**Supplementary Materials**

**<u>Study Population</u>**

The National Center for Biotechnology Information (NCBI) 'pathogen detection' repository (https://www.ncbi.nlm.nih.gov/pathogens) *Salmonella* metadata was first downloaded (13 October 2017) and explored to identify non-typhoidal *Salmonella* (NTS) isolates recovered in the U.S. between 2006 and 2017 and for which data on the 'Run number' and the 'serotype' was available. For isolates with only the 'Biosample' information, 'Run numbers' were retrieved (in batches) separately from the Sequence Read Archive (SRA)-NCBI repository, when available.

A local dictionary of the serotype and equivalent antigenic formula for NTS serotypes was created based on 'SNOMED Clinical Terms' at the national center for biomedical ontology (NCBO) website (http://bioportal.bioontology.org/ontologies/; SNOMED classes: Organism> Microorganism> Prokaryote> Bacteria> *Salmonella*). The dictionary was used to unify the serotype naming in the metadata.

Overall, 28,494 sequences of NTS isolates belonging to 281 serotypes were identified.

Focusing on serotypes of public-health impact, the data was filtered to include only serotypes that had at least 100 isolates and at least one isolate was retrieved from humans. The filtered selection included 25,897 sequences of isolates belonging to 37 serotypes.

Collection sources were categorized into five categories: human, bovine (animal and food products), poultry (animal and food products), porcine (animal and food products) and other (including environmental, wildlife, domestic animals and non-available sources).

Categorization was based on the NCBI metadata in columns 'Host', 'Host disease',

‘Isolation source’ (Included in **Table S1**). In addition, 3,376 samples isolated from stool or other body secretions that were sequenced by the Enteric Diseases Laboratory Branch (EDLB) of the Centers for Disease Control and Prevention (CDC), were categorized as ‘Human’ source unless indicated otherwise.

**Data Analysis**

**Data quality**

For each serotype, paired-end Illumina reads were downloaded from NCBI’s ftp server (75 isolates were excluded at this step as their raw reads were not available at NCBI sequence read archive (SRA) repository) and reads quality was assessed using FASTQC v0.11.6 [1]. De-novo genome assemblies were conducted using ‘SPAdes’ *de novo* assembler v3.12.0 [2] with the "careful" option in order to reduce short indels and minimize the number of mismatches in the final assembly. When required, the ‘repair.sh’ command in BBmap v38.06 [3] was used to fix disordered raw reads before reassembly. The assemblies N50 were calculated using QUAST v4.6.3 [4] and assemblies with N50 lower than 30,000 base pairs (bp; n=188) were excluded from the analysis. The *Salmonella* In Silico Typing Resource (SISTR) v1.0.2 [5] was used, and only sequences with predicted serotypes (by both the serover antigen and cgmlst) that were in agreement with the serotype as defined in the metadata were further analyzed (1,497 isolates were excluded from the analysis in this step). The serotype Typhimurium var. -5 (i.e. Copenhagen) was predicted by SISTR as *S*. Typhimurium and therefore, in this case, the SISTR output was only used to verify that *S*. Copenhagen isolates were identified as *S*. Typhimurium and not any other serotype.

The average coverage depth of the assemblies was estimated after completion of the analysis (i.e. after removal of genetic duplicates, see below). For this purpose, reads were

46  aligned to the contig assemblies using bowtie2 v2.3.4.1 [6] and BBmap [3] was used to

47  calculate the average coverage depth of the contigs. In 18,211/18,282 (99.6%) of the

48  sequences the average coverage was at least 20, in 39/18,282 (0.21%) of the sequences the

49  average coverage was below 20 and in 33/18,282 (0.18%) of the sequences the average

50  coverage could not be calculated as the raw reads were not available for downloading from

51  NCBI anymore. The sequences with the low (below 20) and unknown average coverage

52  were not removed from the analysis, however given the amount of isolates included in

53  these groups (0.39% of all sequences analyzed here) it is not likely that their inclusion has

54  affected the analysis outcomes.

