in each hybridization experiment. Briefly,
452 2 µl of fixed cells were spotted on teflon coated slides (Marienfeld, Lauda-Königshofen,
453 Germany) and dried at 46°C. Then, cells were permeabilized using 10 ml of 15mg/ml
454 lysozyme (Merck KGaA, Darmstadt, Germany, Cat.Nr. 62970) per well at 46°C for 12 min.
455 After dehydration in an ascending ethanol series (50, 80, 96% (v/v) ethanol) cells were
456 covered with 10 µl hybridization buffer (0.9 M NaCl, 20 mM Tris-HCl (pH 8.0), 0.01% SDS,
457 30% formamide) with probes at a concentration of 10 µM and incubated in a humid
458 chamber in the dark at 46°C for 1.5 h. Slides were washed in 50 ml pre-warmed washing
459 buffer (0.1 M NaCl, 20 mM Tris-HCl (pH 8.0), 5 mM EDTA (pH 8.0)), for 10 min at 48°C in a
460 water bath. Finally, slides were dipped in ice-cold ddH₂O and carefully dried with
461 compressed air. For each strain, FISH was performed in duplicate and two pictures were
462 taken per well, so that the resulting fluorescence intensity was the mean of four images. To
463 increase accuracy in the growth curve assay, five pictures were taken per well, so that the
464 resulting fluorescence intensity was the mean of ten images. All images were recorded with
465 a confocal laser-scanning microscope (LSM 710, Zeiss, Jena, Germany). Excitations for FAM,
466 Cy3 and Cy5 were at 490, 560 and 630 nm respectively. Emission was measured within the
467 following ranges: FAM: 493-552 nm, Cy3: 561-630 nm and Cy5: 638-724 nm. Images were
468 processed with Daime²⁵, using the area of the EUB signal as a mask to measure average
469 fluorescence intensity for BA_SNP_Cy3 and BC_SNP_FAM: The EUB images were segmented

470 and unspecific fluorescence excluded with default threshold settings and this object layer
471 was transferred to BA_SNP_Cy3 and BC_SNP_FAM images.

472 **Isolation of nucleic acids**

473 DNA isolation from inactivated cells was carried out using MasterPure™ Gram Positive DNA
474 Purification Kit (Lucigen, Middleton, WI, USA) according to the manufacturer's protocol. DNA
475 samples were quantified using the Qubit dsDNA HS Assay Kit protocol (Thermo Scientific,
476 Dreieich, Germany). For RNA isolation from inactivated cells, RNeasy Protect Bacteria Mini
477 Kit (Qiagen, Hilden, Germany) was used according to the supplier's protocol for enzymatic
478 lysis and proteinase K digestion of bacteria. In order to eliminate residual DNA, RNA samples
479 were purified twice using RNeasy MinElute Cleanup Kit (Qiagen, Hilden, Germany) and
480 quantified using the Qubit RNA HS Assay Kit protocol (Thermo Scientific, Dreieich, Germany).
481 The absence of DNA in the final RNA preparation was verified by conducting PCR on marker
482 *dhp61*²⁶ with negative results.

483 **Tetraplex droplet digital PCR assay for quantification of 16S rRNA gene alleles**

484 Digital PCR (dPCR) allows for absolute quantification of DNA or RNA template
485 concentrations²⁷. For 16S rRNA gene analysis the 20 µl dPCR pre-droplet mix consisted of 10
486 µl dPCR Supermix for Probes (Bio-Rad Laboratories, Munich, Germany), 1 µl 20x 16S SNP
487 Primer mix (final concentrations 900 nM), 0.6 µl of 20x mix of 16S SNP BC probe (final
488 concentration 150 nM), 0.6 µl of 20x mix of 16SSNP BA probe (final concentration 150 nM),
489 0.9 µl of 20x PL3 primer/probe mix (final concentrations: probe 225 nM, primers 810 nM),
490 1.5 µl of GyrA 20x primer/probe mix (final concentrations: probe 375 nM, primers 1350 nM),

