

1 Multi-layered networks of SalmoNet2 2 enable strain comparisons of the 3 *Salmonella* genus on a molecular level

4
5
6

7 Marton Olbei ^{1,2}, Balazs Bohar ^{1,3}, David Fazekas ^{1,3}, Matthew Madgwick ^{1,2}, Padhmanand
8 Sudhakar ^{1,2,4}, Isabelle Hautefort ¹, Aline Métris ², Jozsef Baranyi ^{2,5}, Robert A. Kingsley ^{2,6},
9 Tamas Korcsmaros ^{1,2}

10 ¹ Earlham Institute, Norwich, UK

11 ² Quadram Institute Biosciences, Norwich, UK

12 ³ Eotvos Lorand University, Budapest, Hungary

13 ⁴ KU Leuven, Leuven, Belgium

14 ⁵ Institute of Nutrition, University of Debrecen, Debrecen, Hungary

15 ⁶ School of Biological Sciences, University of East Anglia, Norwich, UK

16

17 Abstract

18 Serovars of the genus *Salmonella* primarily evolved as gastrointestinal pathogens in a wide
19 range of hosts. Some serotypes later evolved further, adopting a more invasive lifestyle in a
20 narrower host range associated with systemic infections. A system-level knowledge of these
21 pathogens has the potential to identify the complex adaptations associated with the evolution
22 of serovars with distinct pathogenicity, host range and risk to human health. This promises to
23 aid the design of interventions and serve as a knowledge base in the *Salmonella* research
24 community. Here we present SalmoNet2, a major update to SalmoNet, the first multi-layered
25 interaction resource for *Salmonella* strains, containing protein-protein, transcriptional
26 regulatory and enzyme enzyme interactions. The new version extends the number of
27 *Salmonella* genomes from 11 to 20, including strains such as *S. Typhimurium* D23580, an
28 epidemic multidrug-resistant strain leading to invasive non-typhoidal *Salmonella* Disease

29 (iNTS), and a strain from *Salmonella bongori*, another species in the *Salmonella* genus. The
30 database now uses strain specific metabolic models instead of a generalised model to
31 highlight differences between strains. This has increased the coverage of high-quality protein-
32 protein interactions, and enhances interoperability with other computational resources by
33 adopting standardised formats. The resource website has been updated with tutorials to help
34 researchers analyse their *Salmonella* data using molecular interaction networks from
35 SalmoNet2. SalmoNet2 is accessible at <http://salmonet.org/>.

36

37 Importance

38 Multi-layered network databases collate information from multiple sources, and are powerful
39 both as a knowledge base and platform for analysis. Here we present SalmoNet2, an
40 integrated network resource of 20 *Salmonella* strains, containing protein-protein,
41 transcriptional regulatory, and metabolic interactions. Key improvements to the update
42 include expanding the number of strains, strain-specific metabolic networks, an increase in
43 high quality protein-protein interactions, community standard computational formats to help
44 interoperability, and online tutorials to help users analyse their data using SalmoNet2.

45

46 Introduction

47

48 Serovars of the genus *Salmonella* are enteric pathogens, capable of causing a self-limiting
49 gastrointestinal inflammatory disease in a variety of animals. The host species, depending on
50 the *Salmonella* subspecies, range from cold-blooded vertebrates to humans. *Salmonella*
51 infection is one of the most common foodborne or waterborne illnesses resulting in
52 approximately 94 million illnesses, and 155,000 deaths each year¹⁻³.

53

54 Of six subspecies of *Salmonella enterica*, a small number of subspecies I serovars have
55 adapted to cause an invasive infection in a restricted host range, instead of a self-limiting
56 gastrointestinal inflammation typical of *Salmonella* serovars. These extraintestinal
57 *Salmonella* strains, including the typhoidal strains that are human adapted, emerged on
58 multiple occasions independently. A hallmark of adaptation is genomic and phenotypic

59 changes, including loss of function mutations in genes related to adaptation to specific niches
60 in their host commonly affecting anaerobic metabolism, virulence genes, chemotaxis or
61 motility ⁴.

62

63 *S. Typhi* is an ancient pathogen and the most common extraintestinal *Salmonella* serovar
64 affecting humans. Over the past decades, another group of *Salmonella* appeared as one of the
65 most commonly isolated pathogens from the blood of patients, particularly in sub-Saharan
66 Africa ⁵. The invasive nontyphoidal *Salmonella* (iNTS) strains, in common with *S. Typhi* and
67 *S. Paratyphi*, cause a systemic infection. Unlike *S. Typhi*, iNTS commonly affects
68 immunocompromised individuals or young children, leading to bacteremia and meningitis.
69 iNTS is most often caused by specific genotypes of *S. Typhimurium* and *S. Enteritidis* that
70 are distinct from genotypes of these serovars commonly associated with gastrointestinal
71 infections outside of sub-Saharan Africa ⁶⁻⁸.

72 The *Salmonella* genus contains pathogens with diverse host range and pathogenicity, and
73 dissecting the specific differences between gastrointestinal and extraintestinal strains have
74 been pursued by a multitude of means ^{4,9,10}. Previously, we constructed SalmoNet, a multi-
75 layered network resource for 10 *Salmonella* serovars that integrated protein-protein,
76 regulatory and metabolic information ¹¹. With its multi-layered networks SalmoNet can serve
77 as a knowledge base for the community and aid in understanding *Salmonella* pathogenesis
78 and evolution by mapping the differences in molecular interactions between *Salmonella*
79 pathovars on multiple biological levels. This systems level information allows researchers to
80 enhance the information content of their own studies, by adding interaction context to the
81 changes observed on a genomic or transcriptome level.

82

83 Here we present SalmoNet2, an update to the first public multi-layered network resource for
84 *Salmonella* research. The new version extends the coverage of strains from 11 to 20,
85 including an important iNTS strain, and strains outside of subspecies *enterica*, from
86 subspecies *arizonae* and *Salmonella bongori*. To aid interoperability in computational
87 biology the database adopted the PSI-MI TAB format, and is now accessible through the
88 NDEx network repository. In addition, we show how rewiring of the network information can
89 be utilised by the research community to understand aspects of *Salmonella* evolution, the
90 step-by-step workflows of which are now accessible through tutorials on our website.

91

92

93

94 Results

95 SalmoNet2 extends out of subspecies I

96 SalmoNet2 adds 9 new multi-layered networks of *Salmonella* strains in the database
97 compared to the first version. Included amongst others are commonly used laboratory strains,
98 additional extraintestinal strains, including *S. Typhimurium* strain D23580, a well
99 characterised pathogen associated with invasive non-typhoidal *Salmonella* (iNTS) disease, a
100 strain of *Salmonella bongori*, and a member of a different subspecies (subsp. *arizonae*)
101 within *Salmonella enterica*. The extended coverage captures a larger variety of the
102 *Salmonella* genus, and for the first time provides interaction networks for strains from outside
103 of subspecies *enterica* (Supplementary Table I). To define the phylogenetic relationship of
104 the strains included in the database we constructed a neighbour-joining tree from variation in
105 the core genome nucleotide sequence, and compared this with hierarchical classification trees
106 based on matrix representation of protein-protein, regulatory and metabolic networks (Figure
107 1).