55  ***Genetic analyses***

56  Serotype phylogenies reconstruction

57  Before using 'FastTree' v2.1.10 [7] for phylogeny reconstruction of each serotype,

58  genetically identical duplicates were removed from the analysis [overall 5,855 sequences

59  were removed; up to 27 duplicates were removed in 50% of the serotypes and the highest

60  and lowest number of duplicates were removed in *S*. Enteritidis (n=3,042) and *S*.

61  Johannesburg (n=3)].

62  Packages 'ape' v5.0 [8] and 'ggtree' v1.10.5 [9] in R software v3.4.3 [10] were used for

63  visualization. In the final (FastTree) phylogenetic trees, 195 sequences were removed to

64  improve the visualization. These sequences (outliers) included individual sequences or

65  small groups of similar sequences (up to five sequences) demonstrating high dissimilarity

66  to other sequences within the serotype (i.e. longer tree branches), yet neither low quality of

67  the data nor misidentification of the serotype was evident. Up to three outliers were

68  identified in 50% of the serotypes; none were found in nine serotypes and the highest

69    number of outlier sequences was found in *S*. Javiana (n=32). The outlier sequences were

70    analyzed as part of each serotype analysis, yet were removed from the final trees to

71    improve the visualization of the genetic subpopulations within each serotype.

72    A scheme of the pipeline used in the analyses: metadata filtration, data quality assessment

73    and genetic analyses is illustrated in **figure S1**.

74    *Data interpretation*

75    <u>Comparison between subtyping by serotypes and genetic subpopulations</u>

76    Data of the genetic characteristics [i.e. presence of acquired antimicrobial resistance genes

77    (AARGs), multilocus sequence types (MLST) and plasmid replicons) were summarized for

78    each serotype and its genetic subpopulations (**Table S2**). For this purpose, the serotypes

79    data was stratified by genetic subpopulations. Then, for each antibiotic class [i.e. beta-

80    lactams; aminoglycosides; folate pathway inhibitors; tetracyclines; macrolide and

81    lincosamide (including sulphonamide and trimethoprim); quinolones; phenicols; and others

82    (including colistin, fosfomycin, fusidic acid, glycopeptide, nitroimidazole, oxazolidinone

83    and rifampicin)] the number of sequences harboring AARGs and the percentage of

84    sequences harboring AARGs within the genetic subpopulation were calculated. In addition,

85    similar information was summarized for predominant AARGs (i.e. found in at least 10% of

86    the sequences within the subpopulation) in each antibiotic class. However, for beta lactams

87    and quinolones, the data was summarized for all AARGs that were found. The MLST and

88    plasmid replicons were summarized for predominant types only (i.e. found in at least 10%

89    of the sequences within the subpopulation).

In addition, to described above and in the text, the following packages in R software v3.4.3 [10] were used: (i) for data manipulation and summarization – dplyr v0.8.0.1 [11], lettercase v0.13.1 [12], stringr v1.2.0 [13], tidyr v0.8.0 [14] and xlsx v0.5.7 [15]; for creating and formatting the figures - cowplot v0.9.2 [16], Hmisc v4.1.1 [17], ggplot2 v2.2.1 [18], ggpubr v0.2 [19] and gridExtra v2.3 [20].

**References**

1.  **Andrews S**. FastQC: a quality control tool for high throughput sequence data. 2010.
    http://www.bioinformatics.babraham.ac.uk/projects/fastqc.

2.  **Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M et al.** SPAdes: a new genome
    assembly algorithm and its applications to single-cell sequencing. *J Comput Biol*
    2012;19(5):455-477.

3.  **Bushnell B**. BBMap short-read aligner, and other bioinformatics tools.
    http://sourceforge.net/projects/bbmap/.

4.  **Gurevich A, Saveliev V, Vyahhi N, Tesler G**. QUAST: quality assessment tool for genome
    assemblies. *Bioinformatics* 2013;29(8):1072-1075.

5.  **Yoshida CE, Kruczkiewicz P, Laing CR, Lingohr EJ, Gannon VP et al.** The *Salmonella* In Silico
    Typing Resource (SISTR): An Open Web-Accessible Tool for Rapidly Typing and Subtyping
    Draft *Salmonella* Genome Assemblies. *PLoS One* 2016;11(1):e0147101.

