

1 *Can genome-based analysis be the new gold-standard*
2 *for routine *Salmonella* serotyping?*

3 *Sangeeta Banerji^{1§}, Sandra Simon^{1§}, Andreas Tille² and Antje Flieger^{1*}*

4 *¹Robert Koch-Institute, Division of Enteropathogenic Bacteria and Legionella (FG11) / National*
5 *Reference Centre for *Salmonella* and other Bacterial Enteric Pathogens, Burgstrasse 37, 38855*
6 *Wernigerode, Germany*

7 *²Robert Koch-Institute, Department for Infectious Disease Epidemiology (FG31), Burgstrasse 37,*
8 *38855 Wernigerode, Germany*

9 *[§]both authors contributed equally*

10 **Corresponding author*

11 *Contact: fliegera@rki.de*

12 **Abstract**

13 *Salmonella enterica* is the second most reported bacterial cause of food-borne infections in Europe.
14 Therefore molecular surveillance activities based on pathogen subtyping are an important measure
15 of controlling *Salmonellosis* by public health agencies. In Germany, at the federal level, this work is
16 carried out by the National Reference Center for *Salmonella* and other Bacterial Enteric Pathogens
17 (NRC). With rise of next generation sequencing techniques, the NRC has introduced whole-genome-
18 based typing methods for *S. enterica* in 2016. In this study we report on the feasibility of genome-
19 based *in silico* serotyping in the German setting using raw sequence reads. We found that SeqSero
20 and seven gene MLST showed 98% and 95% concordance, respectively, with classical serotyping for
21 the here evaluated serotypes, including the most common German serotypes *S. Enteritidis* and *S.*
22 *Typhimurium* as well as less frequently found serotypes. The level of concordance increased to >99%
23 when the results of both *in silico* methods were combined. However, both tools exhibited
24 misidentification of monophasic variants, in particular monophasic *S. Typhimurium* and therefore
25 need to be fine-tuned for reliable detection of this epidemiologically important variant. We conclude
26 that with adjustments *Salmonella* genome-based serotyping might become the new gold standard.

27 **Introduction**

28 Subtyping of bacterial enteric pathogens, such as *Salmonella enterica*, traditionally relies on
29 serotyping. The species *Salmonella enterica* is divided into six subspecies and consists of more than
30 2600 serovars, which are classified according to the White-Kauffmann- Le Minor Scheme ¹.
31 Serotyping is based on determination of somatic O antigens and flagellin H antigens by reaction with
32 specific antisera. Most *S. enterica* serovars have two alternately expressed H antigens, also referred
33 to as 'phases'. The phase-1 and phase-2 flagellin proteins are encoded by *fliC* and *fliB*, respectively.
34 The phase switch is regulated by the invertase *hin* and the *fliC* repressor gene *fliA* ². Therefore, the
35 specific antigenic formula consists of three positions: the first position represents the O antigens, the
36 second and third positions the two different flagellin H antigens. Each antigen position is separated
37 by a colon, i.e. O:H1:H2. The antigenic formula for *S. Typhimurium* for example is accordingly
38 1,4,[5],12:i:1,2. There are variants of *S. Typhimurium*, which express only one flagellin and which
39 therefore are referred to as monophasic *S. Typhimurium*. *S. Enteritidis* on the other hand does not
40 possess a second flagellin per se, which is reflected in the antigenic formula: 1,9,12:g,m:-. It should
41 be noted that some serovars share the same antigenic formula and require additional testing for
42 unambiguous identification, e.g. the clinically important serovar *S. Choleraesuis* shares its antigenic
43 formula 6,7:c:1,5 with serovars *S. Paratyphi C* and *S. Typhisuis*. A differentiation is possible based on
44 biochemical characteristics or PCR ³.

45 With rise of next generation sequencing (NGS) techniques, genomic typing tools have become
46 increasingly popular and effective. Several *in silico* classification tools employing NGS data are
47 available for *Salmonella*. The serotyping tools are either based on identifying and characterizing the
48 serotype-determining genes or derive the serotype from *in silico* Multi Locus Sequence Typing (MLST)
49 or a combination of both methods. MLST-based serotyping was sparked by the observation of
50 Achtmann et al. that the phylogeny derived from MLST sequence types correlates with serotypes ^{4,5}.
51 Achtmann and his group are also the developers of Enterobase, a platform for the phylogenetic
52 analysis of selected bacteria, including *Salmonella* ⁶. A report of the Establishing Next Generation
53 Sequencing Ability for Genomic Analysis in Europe (ENGAGE) consortium identified four serotyping
54 tools, specifically Metric-Oriented Sequence Typer (MOST), SeqSero, *SalmonellaTypeFinder* and
55 SISTR, which were benchmarked for their performance and were found to have correlation rates
56 between 65% and 88% with classical serotyping (<http://www.engage-europe.eu/resources/benchmarking>). MOST is a pipeline developed and employed by Public Health
57 England, which infers an MLST type with a modified version of the program SRST, which was
58 developed for deducing a sequence type from short reads, and utilizes a local database for
59 identification of corresponding serotypes ^{7,8}. SeqSero is an *in silico* serotyping program and
60 determines the presence of O and H antigen loci within the NGS data, which correspond to the
61 antigens involved in classical serotyping ⁹. *SalmonellaTypeFinder* is a pipeline developed by the
62 Danish Technical University, which runs SeqSero and determines the MLST type using an in-house
63 MLST calling tool, and then both results are used for determination of the serotype
64 (<https://bitbucket.org/genomicepidemiology/salmonellatypefinder/src/master/>). Another typing
65 platform is SISTR, which predicts the serotype by a combination of *in silico* hybridization and
66 extended MLST, incorporated into a 'Microbial *in Silico* Typing' engine ¹⁰.