108

109 ***Figure 1. Core genome SNP based phylogenetic tree, and hierarchical classification of***
110 ***network layers. Extraintestinal (EI) serovars labelled with red, gastrointestinal (GI) serovars***
111 ***with blue labels. A., Neighbour-joining tree from core genome SNPs of the strains. B-D.,***
112 ***Hierarchical classification trees based on matrix representation of protein-protein,***
113 ***regulatory and metabolic networks. The five letter labels encode for the names of the***
114 ***different strains (for details of the encoding please refer to Supplementary Table 1).***

115

116 The topology was in accordance with previously published phylogenies¹² with no clear
117 clustering of extraintestinal and gastrointestinal serovars in the phylogenetic tree. This is
118 consistent with observations in previous works in the literature, where the extraintestinal and
119 gastrointestinal strains could not be distinguished based on genomic dendrograms, and
120 consistent with the independent emergence of extraintestinal serovars from gastrointestinal
121 serovars, through a convergent evolutionary process^{13,14}.

122 SalmoNet2 increases the information content of the individual 123 network layers

124 We included a number of methodological improvements to the workflow of the Salmonet1
125 database, leading to an increased number of high quality interactions for the individual
126 network layers. To increase the coverage of the protein-protein interactions without
127 compromising quality, we have used the IntAct MIscore when extrapolating orthologous
128 interaction information from the IntAct database ¹⁵. Instead of relying on one experimental
129 method as in the first version, using the MIscore as a quality filter permitted extending the
130 number of available high-quality protein-protein interactions that we could use to establish
131 orthologous protein-protein interactions from the commensal bacteria *Escherichia coli*
132 (Supplementary Figure 1).

133

134 By utilising a strain-specific genome-scale metabolic model for each strain developed
135 previously (*Seif et al.*), instead of a general model (*Thiele et al.*), the metabolic layer now
136 includes more enzyme-enzyme relationships, where two proteins are connected if a
137 metabolite produced by one is a substrate for another ¹⁶⁻¹⁸, leading to a more complete
138 description of the metabolic capabilities of the strains. The information content of Position-
139 Specific Scoring Matrices (PSSMs) that are required to carry out genome-wide regulatory
140 scans were enhanced with novel binding sites published since the first version of the
141 database, and from new data uploaded to the CollecTF repository ¹⁹. The total number of
142 interactions has been increased from 81,514 to 270,215, primarily due to the expansion of the
143 PPI layer, and the increase in the number of involved strains. The composition of the
144 consensus network, comprised of shared interactions amongst all strains included in the
145 database, slightly changed from the first version of SalmoNet, indicating the shifts caused by
146 the updated data sources and expanded strain repertoire. 24.4% of regulatory interactions (up
147 from 16%), 68.1% of PPI interactions (down from 72%), and 51.8% (down from 69%) of
148 metabolic interactions were shared amongst all strains, forming the core network of
149 *Salmonella* interactions. Figure 2 shows the changes in the size of the networks and
150 individual layers compared to the first version.

151

152 **Figure 2. Comparison of SalmoNet2 with the first version.** A: main data sources and
153 interactions in SalmoNet2. B: comparison of network size in SalmoNet 1 and SalmoNet2. C:
154 comparison of layer size in terms of participating nodes. D: comparison of layer size in terms

155 *of interactions between SalmoNet and SalmoNet2. The five letter codes encoding for the*
156 *different strains can be found in Supplementary Table 1.*

157

158 **Novel formats improve interoperability**

159 In addition to the previously used formats (.csv, .cys.), we extended the output format data to
160 help computational biologists access network information in SalmoNet2. We now provide
161 networks in the community standard PSI-MITAB format as well, which contains a strictly
162 regulated vocabulary for interaction data, helping interoperability between network resources,
163 a prerequisite for inclusion in the PSICQUIC ecosystem ²⁰. Using standardised formats
164 improves the interoperability with other network information repositories, and provides space
165 to maximise each interaction with as much data as possible, in a controlled and transparent
166 manner ²⁰. This further strengthened the information content of the database, and improved
167 the potential use cases beyond network analysis. To enable the networks to be directly
168 accessible from the widely used Cytoscape network analysis program, we have also
169 deposited them to the NDEx network repository ²¹.

170 **Website enhancements for a user-friendly experience**

171 The SalmoNet website was enhanced compared to the previous version. We now carry new
172 locus tag identifiers for all *Salmonella* strains to enable users to map their experimental data
173 to the SalmoNet2 interaction networks. As part of our shift to OMA as the source of
174 orthology for *Salmonella* proteins, SalmoNet2 now directly links to the respective OMA
175 pages and sequence data instead of Uniprot ²². Where possible, Uniprot data is still accessible
176 through OMA ²³.

177 During the lifecycle of SalmoNet1, we identified a bottleneck with our putative users. The
178 interaction network format, while potentially useful for scientists with a microbiology
179 background, proved difficult to use, which led to potentially less user retention. To combat
180 this, we have created a new tab on the website containing tutorials as an introduction to
181 network analysis using the SalmoNet2 database. These tutorials enable analyses shown in this
182 article. We plan to add additional tutorials, workflows and examples to the website in the
183 future, to further increase the usability of the platform.

184 Case study: network rewiring to identify functional differences in 185 *Salmonella enterica*

186 Network rewiring entails many approaches aimed at quantifying changes between interaction
187 networks, and has been used to identify differences between interaction networks^{24,25}. In this
188 work, we compared the degree of interaction rewiring between the interactomes of four host
189 adapted typhoidal *Salmonella* strains and four gastrointestinal *Salmonella* strains to explore
190 the utility of a multi-layered network resource such as SalmoNet2. We compared the most
191 rewired subgraphs of the two types of strains to find the causes of the rewiring.