6.  **Langmead B, Salzberg SL**. Fast gapped-read alignment with Bowtie 2. *Nat Methods*
    2012;9(4):357-359.

7.  **Price MN, Dehal PS, Arkin AP**. FastTree 2--approximately maximum-likelihood trees for
    large alignments. *PLoS One* 2010;5(3):e9490.

113      8.   **Paradis E, Claude J, Strimmer K**. APE: Analyses of Phylogenetics and Evolution in R

114           language. *Bioinformatics* 2004;20(2):289-290.

115      9.   **Yu G, Smith DK, Zhu H, Guan Y, Lam TTY**. ggtree: an R package for visualization and

116           annotation of phylogenetic trees with their covariates and other associated data. *Methods*

117           *in Ecology and Evolution* 2017;8(1):28-36.

118      10.  **R Core Team**. R: A language and environment for statistical computing. Vienna, Austria;

119           2016. https://www.R-project.org/.

120      11.  **Wickham H, Francois, R., Henry, L., Müller, K.** dplyr: A Grammar of Data Manipulation. R

121           package version 0.7.4. 2017. https://CRAN.R-project.org/package=dplyr.

122      12.  **Brown C**. lettercase: Utilities for Formatting Strings with Consistent Capitalization, Word

123           Breaks and White Space version 0.13.1. 2016. https://CRAN.R-

124           project.org/package=lettercase.

125      13.  **Wickham H**. 2015. stringr: Simple, Consistent Wrappers for Common String Operations.

126           https://CRAN.R-project.org/package=stringr [accessed.

127      14.  **Wickham H, Henry, L. ,** . tidyr: Easily Tidy Data with 'spread()' and 'gather()' Functions. R

128           package version 0.8.0. 2018. https://CRAN.R-project.org/package=tidyr.

129      15.  **Dragulescu AA**. xlsx: Read, write, format Excel 2007 and Excel 97/2000/XP/2003 files. R

130           package version 0.5.7. 2014. https://CRAN.R-project.org/package=xlsx.

131      16.  **Wilke CO**. cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2' version

132           0.9.2. 2017. https://CRAN.R-project.org/package=cowplot.

133      17.  **Frank E**. Hmisc: Harrell Miscellaneous. R package version 4.1-1. https://CRAN.R-

134           project.org/package=Hmisc; 2018.

135      18.  **Wickham H**. *Ggplot2 : elegant graphics for data analysis*. New York: Springer; 2009.

136    19.    **Kassambara A**. ggpubr: 'ggplot2' Based Publication Ready Plots. R package version 0.2.

137           2018. https://CRAN.R-project.org/package=ggpubr.

138    20.    **Auguie B**. gridExtra: Miscellaneous Functions for "Grid" Graphics. R package version 2.3.

139           2017. https://CRAN.R-project.org/package=gridExtra.

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

**Supplementary figure captions**

**Figure S1 –** A scheme of the pipeline used in the analyses: metadata filtration (yellow box), data quality assessment (green box) and genetic analyses (purple box). Infographics were created using 'draw.io' website. Notes: [a] average coverage was estimated for the final 18,282 sequences only (see Supplementary text); [b] outliers were only removed from the final phylogenetic trees figures (see Supplementary text); [c] the final number of sequences analyzed after exclusion of duplicates (but not the outliers; see Supplementary text).

**Figure S2 –** The percentage of sequences that were obtained from each source are presented for each genetic subpopulation within serotype. Barplots were generated for each of the 37 serotypes included in the study. Bars are colored according to the source: human (purple), bovine (blue), poultry (brown), porcine (orange) and other (grey). The number of sequences included in each subpopulation/serotype are indicated above each bar column. The 'Total' column includes the sum of all subpopulations within serotype in addition to sequences (of that serotype) that were not grouped in any of the subpopulations (see supplementary material text for details).

**Figure S3 –** The percentage of sequences that were obtained from each collection period in each genetic subpopulation within serotype. Barplots were generated for each of the 37 serotypes included in the study. Bars are colored according to the collection period: 2006-2009 (blue), 2010-2013 (green) and 2014-2017 (red). The number of sequences included in each subpopulation/serotype are indicated above each bar column. The 'Total' column includes the sum of all subpopulations within serotype in addition to sequences (of that serotype) that were not grouped in any of the subpopulations (see supplementary material text for details).