68 Classification of *Salmonella* by serotyping is especially important for epidemiological investigations
69 and is often routinely performed in its full scheme at National Reference Centers or Laboratories. It is
70 also implemented at the German National Reference Center for *Salmonella* and other Bacterial

71 Enteric Pathogens (NRC). The NRC receives around 3,000-4,000 *Salmonella* isolates per year from
72 human infections for further characterization. The most common serotypes submitted are *S.*
73 *Enteritidis* and *S. Typhimurium*, followed by other broad host range serotypes like *S. Infantis* and *S.*
74 *Derby*¹¹. Since 2016 the NRC has been shifting towards NGS-based analysis¹².

75 Our aim for this study was to estimate whether NGS-based serotyping was feasible as a means of
76 replacing traditional serotyping in our setting. The success rate of classical serotyping depends on
77 many factors (e.g. access to high quality antisera, training of staff, and experience with rare
78 serotypes) and was found to average worldwide at 82% and for European countries at 89% correct
79 results in 2007¹³. Whereas the O antigens are determined within a few hours, characterization of the
80 H phases may require up to 7 days. If the NRC replaced classical serotyping with a genome-based *in*
81 *silico* typing method, this method should ideally match the high reported success rate of classical
82 serotyping¹³. Genome-based typing tools have performed well in several studies with maximum
83 reported concordance levels of approximately 92% for SeqSero⁹ and approximately 94 % for SISTR
84^{10,14}. However, previous studies used assembled genomes for *in silico* typing. Only very recently,
85 Ibrahim and Morin also reported results obtained with paired reads using the web-based application
86 of SeqSero 1.0¹⁵. Genome assembly requires additional time and computing resources, which is a
87 drawback for routine analysis of a large number of genomes.

88 Our goal for this study was therefore to directly use raw reads in order to save time and computing
89 resources. Thus, our requirements for the tools were that the input data should need minimal
90 preprocessing and should potentially fit into our existing analysis pipeline (Ridom SeqSphere⁺)¹⁶.
91 Since we need to process a large number of sequences, offline availability was also of major
92 importance. SeqSero fulfilled all of these requirements (when used as a command line tool). The
93 other above mentioned tools did not as they either use different allele detection algorithms for
94 determination of MLST sequence types than Enterobase (MOST and *SalmonellaTypeFinder*) or
95 require an assembled genome (SISTR). Therefore we decided to assess the performance of SeqSero
96 and the Enterobase MLST scheme from Achtman et al. for serotype prediction⁴.

97

98 **Results**

99 The aim of this study was to assess two *in silico* serotype prediction tools, namely SeqSero and MLST
100 via SeqSphere/Enterobase for their performance in routine *Salmonella* typing at the NRC. We chose
101 520 *Salmonella* isolates, mainly of human origin and predominantly from the years 2014-2018 as the
102 data set for analysis. The selection comprised very frequently found serotypes as well as less
103 frequent serotypes (Table 1). We investigated a total of 20 different serotypes and also looked at
104 monophasic variants as well as rough phenotypes.

105 Data quality is an important bottleneck

106 Initially, we did not set a quality threshold for the raw read sequence files. In the course of the
107 analysis we noticed that analysis with SeqSero 1.0 and/or Ridom SeqSphere⁺ failed if the file sizes of
108 the raw sequence reads were lower than average (<50,000 KB). Since Zhang et al. only included data
109 for analysis with SeqSero 1.0 with a minimal coverage of 10-fold, we aimed for the same quality
110 threshold⁹. Given an average genome size of approximately 4.8 Mb for *Salmonella enterica* subsp.
111 *enterica*, we calculated that a theoretical coverage of ≥10-fold could only be achieved by a minimal

112 read number of 100,000 each for paired end reads with a theoretical read length of 250: theoretical
113 coverage = total number of reads x length of each read [bp] / genome size [bp]. Sequencing was
114 repeated for cases not meeting the minimal read number (Fig. S1).

115 SeqSero analysis correctly predicted the serovar in 98% of the isolates

116 SeqSero 1.0 predicted the serotype in 84% of analyzed strains in accordance to the classical serovar.
117 In additional 14 % of the cases the antigenic formula was shared by more than one serovar and
118 SeqSero 1.0 predicted all eligible serotypes, e.g. Choleraesuis, Typhisuis or Paratyphi C for the
119 antigenic formula 6,7:c:1,5, which we rated as ambiguous. These cases require additional testing as
120 they would if determined by classical serotyping. Therefore an ambiguous prediction was counted as
121 a correlating result in the overall summary. The total rate of correlation (correlation + ambiguous
122 prediction) with our laboratory results was therefore 98% (Table 1). Five cases of prediction failure
123 (1%) occurred, all of which involved failed prediction of the O-7 antigen. Additionally, five cases (1%)
124 of miscorrelation were found, which concerned monophasic strains of *S. Typhimurium* and *S.*
125 *Choleraesuis* (Table 1).

126 Monophasic variants are only predicted correctly if they lack the flagellin genes

127 17 out of 19 monophasic *S. Typhimurium* strains were correctly predicted by SeqSero 1.0 using raw
128 reads (Table 1). In two cases, SeqSero 1.0 predicted phenotypically monophasic *S. Typhimurium* as
129 biphasic. In order to investigate this discrepancy we analyzed the respective whole genome
130 sequences by *de novo* assembly. The isolate ERR2003330 lacked approximately 250 nucleotides in
131 the central part of the *fliB* gene as well as the whole *hin* gene (Fig. S2). Expression of the phase-2
132 flagellin gene *fliB* is co-regulated by the invertase gene *hin* and the *fliC* repressor gene *fliA*².
133 Apparently a transposase, *tnpA*, had integrated into this region. This explains why the second phase
134 could not be detected by classical serotyping. Since SeqSero 1.0 only checks whether the *fliC* and *fliB*
135 alleles are present, it would explain why the lack of the *hin* gene was not detected by the program
136 and the partial deletion of the *fliB* gene might have been too small to be detectable when using raw
137 reads. We noted that SeqSero 1.0 correctly predicted the isolate to be monophasic when the analysis
138 was performed with a5-assembled contigs. During preparation of this manuscript a new version of
139 SeqSero called SeqSero 2.0 was available from github (<https://github.com/denglab/SeqSero 2.0>) and
140 we rechecked the two non-correlating results with SeqSero 2.0 in the default k-mer-based mode. The
141 program correctly classified isolate ERR2003330 as monophasic, probably due to the partial deletion
142 in the *fliB* gene. However, when we used SeqSero 2.0 with a5 assembled contigs it classified the
143 isolate wrongly as biphasic *S. Typhimurium*. The second isolate ERR2003327 had a transposon
144 integrated into the *fliB* gene most probably rendering it non-functional. This isolate was identified to
145 be biphasic with both SeqSero 1.0 and the k-mer-based approach of SeqSero 2.0 when using raw
146 reads, because the *fliB* gene is fully present but interrupted. When using SeqSero 1.0 and SeqSero 2.0
147 with a5 assembled contigs isolate ERR2003327 was correctly predicted to be monophasic by both
148 versions of the program.