192

193 In general, the most rewired nodes were global regulators, such as Crp, Fis and Fur. The
194 significantly enriched functions are similar between the compared strains, with a few key
195 differences. For example, the ferric uptake regulator Fur senses metal concentration and
196 redox state of cells, and regulates many operons and genes involved in these processes²⁶.
197 Interestingly, Fur is enriched in the GO function “iron ion homeostasis” in all included
198 gastrointestinal strains, while this enrichment is absent from the typhoidal strains. Upon
199 further inspection of the genes responsible for the enrichment of the term and their
200 orthologous status, Fur is missing interactions present in GI strains towards the genes *fhuA*,
201 *fhuE*, caused by the disruption of coding sequences in these genes in typhoidal serovars, as
202 highlighted previously in the literature^{13,27}. Similarly, Fur is enriched in the term “cell
203 adhesion” in all gastrointestinal strains, whereas this function is not enriched in typhoidal
204 strains, except *S. Paratyphi C*. Inspection of the genes underlying the enrichment result
205 revealed that the culprit behind the mismatch in functional enrichment is the
206 pseudogenization and subsequent missing interactions with the genes *stiH* and *stiA* in the rest
207 of the typhoidal *Salmonella* strains, two genes responsible for the production of fimbriae,
208 highlighted previously in the literature¹³. From the top 50 most rewired nodes, on average 33
209 nodes had at least one pseudogene first neighbour in the typhoidal serovars, and on average
210 4% of the first neighbours of the top 50 most rewired nodes were pseudogenes. In the
211 gastrointestinal strains, on average 7 nodes had pseudogene first neighbours, and only 1% of
212 their first neighbours were pseudogenes.

213 While a large part of the rewiring was due to gene loss in typhoidal and extraintestinal
214 serovars, we found examples where the cause of rewiring was due to the exclusivity of genes
215 to the extraintestinal group. Two proteins, YreP and YjcS, are present in all typhoidal and
216 extraintestinal strains of *Salmonella* included in SalmoNet2. However, they are missing from

217 all gastrointestinal strains bar one. The protein YjcS has an orthologue in *S. Enteritidis*, but
218 the protein is otherwise missing from the gastrointestinal group. The genes share an upstream
219 regulatory region, and are predicted to interact with the regulators HilC, RtsA and Fur. The
220 *yreP* and *yjcS* genes were first described together in *Escherichia coli*, in two analysed strains:
221 *E. coli* SMS-3-5, an environmental pathogenic isolate with multiple antibiotic resistances,
222 and *E. coli* (NMEC) O7:K1 strain CE10, causing neonatal meningitis. The first gene, *yreP*
223 (*dgcY* in *E. coli*), encodes a putative diguanylate cyclase, based on the presence of a GGDEF
224 domain^{28,29}. Diguanylate-cyclases facilitate the production of c-di-GMP, a ubiquitous
225 secondary messenger metabolite in prokaryotes^{28,29}. The second gene, *yjcS* (EcSMS35_1714
226 in *E. coli*), is an alkyl-sulfatase. This enzyme has been first described in *Pseudomonas spp.*,
227 where a strain carrying this enzyme was able to grow on the surfactant sodium dodecyl
228 sulphate (SDS), and the gene has been characterised in *E. coli* as well^{30,31}.

229

230 After noting their presence in the extraintestinal strains included in SalmoNet2, we expanded
231 the search into a more expansive data source, pubMLST, to see if this split was representative
232 of the serovars as a whole, and not just the specific strains in SalmoNet2³². Figure 3 shows
233 the results of the BLAST searches in the pubMLST database.

234

235 **Figure 3. Prevalence of the *yreP* + regulatory region + *yjcS* segment in *Salmonella***
236 ***serovars based on BLAST hits. The top 10 serovars have been described previously as***
237 ***sources of invasive illness.***

238

239 In total 83% of BLAST hits come from well-known extraintestinal serovars, dominated by *S.*
240 Typhi strains (Figure 3). The top 10 serovars in terms of number of hits are mostly invasive
241 serovars: *S. Typhi*, *S. Paratyphi A*, and *S. Paratyphi C* are notable typhoidal serovars adapted
242 to humans, *S. Dublin*, *S. Pullorum* and *S. Choleraesuis* are well-known host adapted serovars
243 of cattle, poultry and pigs^{4,11}. *S. Napoli* is an emerging serovar in Europe, phylogenetically
244 closely related to *S. Paratyphi A*, carrying an almost identical pattern of typhoid-associated
245 genes, and capable of causing a form of invasive non-typhoidal disease^{33,34}. The invasive
246 behaviour is not as clear cut with the rest of the serovars, but there have been several reports
247 of it: *S. Bovismorbificans* is capable of causing bloodstream infections, and has recently been
248 described as an emerging disease in Malawi, converging towards a phenotype resembling a
249 human adapted iNTS variant³⁵. Although not strictly an extraintestinal serovar, *S. Virchow*
250 has been known to cause invasive illness³⁶⁻³⁹. *S. Weltevreden* is an emerging cause of

251 diarrheal and sometimes invasive disease in humans in tropical regions, and may be adapted
252 to life in aquatic hosts^{40,41}. While large in total numbers in the database, *S. Enteritidis* only
253 makes up 2% of the positive hits. Since *S. Enteritidis* is one of the most commonly isolated
254 iNTS strains, there exists a possible link to invasive behaviour^{42,43}. However, more work is
255 needed to uncover whether the two proteins are beneficial to an extraintestinal lifestyle.
256 This brief case study highlights how the information contained in and linked with SalmoNet2
257 can be used to form scientific questions relating the functionality of genes to the behaviour
258 and phylogenetics of *Salmonella*, based on molecular interaction information. SalmoNet 2
259 contains example strains from the most prevalent serovars, and the information can further be
260 extended using the easily accessible sequence data and homology information through OMA
261 and other computational resources.

262

263

264

265 Discussion

266 By increasing the number of strains to 20 from the previous 11, SalmoNet now extends out of
267 subspecies I., adding information on members of other subspecies (subspecies *arizonae*), or
268 an entirely different species (*Salmonella bongori*). Developing a more compatible structure
269 between SalmoNet, and other available large- scale evolutionary genomics tools such as
270 OMA, there is increased potential to generate interaction networks for specific *Salmonella*
271 strains on request, or build similar data resources for other non-model organisms⁴⁴. With the
272 change to OMA as the backbone of SalmoNet interactions, there is a great potential to study
273 the evolutionary history of proteins, and interactions. The on-demand availability of
274 orthologous proteins from outside of the studied organism or clade can make larger scope
275 comparisons possible⁴⁵.

276

277 The programmatic access interfaces implemented into OMA make these integrated analyses
278 reproducible and scalable⁴⁶. Orthology mapping is the most computationally intensive step of
279 the SalmoNet workflow. The OMA standalone software can save a lot of time and resources
280 here, since the all-against-all Smith-Waterman sequence alignments can be parallelised, both
281 on single computers or high-performance clusters²². Adding a new strain or species in the
282 future is also made easier, as OMA Standalone does not require an all-against-all

283 recomputing of the orthologous relationships in these cases, as pre-computed results can be
284 submitted, in which case only the new genomes require computation time. Using OMA is not
285 only beneficial for the orthology mapping, it is also helpful for the annotation work. The first
286 version of SalmoNet was essentially UniProt based, with UniProt IDs serving as the primary
287 identifiers of the database. Currently not all proteins of all strains have a matching UniProt
288 ID, hence the OMA IDs as our new primary identifier.