178 **Figure S4 –** A maximum likelihood phylogenetic tree was reconstructed with RAxML

179 using the single nucleotide polymorphisms (SNPs) found in the core genome of

180 representative sequences from all 37 serotypes (n=370). Ten sequences were selected from

181 each serotype phylogeny to represent the diversity of the genetic subpopulations. The tree

182 was rooted using *S*. Paratyphi type A outgroup (SRR3033248, SRR3277289; not included

183 in the figure). Bootstrap replicates (n=5,000) were used for branch support. Tree tips were

184 colored according to the serotype and the genetic subpopulation number was indicated in

185 the tip name.

186 **Figure S5 –** The percentage of plasmid size groups found in each genetic subpopulation

187 within serotype. Barplots were generated for each of the 37 serotypes included in the study.

188 Bars are colored according to the plasmid size groups: 'small'- up to 6kbp (Purple);

189 'intermediate'- between 6kbp and 100kbp (yellow); and 'large'- more than 100kbp (dark

190 red). The number of sequences in each subpopulation/serotype are indicated above each

191 bar column. The 'Total' column includes the sum of all subpopulations within serotype

192 (see supplementary material text for details).

193

194

195

196

197

198

199    **Figure S1**



**Metadata Filtration**

NCBI - "Pathogen detection"
Salmonella WGS
n=88,048

- Collected before 2006
- Not of U.S. origin
- No human source/ less than 100 WGS available

U.S. Non Typhoid Salmonella
of public health importance 2006-2017
n=25,897

**Data Quality**

- Raw reads not available (n=75)

Genome assemblies (**Spades**)
n=25,822

Read correction
(**BBmap**)

Assembly quality (**Quast**)

- N50 below 30,000 (n=188)

- Map reads to contigs (**bowtie2**)
- Calculate avg. coverage (**BBmap**)

- Average coverage below 20[a]

In Silico Serotyping (**SISTR**)

- Serotyping mismatch (n=1,497)

High quality genome assemblies
n=24,137

**Genetic Analysis**

Detection of AARGs, MLST and
plasmid replicon types
(**CGE - 'bacterial analysis pipeline'**)
n=18,282[c]

Genome Annotation
(**Prokka**)

For selected WGS from
all serotypes (n=370)

For each serotype
separately

Data Analysis

Phylogeny trees
Annotation

Core genome analyses
(**Roary**)

SNP variable sites extraction
(**Snp-sites**)

- Duplicates removal
(n=5,855)

- Outliers removal[b]
(n=195)

Maximum likelihood (ML)
phylogeny tree (**RAxML**)
n=370

Approximate ML phylogeny
trees (**FastTree**)
n=18,087; 37 serotypes

200

## Figure S2

203    **Figure S3**



204                                                                    **Subpopulations**

12

205    **Figure S4**



**Serotypes**

| | |
|---|---|
| ■ AGONA | ■ MELEAGRIDIS |
| ■ ANATUM | ■ MONOPHASIC |
| ■ BAREILLY | ■ MONTEVIDEO |
| ■ BERTA | ■ MUENCHEN |
| ■ BRAENDERUP | ■ MUENSTER |
| ■ CERRO | ■ NEWPORT |
| ■ COPENHAGEN | ■ NORWICH |
| ■ DERBY | ■ OHIO |
| ■ DUBLIN | ■ ORANIENBURG |
| ■ ENTERITIDIS | ■ POONA |
| ■ HADAR | ■ READING |
| ■ HARTFORD | ■ SAINTPAUL |
| ■ HEIDELBERG | ■ SCHWARZENGRUND |
| ■ INFANTIS | ■ SENFTENBERG |
| ■ JAVIANA | ■ TENNESSEE |
| ■ JOHANNESBURG | ■ THOMPSON |
| ■ KENTUCKY | ■ TYPHIMURIUM |
| ■ LONDON | ■ UGANDA |
| ■ MBANDAKA | |

0.03

206

13

207    **Figure S5**



208    Subpopulations

14