491 4.4 µl of nuclease free water (Qiagen, Hilden, Germany) and 1 µl of template DNA freshly
492 diluted to a concentration of 0.05 ng/µl. To ensure independent segregation of the 16S rRNA
493 gene copies from the bacterial chromosome and the reference genes into droplets, template
494 DNA was digested (no cut sites within 16S rRNA genes) prior to dPCR by BsiWI-HF, BsrGI-HF
495 and HindIII-HF (New England Biolabs GmbH, Frankfurt am Main, Germany) in 1 x Cutsmart
496 buffer (New England Biolabs GmbH, Frankfurt am Main, Germany) for 60 min and then the
497 enzymes heat inactivated at 80°C for 20 min according to the manufacturer's protocol.
498 Partitioning of the reaction mixture into up to 20,000 individual droplets was achieved using
499 a QX200 dPCR droplet generator (Bio-Rad Laboratories, Munich, Germany). A two-step PCR-
500 reaction was performed on a Mastercycler Pro instrument (Eppendorf, Wesseling-Berzdorf,
501 Germany) with the following settings: one DNA-polymerase activation step at 95°C for 10
502 min was followed by 40 cycles of denaturation at 94°C for 30 s and annealing/extension at
503 58°C for 1 min. Finally enzyme inactivation was performed at 98°C for 10 min before
504 the samples were cooled down and held at 4°C. All steps were carried out with a
505 temperature ramp rate of 2°C/s. After completion droplets were analyzed using the QX100
506 Droplet Reader (Bio-Rad) and absolute concentrations for each target were quantified using
507 Poisson statistics as implemented in the Quantasoft Pro Software (Bio-Rad Laboratories,
508 Munich, Germany).

509 Then, the absolute concentrations of *PL3* and *gyrA* were compared. To ensure assay
510 integrity samples with a deviation range greater than 10% within the two markers were
511 excluded and had to be repeated. If deviation was below 10% both targets were set as a
512 reference with a copy-number of one. The software then automatically takes the mean
513 concentration of both references to calculate the copy-numbers of BC and BA alleles.

514 According to the recommendations provided by²⁸, all samples with copy-numbers between
515 0.35 and 0.65 deviating from an integer number or with a confidence interval greater than 1
516 were excluded from analysis and were repeated. All valid runs were rounded to the next
517 integer number.

518 **Duplex one-step reverse transcription dPCR to compare expression levels of**

519 **16S BC- and 16S-BA-allele**

520 The 20 µl RT-dPCR reaction mixture consisted of 5 µl One-Step RT-dPCR Advanced Supermix
521 for Probes (Bio-Rad, Laboratories, Munich, Germany), 2 µl of Reverse Transcriptase (Bio-Rad,
522 final concentration 20 U/µl), 0.6 µl of DTT (Bio-Rad, Laboratories, Munich, Germany; final
523 concentration 10 nM), 1.5 µl 20x 16S SNP Primer mix (final concentrations 1350 nM), 1.5 µl
524 of 20x mix of 16S SNP BC probe (final concentration 375 nM), 1.5 µl of 20x mix of 16S SNP BA
525 probe (final concentration 375 nM), 6.9 µl of nuclease free water (Qiagen, Hilden, Germany)
526 and 1 µl of template RNA. Reverse transcription was achieved within droplets prior to dPCR.
527 Partitioning of the reaction mixture into up to 20,000 droplets was carried out using a QX200
528 dPCR droplet generator (Bio-Rad, Laboratories) and PCR was performed on the Mastercycler
529 Pro (Eppendorf, Wesseling-Berzdorf, Germany) with the following settings: The initial
530 reverse transcription step was performed at 48°C for 60 min. An enzyme activation step at
531 95°C was carried out for 10 min followed by 40 cycles of a two-step program of denaturation
532 at 94°C for 30 s and annealing/extension at 58°C for 1 min. Final enzyme inactivation was
533 performed at 98°C for 10 min before the samples were cooled down and held at 4°C. All
534 steps were carried out with a temperature ramp rate of 2°C/s. After completion, droplets

535 were analyzed using the QX100 Droplet Reader (Bio-Rad, Laboratories, Munich, Germany)

536 and results were quantified with the Quantasoft Pro Software (Bio-Rad, Laboratories).