149

150 Phenotypically monophasic variants of other serovars harboring *fliC* and *fliB* were also not recognized
151 by SeqSero 1.0 or SeqSero 2.0, in particular three strains of *S. Choleraesuis* var. Kunzendorf not
152 expressing phase-1 flagellum gene *fliC* (ERR3264001, ERR3264026, ERR3264035). This corroborates
153 the fact that phenotypic traits are sometimes difficult to detect by *in silico* measures. A monophasic

154 variant of *S. Paratyphi B* variant Java was recognized as monophasic *S. Paratyphi B* by SeqSero 1.0
155 and was additionally recognized to be L(+)-tartrate positive by SeqSero 2.0.

156 Serovar prediction of rough strains is possible by means of SeqSero

157 Importantly, SeqSero 1.0 was able to predict a serotype for five out of six isolates with a rough
158 phenotype, where classical serotyping was not successful (ERR3263893: *S. Typhimurium*,
159 ERR3263889: *S. Typhimurium* monophasic, ERR3263894 [*S. e.* subspecies II]: 58:z6:z39, ERR3263880:
160 *S. Typhimurium* monophasic, and ERR3263875: *S. Typhimurium* monophasic). We classified this as a
161 correlation. For the rough strain ERR3264036 SeqSero 1.0 did not generate a full antigenic formula.
162 PCR analysis according to Woods et al. 2008 targeting a 12.8-kb region specific to *S. Choleraesuis*
163 yielded the serotype Choleraesuis ³. In this case, SeqSero 1.0 was only able to predict a partial
164 antigenic formula for Choleraesuis (-:c:1,5). SeqSero 2.0 was likewise not able to provide a complete
165 antigenic formula for this particular strain. When we mapped the raw reads against the respective
166 *wzy* allele (locus tag EL48_RS10980), we found the allele and the surrounding region (EL48_RS10955-
167 EL48_RS11010) missing (Fig. S3). We conclude that the rough phenotype of this particular isolate had
168 a genetic basis.

169 SeqSero does not reliably predict the O-7 antigen

170 We found five cases of prediction failure when using SeqSero 1.0 and all five cases involved failed
171 prediction of the O-7 antigen, which is part of the epidemiologically important serovars *S.*
172 *Choleraesuis* and *S. Infantis* (ERR3264036, ERR3264076, ERR3264063, ERR3264067, and
173 ERR3264066). Except for Isolate ERR3264036, the remaining four cases had an intact *wzy* allele but
174 only few reads mapped to the O-7 locus (Fig. S4 and S5). When we performed the analysis with
175 SeqSero 2.0 the k-mer based approach yielded a complete antigenic formula which correlated with
176 the laboratory phenotypes for all cases except for the rough strain ERR3264036, where the *wzy* allele
177 is missing.

178 SeqSero is well suited for routine high-throughput analysis of raw reads with the exception of
179 atypical monophasic strains

180 In summary, SeqSero 1.0 is an easy to use tool, which is available as free software from the website
181 of the developers, or as an official Debian package from the Debian website. Currently an alpha test
182 version of SeqSero 2.0 is available on Github with additional features, e.g. k-mer based approach and
183 integrated identification of the taxonomic ID with SalmID in the allele based mode for subspecies
184 identification of ambiguous serovars. When using SeqSero 1.0 with Illumina paired end raw reads we
185 achieved a correlation rate of 98%. The reasons for initial miscorrelation were mainly low data
186 quality, which could be resolved by repeating the sequencing (Fig. S1). SeqSero 1.0 was able to
187 predict a serotype for all rough isolates, except one. It correctly predicted monophasic variants if the
188 flagellin genes *fliC* and/or *fliB* were missing. However, if the flagellin genes were only disrupted
189 and/or other genes required for flagellar expression / phase transition were missing, SeqSero 1.0 and
190 SeqSero 2.0 were not always able to reliably recognize monophasic variants. We conclude that with
191 the exception of atypical monophasic variants of *S. Typhimurium* and other serovars and genetically
192 rough strains (i.e. lack of O antigen determining genes) SeqSero is able to correctly predict the vast
193 majority of common serovars circulating in Germany.

194

195 MLST analysis correctly predicted the serovar in 95% of the isolates

196 MLST predicted the serotype in 95% of *Salmonella* isolates in concordance to the classical serovar
197 found by serotyping (Table 1). Notably, all six rough isolates were assigned to a sequence type and a
198 corresponding serotype. The prediction differed in 25 cases (5%) from the phenotypic classification
199 all of which involved second phase miscorrelation. Figure 1 shows an UPGMA (unweighted pair group
200 method with arithmetic mean) Tree based on MLST and color coded according to the serovar
201 obtained by slide agglutination. As expected, there is a clear correlation between serotype and one
202 or more closely related STs for the majority of isolates (Fig. 1). *S. Enteritidis* for example is distributed
203 into the two closely related STs: ST 11 and ST 183. The *S. Typhi* isolates of our collection spread
204 across five different but closely related STs: ST 1, ST2, ST 3677, ST 2173 and ST 2209 (Fig. 1).