289
290 The availability of strain-specific metabolic models, and the increased specificity of PPI data,
291 although still reliant on orthology mapping, increases the resolution of the resulting network
292 models for non-model organisms, and the more interwoven interaction layers get, the more
293 valuable the information content of the database gets. Although there are other resources
294 containing *Salmonella* interaction data, such as STRING for PPI interactions, RegPrecise for
295 regulatory interactions, or BioCyc for metabolic interactions, no other freely available
296 resource combines the listed connection types besides SalmoNet, for multiple *Salmonella*
297 strains⁴⁷⁻⁴⁹.

298
299 SalmoNet2 enables the network analyses as shown with the rewiring analysis. It highlights
300 how the information contained in and linked with SalmoNet2 can be used to inform scientific
301 questions such as relating the functionality of genes to phenotypes and phylogenetics of
302 *Salmonella*, based on molecular interaction information. SalmoNet2 contains example strains
303 from the most prevalent serovars, and the information can further be extended using the
304 easily accessible sequence data and homology information through OMA and other
305 computational resources as shown with the pubMLSt example.

306
307 To increase the usability and interoperability of the generated interaction information, we
308 now utilise the PSI-MITAB format as well, quickly becoming a standard of biological
309 network information^{20,50}. To have the networks be directly accessible from the widely used
310 Cytoscape network analysis program, the NDEX network repository can host them separately
311 from the SalmoNet website, making them directly available for end-users²¹. Beyond their
312 raw information content, databases are as good as their usability and availability, and the
313 potential for SalmoNet2 data to be found and utilised in as many ways as possible is crucial
314 for this effort to be useful for the scientific community⁵¹. To further enhance the accessibility
315 of SalmoNet2 data we wrote detailed step-by-step tutorials describing the computational

316 steps required to perform analyses such as the comparisons involving the gastrointestinal and
317 typhoidal strains above.

318

319 Methods

320 Updated orthology mapping tool

321 Although the main structure of the database remained the same, the underlying workflow
322 changed. As in the first version, we mapped the orthologous proteins across the included
323 strains. In SalmoNet 1 this was done by InParanoid, a well-established tool for this process ⁵².
324 In this update we used the OMA standalone software to construct these relationships,
325 including the available *Salmonella* strains from the OMA browser database. OMA is a large-
326 scale orthology database and toolkit, containing the orthology information and protein
327 sequence data needed for SalmoNet in one place, including the proteomes and genomes of
328 the strains on request, and important annotation data ⁵³.

329 It is important to note, that the outputs of the tools can be slightly different: according to a
330 study comparing these methods the OMA standalone output OMA groups lead to a generally
331 more precise, but also strict mapping, leading to less false positives (and true positives as
332 well) ⁵⁴. We did however get very similar, and in cases better recall than we did in SalmoNet
333 1.0 (between 69-75% overlap with the 4140 proteins from *E. coli*; Supplementary Table 1)
334 using InParanoid.

335

336 Updated and novel data sources

337 Protein-Protein Interaction Networks

338 The construction of the protein-protein interaction (PPI) network follows the same essential
339 steps it did in the first version of the database, collected from multiple databases ⁵⁵⁻⁵⁸. To
340 increase the coverage of the included PPIs without losing quality, we have used the IntAct
341 PSI-MIscore (> 0.50) when importing interaction information from the IntAct database,
342 instead of relying on one experimental method, as in the first version (psi-mi:"MI:0096"(pull
343 down)). Supplementary Figure 1 shows the distribution of the IntAct PSI-MIscores.

344 Metabolic Networks

345 SalmoNet2 uses new, strain specific genome-scale metabolic models for *Salmonella*^{16,17}. The
346 models used the same STM 1.0 model as a starting point SalmoNet1 did¹⁸, but updated it
347 with new genes and reactions, and were made strain specific, leading to the metabolic models
348 of 410 *Salmonella* strains belonging to 64 serovars. Otherwise, the workflow remained
349 identical, resulting in enzyme-enzyme interactions, where two proteins are connected if a
350 metabolite produced by one is a substrate for another⁵⁹. Similarly, as in the first version, we
351 have excluded links connected by metabolites partaking in more than 10 reactions⁵⁹.

352

353 Regulatory Networks

354 The establishment of the transcriptional regulatory networks was done in an identical way to
355 SalmoNet 1. Supplementary Figure 2 shows the workflow for the construction of the
356 regulatory layer. The core of the network, the manually curated interactions, high-throughput
357 data (ChIP-Seq), and low-throughput, experimentally verified interactions and data sources
358 remained the same. The information content of Position-Specific Scoring Matrices (PSSMs)
359 used to carry out the genome-wide scans was enhanced with novel binding sites published
360 since the first version of the database, from new data uploaded to the CollecTF repository¹⁹.
361 RSAT's consensus tool is no longer available on the web server, info-gibbs took its place,
362 which is the tool that was used to construct the matrices. Similarly, as previously, RSAT
363 retrieve-sequence was used to gather the putative promoter regions for the genomes included
364 in SalmoNet, and matrix-scan was used to establish putative transcription factor - target gene
365 (promoter region) pairs⁶⁰.

366

367 Removal of pseudogenes

368 To remove all hypothetically disrupted coding DNA sequences (HDCs), the curation made by
369 Nuccio & Bäumlér was used to remove such entries¹³ and^{61,62} were used to remove them
370 from *S. Typhimurium* D23580.

371

372 Network rewiring

373 To calculate network rewiring we used the DyNet app in Cytoscape to calculate the rewiring
374 value of the nodes in each group separately⁶³. Four typhoidal strains (*S. Paratyphi* A (AKU
375 1261), *S. Paratyphi* A (ATCC 9150), *S. Paratyphi* C (RKS4594), *S. Typhi* (Ty2) and four
376 gastrointestinal strains (*S. Agona* (SL483), *S. Newport* (SL254), *S. Heidelberg* and *S.*
377 *Typhimurium* (LT2)) were compared for interaction differences.

378

379 The level of rewiring was calculated across all strains, and the degree-corrected rewiring
380 values were ordered in a descending list, where the top 50 hits were further analysed. To
381 calculate the enrichment of Gene Ontology terms in the identified subgraphs the the up-to-
382 date Gene Ontology annotation of the target genes was downloaded using the topGO library
383 in R, and following that the R library clusterProfiler was used to calculate Gene Ontology
384 enrichment with the enricher function, from Biological Process terms^{64,65}. P-value
385 adjustment for multiple testing was carried out with the Benjamini-Hochberg approach, using
386 the p.adjust function in R.

387 The statistically significant enrichment results were compared side-by-side between the
388 groups, and the differences in enrichment were further studied by comparing the sets of genes
389 responsible for (underlying) the enriched terms, i.e., if one group was enriched in a specific
390 term, the presence/absence of the orthologous genes responsible for the enrichment was
391 analysed in the members of the other group.