537

538 **Library preparation, sequencing and assembly of genomes**

539 The libraries for the Illumina sequencing were prepared using the NEBNext® Ultra™ II FS DNA

540 Library Prep Kit for Illumina (New England BioLabs GmbH, Frankfurt am Main, Germany)

541 according to the protocol for large fragment sizes >550 bp but with a minimal fragmentation

542 time of only 30 s. Afterwards, libraries were pooled equimolarly and sequenced on an

543 Illumina MiSeq device (Illumina Inc., San Diego, CA, U.S.A.) using the MiSeq Reagent Kit v3

544 (2×300 bp).

545 The libraries for the nanopore sequencing were prepared using the Ligation

546 Sequencing Kit SQK-LSK109 (Oxford Nanopore Technologies, Oxford, U.K.) combined with

547 the Native Barcoding Expansion EXP-NBD104 and sequenced as one pool on a MinION

548 flowcell FLO-MIN106D (Type R9.4.1; Oxford Nanopore Technologies, Oxford, U.K.) for 48 h.

549 Basecalling and demultiplexing was done separately using Guppy v3.2.10 (Oxford Nanopore

550 Technologies, Oxford, U.K.) with the high accuracy basecalling model. Quality (≥ 10 q) and

551 length ($\geq 1,000$ bp) filtering was done using Filtlong version 0.2.0

552 (<https://github.com/rrwick/Filtlong>).

553 Hybrid assemblies were constructed in two stages. First, nanopore reads were assembled

554 using Flye version 2.7²⁹ with default parameters and two iterations of polishing. Second,

555 Illumina reads were assembled together with the nanopore raw reads and the nanopore

556 assembly as trusted contigs using SPAdes version 3.14³⁰ with parameters “-k

557 55,77,99,113,127 --careful". Afterwards, the assembled contigs were reverse
558 complemented, if necessary and rotated to the same start sequence as strain Ames
559 Ancestor. Finally, the contigs were polished once more using Pilon version 1.23³¹.

560 **Bioinformatics analyses**

561 For the long-read assemblies, ribosomal operons were annotated using barrnap version 0.9
562 (<https://github.com/tseemann/barrnap>). SNP alleles were searched using USEARCH
563 version 11³² and the 16S SNP BA/BC probe sequences (see Supplementary Table S2) as an
564 oligo sequence database. To investigate the frequency and distribution of the alleles of 16S
565 rRNA genes in the *B. anthracis* species comprehensively, we downloaded all available short-
566 read Illumina data sets (at the end of 2020) from the NCBI Sequence Read Archive³³. These
567 data sets were then assembled using SPAdes v1.14³⁰ with parameters "-k 55,77,99,113,127 -
568 -careful". The contigs of the resulting assemblies were extended using tadpole from the
569 BBTools package³⁴ and with parameters "el=1000 er=1000 mode=extend". Afterwards,
570 blastn³⁵ with parameters "-evalue 1e-10 -word_size 9" was used to align the 23S rRNA
571 sequence against each extended contig end. For each assembly, the number of contigs
572 ending with a 23S rRNA fragment were counted and CanSNPer³⁶ was used to determine the
573 canonical SNPs and likely position in the CanSNP tree. In a next step, kmercountexact from
574 the BBTools package was used with the parameters "fastadump=f mincount=2 k=16" to
575 count all *k*-mers of size 16 from error-corrected reads. From these *k*-mers, the frequencies of
576 the two allelic *k*-mers (sequences of 16 nt used for the dPCR probes) was extracted.
577 kmercountexact also reports a *k*-mer-based coverage estimation of the sequenced reads
578 which is used to filter the assemblies by coverage (min. 20X), number of contigs (max. 200),

579 number of potential rRNAs (>8) and success of CanSNPer prediction. For each remaining
580 assembly, the number of rRNAs carrying the SNP of 16S BA allele was estimated by
581 determine the ratio of the allelic *k*-mers multiplied with the total number of rRNAs, rounded
582 to a whole number. To validate this estimation, we applied the same algorithm to every
583 assembly where both short- and long-reads and/or dPCR results were available and
584 compared the estimated number of BA alleles to the counted number in the long-read
585 assembly or to the measured number from the dPCR experiments. They were consistent
586 across different sequencing coverage, total number of rRNA operons and known BA allele
587 frequencies.

