# Machine learning identifies signatures of host adaptation

# in the bacterial pathogen *Salmonella enterica*

Nicole E. Wheeler[1,2,5], Paul P. Gardner[2,3], Lars Barquist[4,5]

1. The Wellcome Trust Sanger Institute, Hinxton, United Kingdom.
2. Biomolecular Interaction Centre, School of Biological Sciences, University of Canterbury, Christchurch, New Zealand.
3. Department of Biochemistry, University of Otago, Dunedin, New Zealand.
4. Institute for Molecular Infection Biology, University of Wuerzburg, Wuerzburg, Germany.
5. Correspondence to: nw17@sanger.ac.uk; lars.barquist@uni-wuerzburg.de

## Abstract

Emerging pathogens are a major threat to public health, however understanding how pathogens adapt to new niches remains a challenge. New methods are urgently required to provide functional insights into pathogens from the massive genomic data sets now being generated from routine pathogen surveillance for epidemiological purposes. Here we integrate a method for scoring the functional impact of mutations with a random forest classifier, and apply this to the classification of *Salmonella enterica* strains associated with extraintestinal disease. Members of the species fall along a continuum, from pathovars which cause gastrointestinal infection and low mortality, associated with a broad host-range, to those that cause invasive infection and high mortality, associated with a narrowed host range. By training our random forest classifier to discriminate gastrointestinal and invasive serovars of *Salmonella*, using a small and well-characterised training dataset, we are able to additionally discriminate recently emerged *Salmonella* Enteritidis and Typhimurium lineages associated with invasive disease in immunocompromised populations in sub-Saharan Africa. Importantly, our classifier produces interpretable lists of gene variants associated with extraintestinal disease. This approach accurately identifies patterns of gene degradation specific to invasive serovars that have been captured by more labour-intensive investigations, but can be readily scaled to larger analyses.

## Introduction

Understanding how bacteria adapt to new niches and hosts and thus emerge or re-emerge as a cause of infectious disease in human and animals is of critical importance to anticipating and preventing epidemic disease (Frank and Schmid-Hempel 2008; Fauci and Morens 2012). With the decreasing cost of genome sequencing, comparative genomics has become a rich source of insight into the origins and movement of bacteria in new pathogenic niches. However, translating whole genome sequence databases into mechanistic and functional insights remains a challenge.

38

39   Early expectations were that pathogen evolution would be driven primarily by the acquisition

40   of virulence factors. However, as whole-genome sequencing has become increasingly

41   routine, a decidedly more complex picture has emerged (Pallen and Wren 2007; Loman and

42   Pallen 2015). A pattern of bacterial entrance to a new niche followed by adaptation through

43   the loss of antivirulence loci and reduced metabolic flexibility is now recognised as a

44   paradigm of the emergence of important human pathogens from non-pathogenic bacterial

45   species (McNally et al. 2016; The et al. 2016; Merhej et al. 2013; Reuter et al. 2014). These

46   new niches can be the result of virulence factor acquisition providing access to a previously

47   inaccessible niche in a so-called foothold moment (Reuter et al. 2014), or the emergence of

48   new host niches driven by chronic disease (Marvig et al. 2015; Klemm et al. 2016; Feasey et

49   al. 2012). While pathogen and host requirements for infection vary, there is increasing

50   evidence of parallel evolution in bacteria adapting to the same or similar host niche. This is

51   perhaps nowhere more evident than in the species *Salmonella enterica*.

52

53   *Salmonella enterica* strains that cause disease in warm-blooded mammals lie on a spectrum

54   from those that have a broad host range and cause self-limiting gastrointestinal infection, to

55   those that are more restricted in host range, but cause systemic disease and are typically

56   associated with higher mortality (Rabsch et al. 2002; Feasey et al. 2012). Host-restricted,

57   extraintestinal variants of *Salmonella enterica* have evolved independently multiple times

58   from gastrointestinal ancestors (Bäumler and Fang 2013), and show a greater degree of

59   gene degradation compared to their generalist relatives (Parkhill et al. 2001; McClelland et

60   al. 2004; Thomson et al. 2008). There are common patterns in the genes that undergo

61   pseudogenization in invasive *Salmonella*, most obviously an extensive network of genes

62   required for anaerobic metabolism in the inflamed host gut (Nuccio and Bäumler 2014;

63   Langridge et al. 2015), a pattern with parallels in other host-adapting enteropathogens

64   (McNally et al. 2016).

65

66   Identifying these signals of parallel evolution has been challenging, relying mainly on manual

67   annotation and comparison of pseudogenes (Nuccio and Bäumler 2014; Langridge et al.