392

393 To study the relationship of YreP and YjcS to the extraintestinal pathovar, network rewiring
394 was calculated in an identical manner as above, but all extraintestinal and gastrointestinal
395 strains from SalmoNet2 were involved in the comparisons. BLAST searches for the *yreP* and
396 *yjcS* genes was done through the pubMLST website, with default parameters³². The entire
397 genomic sequence of the genes and their shared regulatory region was queried, as taken from
398 *S. Gallinarum* strain 287/91 (see Supplementary File 1). The hits were filtered for above 95%
399 sequence identity, and the top 10% of bitscores to make sure the compared sequences contain
400 both the genes and the shared regulatory region.

401

402

403 Data availability

404 The data generated for this study is available at the database website, <http://salmonet.org>.

405

406 Acknowledgements

407 The work of MO, PS, IH, TK were supported by the UKRI BBSRC Gut Microbes and Health

408 Institute Strategic Programme BB/R012490/1 and its constituent projects

409 BBS/E/F/000PR10353 and BBS/E/ F/000PR10355. MO, BB, DF, PS, IH, TK were also

410 supported by a BBSRC Core Strategic Programme Grant for Genomes to Food Security

411 (BB/CSP1720/1) and its constituent work packages, BBS/E/T/000PR9819 and

412 BBS/E/T/000PR9817. PS was supported by the European Research Council Advanced Grant

413 (ERC-2015-AdG, 694679, CrUCCial). MO and MM were supported by a BBSRC - Norwich

414 Research Park Biosciences Doctoral Training Partnership grant (BB/M011216/1 and

415 BB/S50743X/1). RK was supported by the UKRI Institute Strategic Programme Microbes in

416 the Food Chain BB/R012504/1 and its constituent project(s) BBS/E/F/000PR10348 and

417 BBS/E/F/000PR10349.

418

419

420

421

422 Bibliography

423 1. Coburn, B., Grassl, G. A. & Finlay, B. B. Salmonella, the host and disease: a brief

424 review. *Immunol. Cell Biol.* **85**, 112–118 (2007).

425 2. Hohmann, E. L. Nontyphoidal salmonellosis. *Clin. Infect. Dis.* **32**, 263–269 (2001).

426 3. Majowicz, S. E. *et al.* The global burden of nontyphoidal Salmonella gastroenteritis.

427 *Clin. Infect. Dis.* **50**, 882–889 (2010).

- 428 4. Tanner, J. R. & Kingsley, R. A. Evolution of Salmonella within Hosts. *Trends*
429 *Microbiol.* **26**, 986–998 (2018).
- 430 5. Tsai, C. N. & Coombes, B. K. Emergence of invasive Salmonella in Africa. *Nat.*
431 *Microbiol.* (2021) doi:10.1038/s41564-021-00864-5.
- 432 6. Gilchrist, J. J. & MacLennan, C. A. Invasive nontyphoidal salmonella disease in africa.
433 *Ecosal Plus* **8**, (2019).
- 434 7. Feasey, N. A., Dougan, G., Kingsley, R. A., Heyderman, R. S. & Gordon, M. A.
435 Invasive non-typhoidal salmonella disease: an emerging and neglected tropical disease in
436 Africa. *Lancet* **379**, 2489–2499 (2012).
- 437 8. Okoro, C. K. *et al.* Intracontinental spread of human invasive Salmonella Typhimurium
438 pathovariants in sub-Saharan Africa. *Nat. Genet.* **44**, 1215–1221 (2012).
- 439 9. Perez-Sepulveda, B. M. & Hinton, J. C. D. Functional transcriptomics for bacterial gene
440 detectives. *Microbiol. Spectr.* **6**, (2018).
- 441 10. Langridge, G. C., Nair, S. & Wain, J. Nontyphoidalsalmonella serovars cause different
442 degrees of invasive disease globally. *J. Infect. Dis.* **199**, 602–603 (2009).
- 443 11. Métris, A. *et al.* SalmoNet, an integrated network of ten Salmonella enterica strains
444 reveals common and distinct pathways to host adaptation. *NPJ Syst. Biol. Appl.* **3**, 31
445 (2017).
- 446 12. Branchu, P., Bawn, M. & Kingsley, R. A. Genome Variation and Molecular
447 Epidemiology of Salmonella enterica Serovar Typhimurium Pathovariants. *Infect.*
448 *Immun.* **86**, (2018).
- 449 13. Nuccio, S.-P. & Bäumler, A. J. Comparative analysis of Salmonella genomes identifies a
450 metabolic network for escalating growth in the inflamed gut. *MBio* **5**, e00929-14 (2014).
- 451 14. Timme, R. E. *et al.* Phylogenetic diversity of the enteric pathogen Salmonella enterica
452 subsp. enterica inferred from genome-wide reference-free SNP characters. *Genome Biol.*

- 453 *Evol.* **5**, 2109–2123 (2013).
- 454 15. Villaveces, J. M. *et al.* Merging and scoring molecular interactions utilising existing
455 community standards: tools, use-cases and a case study. *Database (Oxford)* **2015**,
456 (2015).
- 457 16. Seif, Y., Monk, J. M., Machado, H., Kavvas, E. & Palsson, B. O. Systems Biology and
458 Pangenome of Salmonella O-Antigens. *MBio* **10**, (2019).
- 459 17. Seif, Y. *et al.* Genome-scale metabolic reconstructions of multiple Salmonella strains
460 reveal serovar-specific metabolic traits. *Nat. Commun.* **9**, 3771 (2018).
- 461 18. Thiele, I. *et al.* A community effort towards a knowledge-base and mathematical model
462 of the human pathogen Salmonella Typhimurium LT2. *BMC Syst. Biol.* **5**, 8 (2011).
- 463 19. Kılıç, S. *et al.* From data repositories to submission portals: rethinking the role of
464 domain-specific databases in CollecTF. *Database (Oxford)* **2016**, (2016).
- 465 20. Perfetto, L. *et al.* CausalTAB: the PSI-MITAB 2.8 updated format for signalling data
466 representation and dissemination. *Bioinformatics* **35**, 3779–3785 (2019).
- 467 21. Pillich, R. T., Chen, J., Rynkov, V., Welker, D. & Pratt, D. Ndex: A community resource
468 for sharing and publishing of biological networks. *Methods Mol. Biol.* **1558**, 271–301
469 (2017).
- 470 22. Altenhoff, A. M. *et al.* The OMA orthology database in 2018: retrieving evolutionary
471 relationships among all domains of life through richer web and programmatic interfaces.
472 *Nucleic Acids Res.* **46**, D477–D485 (2018).
- 473 23. The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic Acids*
474 *Res.* **45**, D158–D169 (2017).
- 475 24. Mehta, T. K. *et al.* Evolution of regulatory networks associated with traits under
476 selection in cichlids. *Genome Biol.* **22**, 25 (2021).
- 477 25. Treveil, A. *et al.* Regulatory network analysis of Paneth cell and goblet cell enriched gut