588 **Data availability**

589 All genomic data generated or analyzed prior or during this study can be accessed via the
590 NCBI BioProject PRJNA695105. Individual accession numbers are listed in Supplementary
591 Tables S3 and S4.

592 **References**

- 593 1 Takahashi, H. *et al.* *Bacillus anthracis* incident, Kameido, Tokyo, 1993. *Emerg Infect Dis*
594 **10**, 117-120, doi:10.3201/eid1001.030238 (2004).
- 595 2 Turnbull, P. C. (ed Peter Turnbull (Editor)) (WHO Press, Geneva (CH), 2008).
- 596 3 Candelon, B., Guilloux, K., Ehrlich, S. D. & Sorokin, A. Two distinct types of rRNA
597 operons in the *Bacillus cereus* group. *Microbiology* **150**, 601-611,
598 doi:10.1099/mic.0.26870-0 (2004).
- 599 4 Hakovirta, J. R., Prezioso, S., Hodge, D., Pillai, S. P. & Weigel, L. M. Identification and
600 analysis of informative single nucleotide polymorphisms in 16S rRNA gene sequences
601 of the *Bacillus cereus* group. *J Clin Microbiol*, doi:JCM.01267-16 [pii];
602 10.1128/JCM.01267-16 (2016).
- 603 5 Weerasekara, M. L. *et al.* Double-color fluorescence in situ hybridization (FISH) for the
604 detection of *Bacillus anthracis* spores in environmental samples with a novel
605 permeabilization protocol. *J Microbiol Methods* **93**, 177-184,
606 doi:10.1016/j.mimet.2013.03.007, S0167-7012(13)00101-2 [pii] (2013).