205 Sequence types do not consistently correlate with detection of flagellin antigens

206 It is notable that for the majority of isolates in Enterobase the antigenic formula is not provided by
207 the user. Nevertheless, the majority of Enterobase strains belonging to ST 34, which had an antigenic
208 formula provided, represented monophasic *S. Typhimurium* (203 out of 209 isolates as of May 2019).
209 Therefore we assigned all ST 34 strains to monophasic *Typhimurium*. Enterobase strains belonging to
210 ST 19 were a mix of monophasic and biphasic *Typhimurium*. We opted to classify all ST 19 isolates as
211 biphasic *Typhimurium* although this would result in a high error rate. We preferred this to no
212 classification at all. We obtained correlating results between MLST and classical serotyping for 17 out
213 of 19 (89.5%) of our monophasic *S. Typhimurium* strains. Only 32 out of 52 biphasic *S. Typhimurium*
214 belonged to ST 19 (61.5%) and were therefore also classified as biphasic with MLST. 20 out of 52
215 (38.5%) phenotypically biphasic *Typhimurium* belonged to ST 34 and were therefore wrongly
216 classified as monophasic by MLST. We also checked whether the classification of monophasic and
217 biphasic *S. Typhimurium* would be improved by clustering according to core genome MLST. Figure 2
218 depicts a minimum spanning tree of only monophasic and biphasic *S. Typhimurium* isolates (including
219 three rough isolates) based on the Enterobase core genome MLST scheme. The isolates cluster
220 according to their ST rather than to their flagellin expression.

221 The *S. Choleraesuis* isolates of our collection, phenotypically lacking FliC were also not correctly
222 classified by MLST typing. In Enterobase monophasic *S. Choleraesuis* var. Kunzendorf predominantly
223 belonged to ST 66, whereas our isolates belonged to ST 145. Interestingly, MLST distinguished the
224 monophasic *S. Paratyphi* B var. Java as such, since ST 42 mostly consists of monophasic var. Java
225 entries in Enterobase.

226 MLST-based serotype prediction additionally provides phylogenetic context

227 The majority of our serotypes could each be assigned to a single eBG: e.g. *S. Typhimurium* to eBG 1,
228 *S. Enteritidis* to eBG 4, *S. Typhi* to eBG 13 and *S. Choleraesuis* to eBG 6 (Table 2). This is also reflected
229 in the phylogenetic tree, where the different STs, which comprise the same serovar and belong to
230 the same eBG are located in the same branch (Fig. 1). This indicates that German strains belonging to
231 these serovars stem from a common ancestor ^{4,17}. One advantage of MLST serotype prediction
232 compared to SeqSero was that there was no ambiguous serotype prediction. Different serovars with
233 the same antigenic formula split into distinct eBurst groups (e.g. *S. Choleraesuis* eBG 6 and *S.*
234 *Paratyphi* C eBG 20). MLST additionally provided important phylogenetic information, e.g. the *S.*
235 *Derby* strains in our collection were of a polyphyletic nature as they split into three different eBGs
236 (Table 2 and Fig. 1). In conclusion, MLST-based serotype prediction also proved to be very successful

237 with the draw-back of not being able to distinguish between monophasic and biphasic *S.*
238 *Typhimurium* as well as between *S. Choleraesuis* and monophasic *S. Choleraesuis* var. *Kunzendorf*.

239

240 Combination of SeqSero and MLST increases robustness of prediction

241 After performing both analyses independently, we combined SeqSero 1.0 and MLST and used both
242 results for predicting the serotype. In general, there was good agreement between the two methods.
243 In case of disagreement, we evaluated the sequences individually. There were 24 cases of
244 disagreement between SeqSero and MLST all of which concerned phase variation. Since our findings
245 indicated that MLST was not suited for identification of phase variation and SeqSero generally
246 performed better in this regard, we rated the SeqSero result as more adequate. There was
247 disagreement between SeqSero and MLST regarding 20 *S. Typhimurium* isolates of ST 34, which were
248 classified as monophasic by MLST and biphasic by SeqSero. Since biphasic ST 34 isolates cannot be
249 correctly classified by MLST we chose the SeqSero prediction for these cases. The same applied for
250 monophasic ST 19 *S. Typhimurium* isolates, which were also not correctly classified by MLST. The two
251 isolates, which carried a transposase in *fjB* (ERR2003330 and ERR2003327), were correctly predicted
252 as monophasic by MLST and here we opted for the MLST prediction because we had already
253 analyzed these isolates by mapping. In the 5 cases of prediction failure by SeqSero, we chose the
254 MLST prediction as the serovar. This way, the percentage of correlation was increased to >99%. In
255 summary the combination of both independent methods enabled the identification of potential
256 misclassifications where a closer analysis was necessary and thus reduced the rate of error.

257

258 **Discussion**

259 In this study we evaluated two genome-based *in silico* approaches and their combination for
260 predicting *Salmonella* serotypes and their suitability for replacing classical serotyping. Table 3
261 summarizes the advantages and drawbacks of the three typing methods. We found that both tested
262 prediction methods, the *in silico* serotyping approach by SeqSero 1.0 and the indirect serotype
263 prediction with MLST yielded excellent correlation with our laboratory-based results analyzing 520
264 isolates from our strain collection (98% SeqSero, 95% MLST). Since our collection lacked a
265 representative selection of strains of rare serotypes or higher subspecies we cannot rate the
266 performance in this regard. Nonetheless it was representative of the most common human strains in
267 Germany.