- 478 organoids using transcriptomics approaches. *Mol. Omics* **16**, 39–58 (2020).
- 479 26. Troxell, B., Fink, R. C., Porwollik, S., McClelland, M. & Hassan, H. M. The Fur regulon
480 in anaerobically grown *Salmonella enterica* sv. Typhimurium: identification of new Fur
481 targets. *BMC Microbiol.* **11**, 236 (2011).
- 482 27. Wang, Y. *et al.* Evolution and sequence diversity of fhua in salmonella and escherichia.
483 *Infect. Immun.* **86**, (2018).
- 484 28. Povolotsky, T. L. & Hengge, R. Genome-Based Comparison of Cyclic Di-GMP
485 Signaling in Pathogenic and Commensal *Escherichia coli* Strains. *J. Bacteriol.* **198**, 111–
486 126 (2016).
- 487 29. Ryjenkov, D. A., Tarutina, M., Moskvina, O. V. & Gomelsky, M. Cyclic diguanylate is a
488 ubiquitous signaling molecule in bacteria: insights into biochemistry of the GGDEF
489 protein domain. *J. Bacteriol.* **187**, 1792–1798 (2005).
- 490 30. Liang, Y., Gao, Z., Dong, Y. & Liu, Q. Structural and functional analysis show that the
491 *Escherichia coli* uncharacterized protein YjcS is likely an alkylsulfatase. *Protein Sci.* **23**,
492 1442–1450 (2014).
- 493 31. Williams, J. & Payne, W. J. Enzymes induced in a bacterium by growth on sodium
494 dodecyl sulfate. *Appl. Microbiol.* **12**, 360–362 (1964).
- 495 32. Jolley, K. A., Bray, J. E. & Maiden, M. C. J. Open-access bacterial population genomics:
496 BIGSdb software, the PubMLST.org website and their applications. [version 1; peer
497 review: 2 approved]. *Wellcome Open Res.* **3**, 124 (2018).
- 498 33. Gori, M. *et al.* High-resolution diffusion pattern of human infections by *Salmonella*
499 *enterica* serovar Napoli in Northern Italy explained through phylogeography. *PLoS ONE*
500 **13**, e0202573 (2018).
- 501 34. Huedo, P. *et al.* *Salmonella enterica* Serotype Napoli is the First Cause of Invasive
502 Nontyphoidal Salmonellosis in Lombardy, Italy (2010-2014), and Belongs to Typhi

- 503 Subclade. *Foodborne Pathog. Dis.* **14**, 148–151 (2017).
- 504 35. Bronowski, C. *et al.* Genomic characterisation of invasive non-typhoidal Salmonella
505 enterica Subspecies enterica Serovar Bovismorbificans isolates from Malawi. *PLoS*
506 *Negl. Trop. Dis.* **7**, e2557 (2013).
- 507 36. Eckerle, I., Zimmermann, S., Kapaun, A. & Junghanss, T. Salmonella enterica serovar
508 Virchow bacteremia presenting as typhoid-like illness in an immunocompetent patient. *J.*
509 *Clin. Microbiol.* **48**, 2643–2644 (2010).
- 510 37. Mani, V., Brennand, J. & Mandal, B. K. Invasive illness with Salmonella virchow
511 infection. *Br. Med. J.* **2**, 143–144 (1974).
- 512 38. Messer, R. D., Warnock, T. H., Heazlewood, R. J. & Hanna, J. N. Salmonella meningitis
513 in children in far north Queensland. *J. Paediatr. Child Health* **33**, 535–538 (1997).
- 514 39. Todd, W. T. & Murdoch, J. M. Salmonella virchow: a cause of significant bloodstream
515 invasion. *Scott. Med. J.* **28**, 176–178 (1983).
- 516 40. Hounmanou, Y. M. G. *et al.* Molecular characteristics and zoonotic potential of
517 salmonella weltevreden from cultured shrimp and tilapia in vietnam and china. *Front.*
518 *Microbiol.* **11**, 1985 (2020).
- 519 41. Makendi, C. *et al.* A Phylogenetic and Phenotypic Analysis of Salmonella enterica
520 Serovar Weltevreden, an Emerging Agent of Diarrheal Disease in Tropical Regions.
521 *PLoS Negl. Trop. Dis.* **10**, e0004446 (2016).
- 522 42. Feasey, N. A. *et al.* Distinct Salmonella Enteritidis lineages associated with enterocolitis
523 in high-income settings and invasive disease in low-income settings. *Nat. Genet.* **48**,
524 1211–1217 (2016).
- 525 43. Gordon, M. A. Invasive nontyphoidal Salmonella disease: epidemiology, pathogenesis
526 and diagnosis. *Curr. Opin. Infect. Dis.* **24**, 484–489 (2011).
- 527 44. Olbei, M., Kingsley, R. A., Korcsmaros, T. & Sudhakar, P. Network Biology

- 528 Approaches to Identify Molecular and Systems-Level Differences Between Salmonella
529 Pathovars. *Methods Mol. Biol.* **1918**, 265–273 (2019).
- 530 45. Demeter, A. *et al.* ULK1 and ULK2 are less redundant than previously thought:
531 computational analysis uncovers distinct regulation and functions of these autophagy
532 induction proteins. *Sci. Rep.* **10**, 10940 (2020).
- 533 46. Kaleb, K., Warwick Vesztröcy, A., Altenhoff, A. & Dessimoz, C. Expanding the
534 Orthologous Matrix (OMA) programmatic interfaces: REST API and the *OmaDB*
535 packages for R and Python [version 2; peer review: 2 approved]. *F1000Res.* **8**, 42
536 (2019).
- 537 47. Caspi, R. *et al.* BioCyc: A Genomic and Metabolic Web Portal with Multiple Omics
538 Analytical Tools. *The FASEB Journal* (2019).
- 539 48. Novichkov, P. S. *et al.* RegPrecise 3.0--a resource for genome-scale exploration of
540 transcriptional regulation in bacteria. *BMC Genomics* **14**, 745 (2013).
- 541 49. Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with increased
542 coverage, supporting functional discovery in genome-wide experimental datasets.
543 *Nucleic Acids Res.* **47**, D607–D613 (2019).
- 544 50. Kerrien, S. *et al.* IntAct--open source resource for molecular interaction data. *Nucleic*
545 *Acids Res.* **35**, D561-5 (2007).
- 546 51. Merali, Z. & Giles, J. Databases in peril. *Nature* **435**, 1010–1011 (2005).
- 547 52. O'Brien, K. P., Remm, M. & Sonnhammer, E. L. L. Inparanoid: a comprehensive
548 database of eukaryotic orthologs. *Nucleic Acids Res.* **33**, D476-80 (2005).
- 549 53. Altenhoff, A. M. *et al.* OMA orthology in 2021: website overhaul, conserved isoforms,
550 ancestral gene order and more. *Nucleic Acids Res.* (2020) doi:10.1093/nar/gkaa1007.
- 551 54. Altenhoff, A. M. *et al.* Standardized benchmarking in the quest for orthologs. *Nat.*
552 *Methods* **13**, 425–430 (2016).