- 607 6 Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data
608 processing and web-based tools. *Nucleic Acids Res* **41**, D590-596,
609 doi:10.1093/nar/gks1219 (2013).
- 610 7 Koshkin, A. A. *et al.* LNA (Locked Nucleic Acids): Synthesis of the adenine, cytosine,
611 guanine, 5-methylcytosine, thymine and uracil bicyclonucleoside monomers,
612 oligomerisation, and unprecedented nucleic acid recognition. *Tetrahedron* **54**, 3607-
613 3630, doi:[https://doi.org/10.1016/S0040-4020\(98\)00094-5](https://doi.org/10.1016/S0040-4020(98)00094-5) (1998).
- 614 8 Ravel, J. *et al.* The complete genome sequence of *Bacillus anthracis* Ames "Ancestor".
615 *Journal of Bacteriology* **191**, 445-446, doi:10.1128/jb.01347-08 (2008).
- 616 9 Johnson, S. L. *et al.* Complete genome sequences for 35 biothreat assay-relevant
617 bacillus species. *Genome Announc* **3**, e01501-01515, doi:10.1128/genomeA.00151-15
618 (2015).
- 619 10 Rasko, D. A., Altherr, M. R., Han, C. S. & Ravel, J. Genomics of the *Bacillus cereus*
620 group of organisms. *FEMS Microbiol Rev* **29**, 303-329, doi:S0168-6445(05)00003-3
621 [pii]; 10.1016/j.femsre.2004.12.005 (2005).
- 622 11 Ellerbrok, H. *et al.* Rapid and sensitive identification of pathogenic and apathogenic
623 *Bacillus anthracis* by real-time PCR. *FEMS Microbiol Lett* **214**, 51-59,
624 doi:S0378109702008376 [pii] (2002).
- 625 12 Van Ert, M. N. *et al.* Global genetic population structure of *Bacillus anthracis*. *PLoS*
626 *One* **2**, e461, doi:10.1371/journal.pone.0000461 (2007).
- 627 13 Cilia, V., Lafay, B. & Christen, R. Sequence heterogeneities among 16S ribosomal RNA
628 sequences, and their effect on phylogenetic analyses at the species level. *Mol Biol*
629 *Evol* **13**, 451-461, doi:10.1093/oxfordjournals.molbev.a025606 (1996).
- 630 14 Acinas, S. G., Marcelino, L. A., Klepac-Ceraj, V. & Polz, M. F. Divergence and
631 redundancy of 16S rRNA sequences in genomes with multiple *rrn* operons. *J Bacteriol*
632 **186**, 2629-2635, doi:10.1128/jb.186.9.2629-2635.2004 (2004).
- 633 15 Větrovský, T. & Baldrian, P. The variability of the 16S rRNA gene in bacterial genomes
634 and its consequences for bacterial community analyses. *PLoS ONE* **8**, e57923-e57923,
635 doi:10.1371/journal.pone.0057923 (2013).
- 636 16 Werner, G. *et al.* Detection of mutations conferring resistance to linezolid in
637 *Enterococcus* spp. by fluorescence in situ hybridization. *Journal of Clinical*
638 *Microbiology* **45**, 3421-3423, doi:10.1128/jcm.00179-07 (2007).
- 639 17 Kurylo, C. M. *et al.* Endogenous rRNA sequence variation can regulate stress response
640 gene expression and phenotype. *Cell reports* **25**, 236-248.e236,
641 doi:10.1016/j.celrep.2018.08.093 (2018).
- 642 18 Furuta, Y. *et al.* Loss of bacitracin resistance due to a large genomic deletion among
643 *Bacillus anthracis* strains. *mSystems* **3**, doi:10.1128/mSystems.00182-18 (2018).
- 644 19 Heydari, M., Miclotte, G., Demeester, P., Van de Peer, Y. & Fostier, J. Evaluation of the
645 impact of Illumina error correction tools on de novo genome assembly. *BMC*
646 *Bioinformatics* **18**, 374, doi:10.1186/s12859-017-1784-8 (2017).
- 647 20 Salzberg, S. L. *et al.* GAGE: A critical evaluation of genome assemblies and assembly
648 algorithms. *Genome Res* **22**, 557-567, doi:10.1101/gr.131383.111 (2012).
- 649 21 Braun, P. *et al.* Microevolution of anthrax from a young ancestor (M.A.Y.A.) suggests a
650 soil-borne life cycle of *Bacillus anthracis*. *PLoS ONE* **10**, e0135346,
651 doi:10.1371/journal.pone.0135346 (2015).