268 Our collection also included a novel serovar, derived from an outbreak related to sesame seeds ¹⁸.
269 Interestingly, the antigenic formula of this novel serovar was correctly identified by SeqSero
270 demonstrating its effectiveness for classifying novel serovars. Our correlation rate of 98% using raw
271 reads matches very well the correlation rate determined by the developers of SeqSero of 98.7% using
272 308 CDC strains ⁹. However, the correlation rate found by Zhang and colleagues dropped to 92.6%
273 when using a higher number of isolates, i.e 3306 isolates from GenomeTrakr. Likewise, a recent study
274 of 1041 environmental *Salmonella* isolates including a wider variety than our study yielded a
275 correlation of 86% to classical serotyping ¹⁵. Recently, the developers of SeqSero presented a new
276 version of the program named SeqSero 2.0 at the International Symposium on *Salmonella* and
277 *salmonellosis* 2018 ¹⁹. SeqSero 2.0 can use SalmID in the assembly mode for subspecies identification

278 of ambiguous serovars (www.github.com/hcdeebakker/salmID). We did not test the assembly
279 mode since it required the additional program SalmID, which we did not include in our assessment.
280 We tested SeqSero 2.0 in its default k-mer based mode for reassessment of the ten cases where
281 SeqSero failed. We found that with the default settings, SeqSero 2.0 also did not consistently detect
282 monophasic variants of *S. Typhimurium* but showed improved performance in cases of high
283 sequence variability.

284 Our results indicate that SeqSero does not reliably predict monophasic variants, in particular
285 monophasic *S. Typhimurium*. Monophasic *S. Typhimurium* lacking *fliB* are correctly classified by
286 SeqSero but atypical monophasic variants where *fliB* is present may be misclassified as biphasic. This
287 is a potentially crucial limitation of the program as monophasic variants, especially of *S.*
288 *Typhimurium*, are epidemiologically important and the latter comprise approximately 2/3 of the *S.*
289 *Typhimurium* received at the NRC²⁰⁻²⁴. We suggest including the detection of additional factors to the
290 *fliB* allele, which determine integrity of the second phase flagellar antigen. Also the algorithm for
291 phase determination when using raw reads should be refined so that disruptions in the *fliC* / *fliB*
292 genes can be detected in spite of the fact that the gene is fully present.

293 Regarding MLST, it was foreseeable by examining the strains in Enterobase that a clear classification
294 between monophasic and biphasic *S. Typhimurium* based on ST would not be possible. Achtman et
295 al. did not find a correlation between ST and monophasic *S. Typhimurium* when they analyzed a large
296 and diverse collection⁴. On the other hand, it was reported that Italian and UK monophasic *S.*
297 *Typhimurium* strains belonged to ST 34^{20,21}. Petrovska et al. showed that the current monophasic
298 epidemic *S. Typhimurium* strains evolved from at least three independent events²¹. The monophasic
299 strains of our collection predominantly belong to the current European ST 34 epidemic clone,
300 therefore a good correlation for monophasic strains of ST 34 was obtained with MLST. On the other
301 hand, biphasic strains of ST 34 were misclassified as monophasic. We therefore conclude that the
302 classical MLST scheme alone is not able to clearly distinguish between monophasic and biphasic *S.*
303 *Typhimurium* due to their polyphyletic nature. Our results further indicate that clustering by core
304 genome MLST does also not improve classification according to flagellin expression. Since recent
305 studies have found *S. Typhimurium* regions, which seem characteristic for certain monophasic
306 variants it may be possible to develop an additional scheme based on the presence/absence of such
307 specific genes to reliably identify monophasic variants of *S. Typhimurium*^{21,24}.

308 We obtained the highest correlation to classical serotyping when we combined the predictions of
309 SeqSero and MLST because the two methods use independent approaches for serotype
310 determination and thereby complemented each other. Since SeqSero directly generates an antigenic
311 formula, we rated its output as more adequate than the indirect determination by MLST.
312 Nonetheless, with the additional information provided by MLST, it was possible to clarify all
313 ambiguous predictions by SeqSero because the serovars, which shared the same antigenic formula,
314 had different STs. Our results also indicate that MLST might even perform better in classifying rough
315 strains than SeqSero. The combined prediction increased robustness because miscorrelating
316 predictions of the two programs gave rise to more detailed analysis. Currently there are two tools
317 available, which use the combined prediction of *in silico* serotyping and MLST. One is
318 SalmonellaTypeFinder, which uses SeqSero and MLST and thus has the potential of performing well
319 (<https://bitbucket.org/genomicepidemiology/salmonellatypewriter/src/master/>). We did not
320 evaluate this tool in our study because it uses a different MLST calling algorithm than we routinely do
321 and has not been published yet. The second tool is SISTR, which predicts the serotype with the help

322 of *in silico* genoserotyping and validates the results with core genome MLST¹⁰. We did not evaluate
323 SISTR because it requires assembled genomes. However, it performed very well in a previous
324 report¹⁴. A combination of genome-based serotyping and MLST is also advocated by other
325 governmental agencies like Public Health England, who use MOST and Public Health Agency of
326 Canada who use SISTR¹⁴.

327

328 Conclusion

329 SeqSero is an *in silico* serotyping tool generating an antigenic formula directly comparable to classical
330 serotyping. MLST provides important phylogenetic information and is able to distinguish serovars
331 with the same antigenic formula. The concomitant use of both tools seems best suited for *in silico*
332 strain characterization to obtain the utmost information and a robust prediction. Nevertheless, some
333 improvements are necessary to differentiate monophasic from biphasic strains. If the serotype is
334 predicted by these two independent methods, a disagreement could indicate a potential problem
335 requiring further investigation. Since we obtained a correlation rate of >99% for SeqSero in
336 combination with MLST, we conclude that the here investigated *in silico* typing tools could in
337 combination outperform the current gold standard of phenotypic serotyping and could become the
338 new gold standard.

339

340 Methods

341 Short read sequencing

342 Whole genome sequencing was performed at the NRC or at the Robert Koch-Institute's sequencing
343 core facility on a MiSeq benchtop sequencer using Illumina's MiSeq Reagent Kit v3, yielding 2 x 300
344 bp paired end reads. Adapter-clipped reads were obtained from the sequencing unit and used in this
345 study without additional processing unless stated otherwise. Sequencing was repeated for cases not
346 meeting the minimal read number of 100,000 (Fig. S1). The fastq files of the paired-end sequence
347 reads are available from the European Nucleotide Archive under the project numbers PRJEB30317 &
348 PRJEB16326. Project PRJEB16326 is part of EU COMPARE (<https://www.compare-europe.eu/>) and a
349 subset of the German samples of that project have been included in this study.