- 553 55. Calderone, A., Castagnoli, L. & Cesareni, G. mentha: a resource for browsing integrated
554 protein-interaction networks. *Nat. Methods* **10**, 690–691 (2013).
- 555 56. Mosca, R., Céol, A. & Aloy, P. Interactome3D: adding structural details to protein
556 networks. *Nat. Methods* **10**, 47–53 (2013).
- 557 57. Orchard, S. *et al.* The MIntAct project - IntAct as a common curation platform for 11
558 molecular interaction databases. *Nucleic Acids Res.* **42**, D358-63 (2014).
- 559 58. Oughtred, R. *et al.* The BioGRID interaction database: 2019 update. *Nucleic Acids Res.*
560 **47**, D529–D541 (2019).
- 561 59. Kreimer, A., Borenstein, E., Gophna, U. & Ruppin, E. The evolution of modularity in
562 bacterial metabolic networks. *Proc Natl Acad Sci USA* **105**, 6976–6981 (2008).
- 563 60. Nguyen, N. T. T. *et al.* RSAT 2018: regulatory sequence analysis tools 20th anniversary.
564 *Nucleic Acids Res.* **46**, W209–W214 (2018).
- 565 61. Kingsley, R. A. *et al.* Epidemic multiple drug resistant Salmonella Typhimurium causing
566 invasive disease in sub-Saharan Africa have a distinct genotype. *Genome Res.* **19**, 2279–
567 2287 (2009).
- 568 62. Canals, R. *et al.* Adding function to the genome of African Salmonella Typhimurium
569 ST313 strain D23580. *PLoS Biol.* **17**, e3000059 (2019).
- 570 63. Goenawan, I. H., Bryan, K. & Lynn, D. J. DyNet: visualization and analysis of dynamic
571 molecular interaction networks. *Bioinformatics* **32**, 2713–2715 (2016).
- 572 64. Alexa, A. & Rahnenfuhrer, J. topGO: Enrichment Analysis for Gene Ontology.
573 <https://www.bioconductor.org/packages/release/bioc/html/topGO.html> (2021).
- 574 65. Wu, T. *et al.* clusterProfiler 4.0: A universal enrichment tool for interpreting omics data.
575 *Innovation (N Y)* **2**, 100141 (2021).

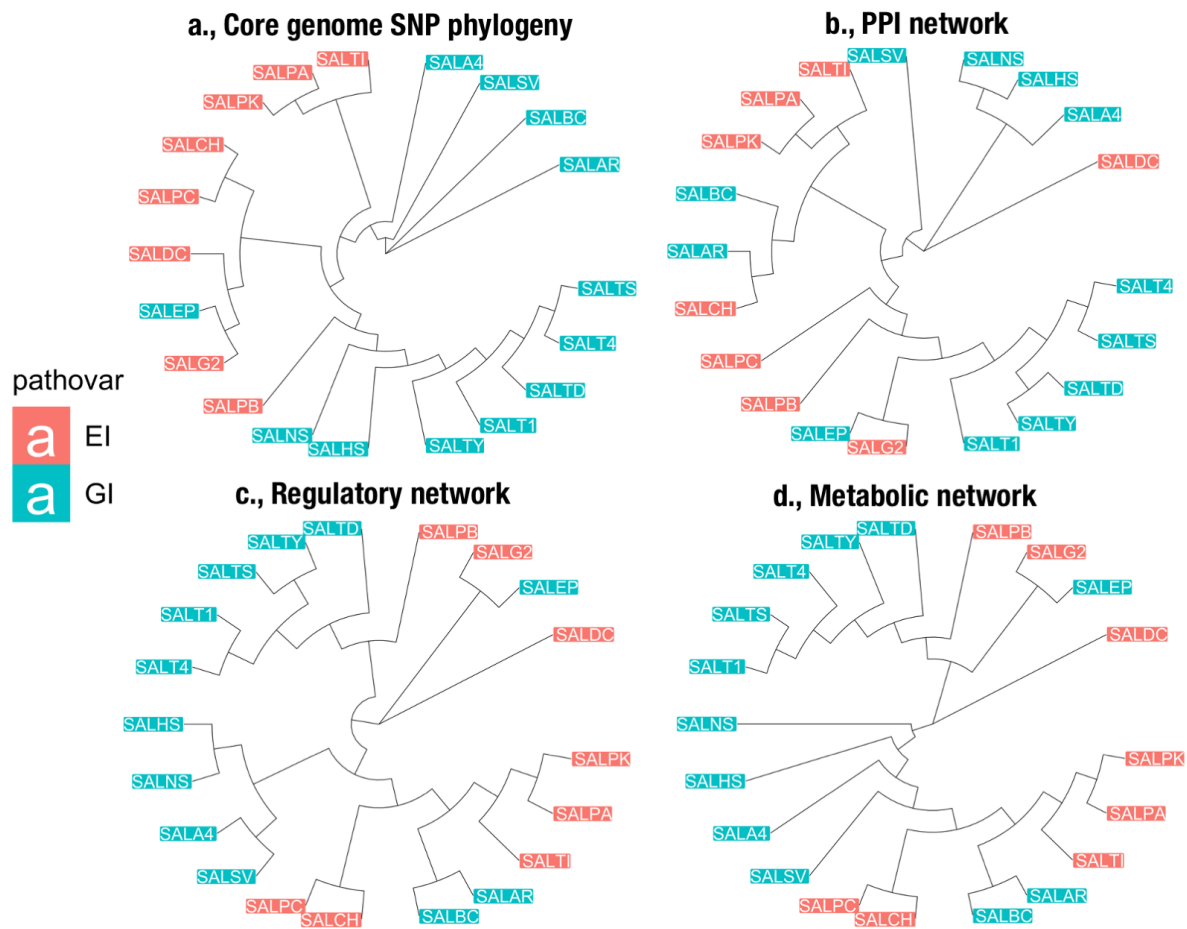


Figure 1. Core genome SNP based phylogenetic tree, and hierarchical classification of network layers. Extraintestinal (EI) serovars labelled with red, gastrointestinal (GI) serovars with blue labels. A., Neighbour-joining tree from core genome SNPs of the strains. B-D., Hierarchical classification trees based on matrix representation of protein-protein, regulatory and metabolic networks. The five letter labels encode for the names of the different strains (for details of the encoding please refer to Supplementary Table 1).