- 652 22 Stahl, D. A., Flesher, B., Mansfield, H. R. & Montgomery, L. Use of phylogenetically
653 based hybridization probes for studies of ruminal microbial ecology. *Appl Environ*
654 *Microbiol* **54**, 1079-1084 (1988).
- 655 23 Manz, W., Amann, R., Ludwig, W., Wagner, M. & Schleifer, K.-H. Phylogenetic
656 oligodeoxynucleotide probes for the major subclasses of proteobacteria: Problems
657 and solutions. *Systematic and Applied Microbiology* **15**, 593-600,
658 doi:[https://doi.org/10.1016/S0723-2020\(11\)80121-9](https://doi.org/10.1016/S0723-2020(11)80121-9) (1992).
- 659 24 Daims, H., Stoecker, K. & Wagner, M. in *Molecular microbial ecology* (eds Mark
660 Osborn & Cindy Smith) 213-239 (Taylor & Francis, 2005).
- 661 25 Daims, H., Lückner, S. & Wagner, M. Daime, a novel image analysis program for
662 microbial ecology and biofilm research. *Environ Microbiol* **8**, 200-213,
663 doi:10.1111/j.1462-2920.2005.00880.x (2006).
- 664 26 Antwerpen, M. H., Zimmermann, P., Bewley, K., Frangoulidis, D. & Meyer, H. Real-time
665 PCR system targeting a chromosomal marker specific for *Bacillus anthracis*. *Mol Cell*
666 *Probes* **22**, 313-315, doi:S0890-8508(08)00041-8 [pii]; 10.1016/j.mcp.2008.06.001
667 (2008).
- 668 27 Kuypers, J. & Jerome, K. R. Applications of Digital PCR for Clinical Microbiology. *J Clin*
669 *Microbiol* **55**, 1621-1628, doi:10.1128/jcm.00211-17 (2017).
- 670 28 Bell, A. D., Usher, C. L. & McCarroll, S. A. Analyzing copy number variation with
671 droplet digital PCR. *Methods Mol Biol* **1768**, 143-160, doi:10.1007/978-1-4939-7778-
672 9_9 (2018).
- 673 29 Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of long, error-prone reads
674 using repeat graphs. *Nat Biotechnol* **37**, 540-546, doi:10.1038/s41587-019-0072-8
675 (2019).
- 676 30 Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to
677 single-cell sequencing. *J Comput Biol* **19**, 455-477, doi:10.1089/cmb.2012.0021
678 (2012).
- 679 31 Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant
680 detection and genome assembly improvement. *PLoS ONE* **9**, e112963,
681 doi:10.1371/journal.pone.0112963 (2014).
- 682 32 Edgar, R. C. Search and clustering orders of magnitude faster than BLAST.
683 *Bioinformatics* **26**, 2460-2461, doi:10.1093/bioinformatics/btq461 (2010).
- 684 33 Leinonen, R., Sugawara, H., Shumway, M. & International Nucleotide Sequence
685 Database, C. The sequence read archive. *Nucleic acids research* **39**, D19-D21,
686 doi:10.1093/nar/gkq1019 (2011).
- 687 34 Bushnell, B., Rood, J. & Singer, E. BBMerge - Accurate paired shotgun read merging via
688 overlap. *PLoS ONE* **12**, e0185056, doi:10.1371/journal.pone.0185056 (2017).
- 689 35 Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein
690 database search programs. *Nucleic Acids Res* **25**, 3389-3402,
691 doi:10.1093/nar/25.17.3389 (1997).
- 692 36 Lärkeryd, A. *et al.* CanSNPer: a hierarchical genotype classifier of clonal pathogens.
693 *Bioinformatics* **30**, 1762-1764, doi:10.1093/bioinformatics/btu113 (2014).
- 694
- 695

696 **Acknowledgements**

697 We thank Linda Dobrzykowski and Josua Zinner for technical assistance. For providing
698 bacterial strains we are grateful to Fabian Leendertz, Silke Klee and Roland Grunow (RKI),
699 Wolfgang Beyer (University of Hohenheim) and Paul Keim (Northern Arizona University). We
700 also thank Olfert Landt and his team at TIB Molbiol (Berlin) for technical support in designing
701 LNA-based probes.

702 This study was supported by funds from the German Federal Ministry of Defense
703 (Sonderforschungsprojekt 36Z1-S-431618 and STAN 48-2009-23).

704 The research described herein is part of the Medical Biological Defense Research Program of
705 the Bundeswehr Joint Medical Service. Opinions, interpretations, conclusions, and
706 recommendations are those of the authors and are not necessarily endorsed by any
707 governmental agency, department or other institutions.

708 **Contributions**

709 P.B., F.Z., G.G. and K.S. designed the study and interpreted the results. M.W. contributed the
710 bioinformatics analysis. I.S., F.Z., K.A. and S.M. performed FISH experiments. P.B. and I.S.
711 performed digital PCR experiments. G.G., P.B., F.Z. and K.S. wrote the first draft manuscript
712 and all authors edited the manuscript. The authors dedicate this work to our dear late
713 colleague Karin Aistleitner who left us too soon and unexpectedly.

714 **Competing interests**

715 The authors declare no competing interests.