350 SeqSero 1.0/SeqSero 2.0

351 The SeqSero 1.0 command line tool was downloaded from Github
352 (<https://github.com/denglab/SeqSero>) and an official Debian package was created, which is available
353 from <https://blends.debian.org/med/tasks/bio>. The installed program was then embedded into a
354 script for batch analysis. Illumina MiSeq paired-end reads were directly used for serotype prediction.
355 Apart from choosing the correct mode for the input data, i.e. single-end, paired-end, interleaved or
356 assembled, the program offers no additional options. During drafting of this manuscript an alpha test
357 version of SeqSero 2.0 became available from Github (<https://github.com/denglab/SeqSero 2.0>). We
358 used SeqSero 2.0 with its default setting (k-mer based mode) only to analyze isolates where SeqSero
359 1.0 did not produce a correlating result to classical serotyping.

360 Ridom SeqSphere settings and allele calling procedure

361 For MLST analysis we used the 7 gene MLST scheme from Achtman *et al.* embedded in the Ridom
362 SeqSphere⁺ software (Ridom GmbH, Münster, Germany)⁴. Please note, that in spite of the fact that
363 the scheme recommends *de novo* assembly of raw reads, we used mapping in order to save time and
364 resources. Using the raw reads, the pipeline quality-trimmed and mapped the Illumina MiSeq reads
365 against the reference genome *S. Typhimurium* LT2 (GenBank AE006468.2) using the build-in
366 Burrows-Wheeler Aligner in the default mode. This ideally yielded allele numbers for the seven
367 housekeeping genes and the corresponding sequence type (ST). If Ridom SeqSphere⁺ was not able to
368 assign a ST there were generally two reasons: either low data quality ('Target QC procedure failure')
369 or it was a potential new ST. For cases of low sequence quality sequencing was repeated (Fig. S1).
370 For phylogenetic analysis of monophasic and biphasic *S. Typhimurium* isolates the Enterobase core
371 genome MLST scheme was used in SeqSphere⁺.

372 Assigning Sequence types and corresponding serotypes with Enterobase

373 The obtained MLST sequence types were entered into Enterobase to find corresponding serotypes
374 from the database and if available the e-burst groups (eBGs). eBGs determination is based on an
375 algorithm, which identifies the relationship of isolates with similar genotypes¹⁷. Enterobase
376 periodically confers official eBG numbers to new eBGs.

377 If Ridom SeqSphere⁺ reported a potential new ST we uploaded the NGS data of the respective
378 isolates to Enterobase in order to obtain an official ST.

379 De novo assembly and mapping

380 For isolates with non-correlating results *de novo* assembly was performed using A5 or SPAdes^{25,26}
381 Some isolates were further analyzed by mapping the raw reads against specific loci using the
382 Geneious mapper or Bowtie2 in Geneious (www.geneious.com).

383

384 **Data Availability**

385 The raw sequence reads analyzed in this study are publicly available at the European Nucleotide
386 Archive under the project accession numbers PRJEB30317 and PRJEB16326. PRJEB16326 is part of
387 COMPARE and a subset of the German samples has been included in the current study. An overview
388 of all strains and metadata is given in Table S1.

389

390 **References**

- 391 1 WHO Collaborating Centre for Reference and Research on Salmonella. *Antigenic formulae of*
392 *the Salmonella serovars*. 9edn, (1997).
- 393 2 Barco, L. *et al.* Molecular Characterization of "Inconsistent" Variants of *Salmonella*
394 *Typhimurium* Isolated in Italy. *Foodborne Pathog Dis* **11**, 497-499,
395 doi:10.1089/fpd.2013.1714 (2014).
- 396 3 Woods, D. F. *et al.* Rapid Multiplex PCR and Real-Time TaqMan PCR Assays for Detection of
397 *Salmonella enterica* and the Highly Virulent Serovars *Choleraesuis* and *Paratyphi C*. *J Clin*
398 *Microbiol* **46**, 4018-4022, doi:10.1128/Jcm.01229-08 (2008).
- 399 4 Achtman, M. *et al.* Multilocus Sequence Typing as a Replacement for Serotyping in
400 *Salmonella enterica*. *Plos Pathog* **8**, doi:10.1371/journal.ppat.1002776 (2012).