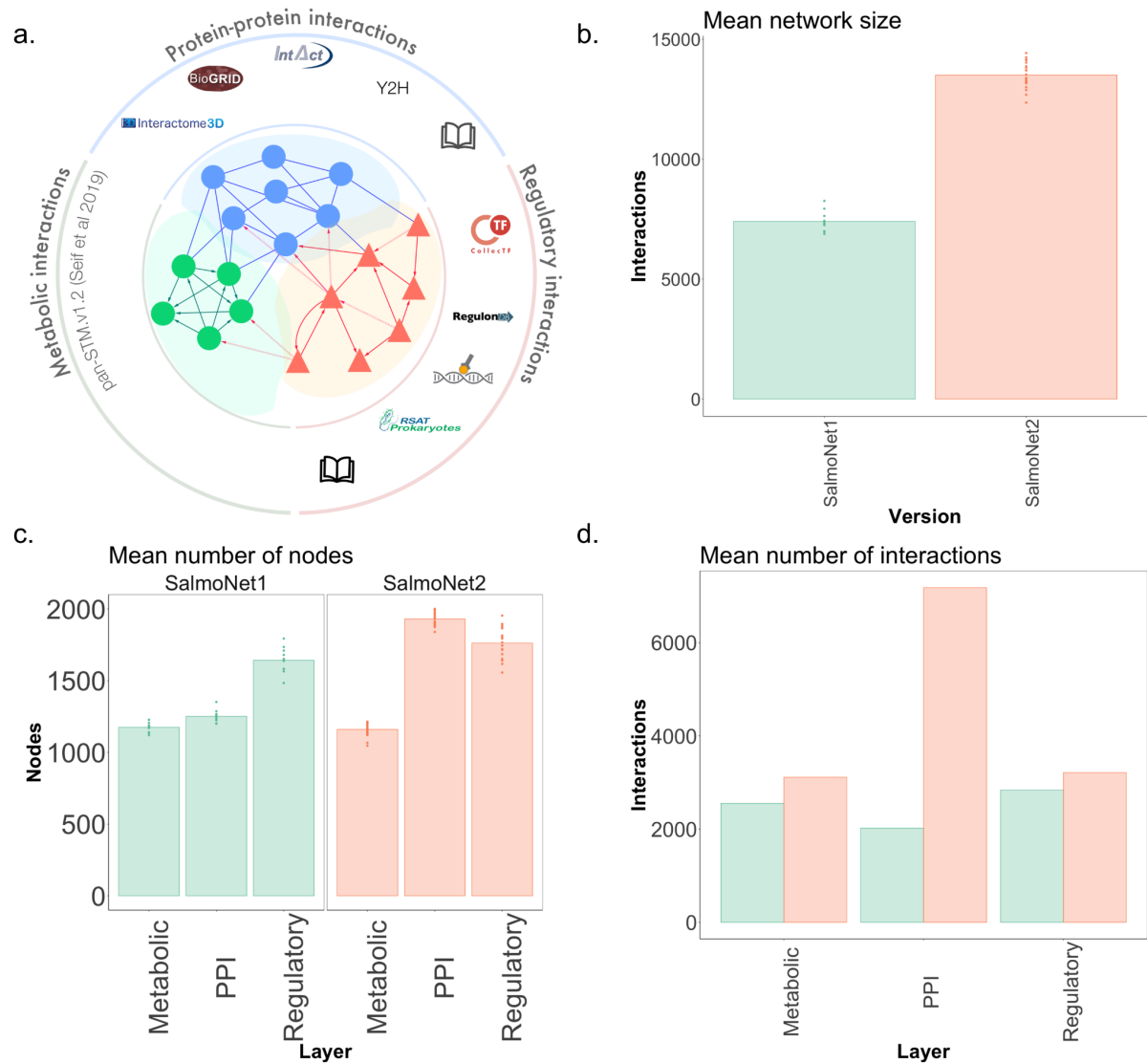


Figure 2. Comparison of SalmoNet2 with the first version. *A: main data sources and interactions in SalmoNet2. B: comparison of network size in SalmoNet 1 and SalmoNet2. C: comparison of layer size in terms of participating nodes. D: comparison of layer size in terms of interactions between SalmoNet and SalmoNet2. The five letter codes encoding for the different strains can be found in Supplementary Table 1.*

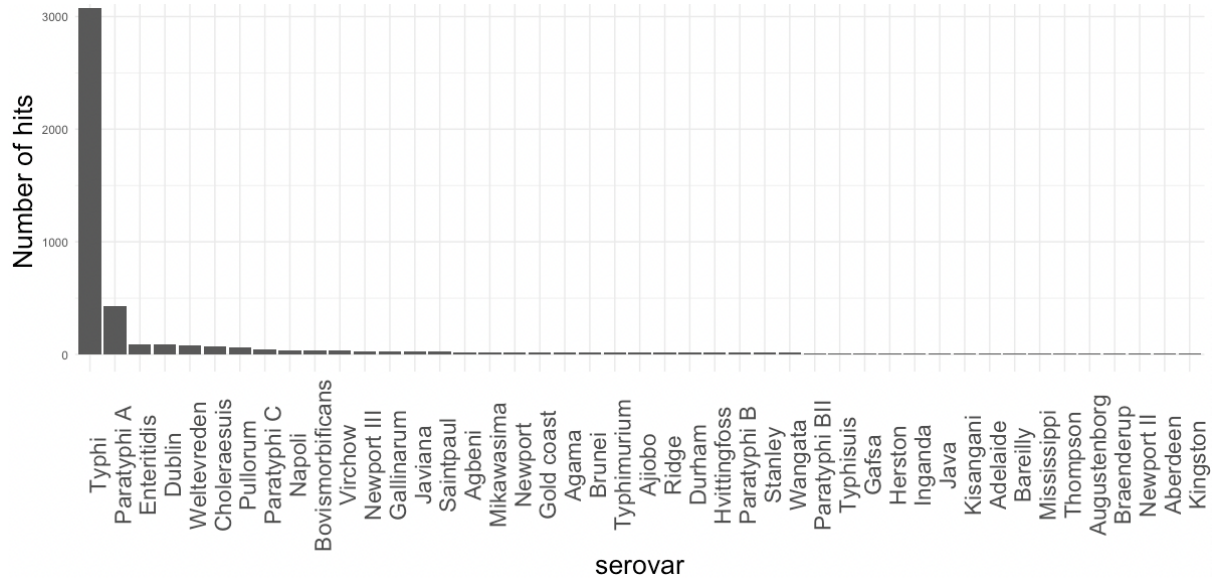


Figure 3. Prevalence of the *yreP* + regulatory region + *yjcS* segment in *Salmonella* serovars based on BLAST hits. The top 10 serovars have been described previously as sources of invasive illness. Serovars containing < 5 isolates were removed from this figure for clarity.