- 401 5 Alikhan, N. F., Zhou, Z. M., Sergeant, M. J. & Achtman, M. A genomic overview of the
402 population structure of *Salmonella*. *Plos Genet* **14**, doi:10.1371/journal.pgen.1007261 (2018).
- 403 6 Zhou Z., Alikhan N.F., Mohamed K., the Agama Study Group, Achtman M. The user's guide to
404 comparative genomics with Enterobase. Three case studies: micro-clades within *Salmonella*
405 *enterica* serovar Agama, ancient and modern populations of *Yersinia pestis*, and core
406 genomic diversity of all *Escherichia*. <https://doi.org/10.1101/613554> (2019).
- 407 7 Tewolde, R. *et al.* MOST : a modified MLST typing tool based on short read sequencing. *Peerj*
408 **4**, doi:10.7717/peerj.2308 (2016).
- 409 8 Inouye, M., Conway, T. C., Zobel, J. & Holt, K. E. Short read sequence typing (SRST): multi-
410 locus sequence types from short reads. *Bmc Genomics* **13**, doi:10.1186/1471-2164-13-338
411 (2012).
- 412 9 Zhang, S. K. *et al.* Salmonella Serotype Determination Utilizing High-Throughput Genome
413 Sequencing Data. *J Clin Microbiol* **53**, 1685-1692, doi:10.1128/Jcm.00323-15 (2015).
- 414 10 Yoshida, C. E. *et al.* The Salmonella In Silico Typing Resource (SISTR): An Open Web-
415 Accessible Tool for Rapidly Typing and Subtyping Draft *Salmonella* Genome Assemblies. *Plos*
416 *One* **11**, doi:10.1371/journal.pone.0147101 (2016).
- 417 11 Hauser, E. *et al.* Clonal Dissemination of *Salmonella enterica* Serovar Infantis in Germany.
418 *Foodborne Pathog Dis* **9**, 352-360, doi:10.1089/fpd.2011.1038 (2012).
- 419 12 Simon, S. *et al.* Evaluation of WGS based approaches for investigating a food-borne outbreak
420 caused by *Salmonella enterica* serovar Derby in Germany. *Food Microbiol* **71**, 46-54,
421 doi:10.1016/j.fm.2017.08.017 (2018).
- 422 13 Hendriksen, R. S. *et al.* WHO Global Salm-Surv external quality assurance system for
423 serotyping of *Salmonella* isolates from 2000 to 2007. *J Clin Microbiol* **47**, 2729-2736,
424 doi:10.1128/JCM.02437-08 (2009).
- 425 14 Yachison, C. A. *et al.* The Validation and Implications of Using Whole Genome Sequencing as
426 a Replacement for Traditional Serotyping for a National *Salmonella* Reference Laboratory.
427 *Front Microbiol* **8**, doi:10.3389/fmicb.2017.01044 (2017).
- 428 15 Ibrahim, G. M. & Morin, P. M. *Salmonella* Serotyping Using Whole Genome Sequencing. *Front*
429 *Microbiol* **9**, 2993, doi:10.3389/fmicb.2018.02993 (2018).
- 430 16 Junemann, S. *et al.* Updating benchtop sequencing performance comparison. *Nat Biotechnol*
431 **31**, 294-296, doi:10.1038/nbt.2522 (2013).
- 432 17 Feil, E. J., Li, B. C., Aanensen, D. M., Hanage, W. P. & Spratt, B. G. eBURST: Inferring patterns
433 of evolutionary descent among clusters of related bacterial genotypes from multilocus
434 sequence typing data. *J Bacteriol* **186**, 1518-1530, doi:10.1128/Jb.186.5.1518-1530.2004
435 (2004).
- 436 18 Schielke, A.*., Simon, S *et al.* European-wide salmonellosis outbreak with a novel serotype
437 (11:z41:e,n,z15) attributable to sesame products, 2016-2017. *Eurosurveillance* (accepted for
438 publication).
- 439 19 Colin, P. International Symposium on *Salmonella* and salmonellosis. *Food Microbiol* **71**, 1,
440 doi:10.1016/j.fm.2017.12.009 (2018).
- 441 20 Mastorilli, E. *et al.* A Comparative Genomic Analysis Provides Novel Insights Into the
442 Ecological Success of the Monophasic *Salmonella* Serovar 4,[5],12:i:-. *Front Microbiol* **9**,
443 doi:10.3389/fmicb.2018.00775 (2018).
- 444 21 Petrovska, L. *et al.* Microevolution of Monophasic *Salmonella* Typhimurium during Epidemic,
445 United Kingdom, 2005-2010. *Emerging infectious diseases* **22**, 617-624,
446 doi:10.3201/eid2204.150531 (2016).

- 447 22 Alt, K. *et al.* Outbreak of Uncommon O4 Non-Agglutinating *Salmonella* Typhimurium Linked
448 to Minced Pork, Saxony-Anhalt, Germany, January to April 2013. *Plos One* **10**,
449 doi:10.1371/journal.pone.0128349 (2015).
- 450 23 Hauser, E. *et al.* Pork Contaminated with *Salmonella enterica* Serovar 4,[5],12:i:-, an
451 Emerging Health Risk for Humans. *Appl Environ Microb* **76**, 4601-4610,
452 doi:10.1128/Aem.02991-09 (2010).
- 453 24 Palma, F. *et al.* Genome-wide identification of geographical segregated genetic markers in
454 *Salmonella enterica* serovar Typhimurium variant 4,[5],12:i:-. *Sci Rep-Uk* **8**,
455 doi:10.1038/s41598-018-33266-5 (2018).
- 456 25 Tritt, A., Eisen, J. A., Facciotti, M. T. & Darling, A. E. An Integrated Pipeline for de Novo
457 Assembly of Microbial Genomes. *Plos One* **7**, doi:10.1371/journal.pone.0042304 (2012).
- 458 26 Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its Applications to
459 Single-Cell Sequencing. *J Comput Biol* **19**, 455-477, doi:10.1089/cmb.2012.0021 (2012).
- 460

461 **Acknowledgements**

462 We thank Monique Duwe, Marita Wahnfried and Susanne Kulbe for excellent technical assistance.
463 We thank the RKI NGS sequencing team in Berlin and Wernigerode. We are also thankful to all
464 laboratories that sent strains to the NRC. The authors acknowledge funding by the EU COMPARE
465 project within the Horizon 2020 (AF).

466 **Author contributions**

467 SB performed the *in silico* analysis and drafted the manuscript. SS supervised the serotyping and
468 sequencing, advised in cgMLST analysis and provided input to the manuscript. AT wrote scripts for
469 batch analysis and extraction of results. AF conceived and supervised the project and provided input
470 to the manuscript.

471 **Competing interests**

472 The authors declare no competing interests.

473 **Figure Legends**

474 Fig.1 Unweighted Pair Group Method with Arithmetic mean (UPGMA) tree of all investigated isolates
475 based on 7-gene MLST. The tree shows that serovars correlate to STs. Colors are based on ST.

476 Fig.2 Minimal Spanning tree of monophasic and biphasic Typhimurium isolates based on the
477 Enterobase core genome MLST scheme and 7-gene MLST. The tree reveals that *S. Typhimurium*
478 isolates cluster according to ST rather than expression of flagellin. Colors are based on phase and STs.

479

480 Table 1. Overview of serotype prediction with SeqSero and MLST. Serotype was first determined by
 481 classical serotyping. Whole genome sequences were then analyzed with SeqSero or MLST.
 482 Correlation means that the predicted serotype was the same as the classically determined serovar.
 483 Ambiguous means that the correct serotype was listed among others. Prediction failure means that
 484 no complete antigenic formula was derived. Misclassification means that a wrong antigenic formula
 485 was derived. Overall (%) is the same as Total (%) except that Ambiguous (%) is added to Correlation
 486 (%). Final results are shown, i.e. after resequencing if data quality was not met.

Serotype	Sequen- ced Isolates	Correlation			Ambiguous			Prediction failure			Misclassification		
		Seq- Sero	MLST	Seq- Sero + MLST	Seq- Sero	MLST	Seq- Sero + MLST	Seq- Sero	MLST	Seq- Sero + MLST	Seq- Sero	MLST	Seq- Sero + MLST
Agona	3	3	3	3	0	0	0	0	0	0	0	0	0
Choleraesuis	33	0	33	33	30 or Typhisuis or Paratyphi C	0	0	3	0	0	0	0	0
Choleraesuis monophasic	3	0	0	0	0	0	0	0	0	0	3	3	3
Derby	55	55	55	55	0	0	0	0	0	0	0	0	0
11:z41:e,n,z15 (novel serovar)	10	10	10	10	0	0	0	0	0	0	0	0	0
Enteritidis	115	115	115	115	0	0	0	0	0	0	0	0	0
Infantis	50	49	50	50	0	0	0	1	0	0	0	0	0
Kentucky	7	7	7	7	0	0	0	0	0	0	0	0	0
Kintambo	3	0	3	3	3 or Washington	0	0	0	0	0	0	0	0
Kottbus	12	0	12	12	12 or Ferruch	0	0	0	0	0	0	0	0
Mbandaka	15	15	15	15	0	0	0	0	0	0	0	0	0
Mikawasima	10	10	10	10	0	0	0	0	0	0	0	0	0
Muenchen	25	0	25	25	25 or Virginia	0	0	0	0	0	0	0	0
Paratyphi B	6	6	6	6	0	0	0	0	0	0	0	0	0
Paratyphi B monophasic	1	1	1	1	0	0	0	0	0	0	0	0	0
Paratyphi C	2	0	2	2	2 or Cholerae- suis or Typhisuis	0	0	0	0	0	0	0	0
Poano	2	2	2	2	0	0	0	0	0	0	0	0	0
Strathcona	2	2	2	2	0	0	0	0	0	0	0	0	0
Stourbridge	14	14	14	14	0	0	0	0	0	0	0	0	0
Sundsvall	1	0	1	1	1 or Soohanina or Sundvall	0	0	0	0	0	0	0	0
Typhi	74	74	74	74	0	0	0	0	0	0	0	0	0
Typhimurium biphasic	52	52	32	52	0	0	0	0	0	0	0	20	0
Typhimurium monophasic	19	17	17	19	0	0	0	0	0	0	2	2	0
Serologically rough	6	5	6	6	0	0	0	1	0	0	0	0	0
Total number	520	437	495	517	73	0	0	5	0	0	5	25	3
Total (%)	100.0	84.0	95.2	99.4	14.0	0	0	1.0	0	0	1.0	4.8	0.6
Overall (%)	100.0	98.0	95.2	99.4	-	-	-	1.0	0	0	1.0	4.8	0.6

487

488

489 Table 2. Overview of Serovars with corresponding MLST sequence types and e-Burst groups

<i>Salmonella</i> Serotype (Enterobase)	Sequence type	e-Burst Group	Number of Isolates
Agona	13	54	3
Choleraesuis	139	6	1
Choleraesuis	145	6	36
Derby	39	57	6
Derby	774	57	1
Derby	40	57	5
Derby	71	244	2
Derby	682	264	41
Enteritidis	11	4	110
Enteritidis	183	4	5
11:z41:e,n,z15	2914	472	10
Infantis	32	31	49
Infantis	2283	31	1
Kentucky	198	56	7
Kintambo	407	400	1
Kintambo	2839	ND	1
Kintambo	5841	ND	1
Kottbus	212	64	11
Kottbus	1669	63	1
Mikawasima	1815	247	10
Mbandaka	413	62	15
Muenchen	82	8	25
Paratyphi B	86	5	6
Paratyphi B mono (var Java)	42	32	1
Paratyphi C	146	20	2
Poano	557	87	2
Strathcona	2559	ND	2
Stourbridge	736	438	8
Stourbridge (only RKI data)	3736	464	6
Sundsvall (first typed as Poano)	488	305.2	1
Subsp. II	781	340	1
Typhi	1	13	38
Typhi	2	13	32
Typhi	2173	13	1
Typhi	2209	13	1
Typhi	3677	13	2
Typhimurium & monophasic var.	19	1	36
Typhimurium & monophasic var.	34	1	39
Total	39	>26	520

490

491

492 Table 3. Overview of advantages and drawbacks of the investigated typing methods and their sources
493 of errors. Concerning classical serotyping we also referred to Hendriksen et al. 2009¹³.

Typing Method	Advantage	Drawback	Main reasons for Errors	How to address sources of errors
Serotyping	Directly determines phenotype	No typing of rough strains possible	Lack of experience with serotyping	Intensively trained staff
	Well established method	Requires high quality antisera		Quality control mechanism
SeqSero	Classification analogous to classical serotyping	Genotype may not correspond to phenotype due to undetected mutations	Low sequence data quality	Quality control mechanism, e.g. of sequencing process
	No assembly required	High quality sequencing data required (e.g. coverage, contamination)	Monophasic variants are only determined by lack of <i>fljB</i>	Improve detection method for monophasic variants
	Can be automated			
MLST-based typing	Provides phylogenetic information	High quality sequencing data required (e.g. coverage, contamination)	Low sequence data quality	Quality control mechanism, e.g. of sequencing process
	Can be automated	Assembly recommended		

494

495

496