68   2015). Detection of pseudogenes in particular relies on ad-hoc criteria to identify large

69   truncations, deletions, or frameshifts (Lerat and Ochman 2005; Kuo and Ochman 2010). It is

70   rare that the same genes or complete pathways are pseudogenized in host-adapted species;

71   rather interpretation has relied on identifying overrepresentation of independent

72   pseudogenization events clustered in certain pathways (Nuccio and Bäumler 2014). If

73   pseudogenization leads to pathway attenuation or inactivation, it seems likely that reduced

74   selective pressure will lead to a higher incidence of detrimental mutation fixation in other

75   genes in these pathways. Indeed, we have previously shown that functional variant calling,

76   based on sequence deviation from patterns of conservation observed in deep sequence

77   alignments, shows a similar functional signal in host-restricted *Salmonella enterica* serovar

78   Gallinarum to pseudogene analysis (Wheeler et al. 2016), identifying a larger cohort of

79   genes where constraints on drift appear to have been lifted during host-adaptation.

80

81   In previous work we developed DeltaBS, a profile hidden Markov model (HMM) based

82   approach to functional variant calling (Wheeler et al. 2016). The basic assumption of this

83   approach is that variation in conserved positions of a protein sequence is more likely to

84   affect protein function than variation in less conserved regions. This approach can integrate

85   information about nonsynonymous mutations, indels, and truncations. We have previously

86   shown that DeltaBS can successfully identify functional changes in genes that would be

87   missed by standard pseudogene analysis (Kingsley et al. 2013), and that a subset of genes

88   in host-adapted strains appear to accumulate large DeltaBS values (Wheeler et al. 2016).

89   Additionally, others have observed similar changes in DeltaBS distributions during

90   adaptation of *Salmonella* to a single immunocompromised host (Klemm et al. 2016).  We

91   generally assume that a large DeltaBS value is indicative of a decay in protein function. We

92   cannot rule out that a large DeltaBS may rather indicate a change in protein function, though

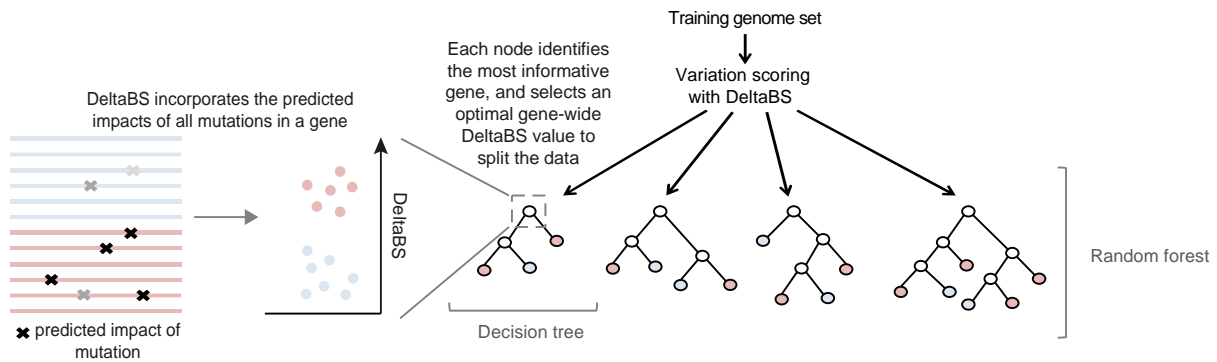93   we expect this to be relatively rare.

94

95   Here, we have leveraged these previous observations to identify signatures of mutational

96   burden consistent with adaptation to an invasive lifestyle. We have developed a random

97   forest classifier using delta bitscore (DeltaBS) functional variant calling (Wheeler et al. 2016)

98   that can perfectly separate intestinal *Salmonella* serovars from host-adapted, extraintestinal

99   serovars. We use random forest models because they perform well on datasets with few

100   informative variables (Dutilh et al. 2013; Pappu and Pardalos 2014), and have the potential

101   to detect functional relationships (i.e. epistasis) between genes with a decision tree structure

102   (Touw et al. 2013; Wei et al. 2014). They have been applied successfully in the past to

103   predict microbial phenotype using gene presence/absence data (Bayjanov et al. 2012), and

104   SNPs already known to be associated with phenotype (Laabei et al. 2014; Alam et al. 2014).

105   We show that these models produce interpretable signatures of host-adaptation, and

106   furthermore that these signatures can be detected in strains of *Salmonella* associated with

107   invasive disease in immunocompromised populations in sub-Saharan Africa.


108   **Results**


109   ***Constructing a random forest classifier for extraintestinal Salmonellae***

110   The approach taken in this investigation is summarised in Fig 1, and described below. We

111   built our model using a collection of genomes from well-characterised reference strains of

112   gastrointestinal and extraintestinal *Salmonella* serovars (Supplemental Table S1), drawing

113   on the extensive curation of orthology relationships performed by Nuccio and Bäumler

114   (2014). These strains were originally characterised as "gastrointestinal" or "extraintestinal"

115   based on common patterns of gene degradation, host restriction and clinical characteristics

116   observed among the extraintestinal strains (Nuccio and Bäumler 2014), and we have

117   employed this same categorisation our analysis. We scored the functional importance of

118   sequence variation by comparing the protein coding genes of each serovar to profile HMMs

119   from the eggNOG database (Huerta-Cepas et al. 2016), designed to capture patterns of

5

120    sequence variation typically seen in the protein coding genes of Gammaproteobacteria (see

121    Methods).
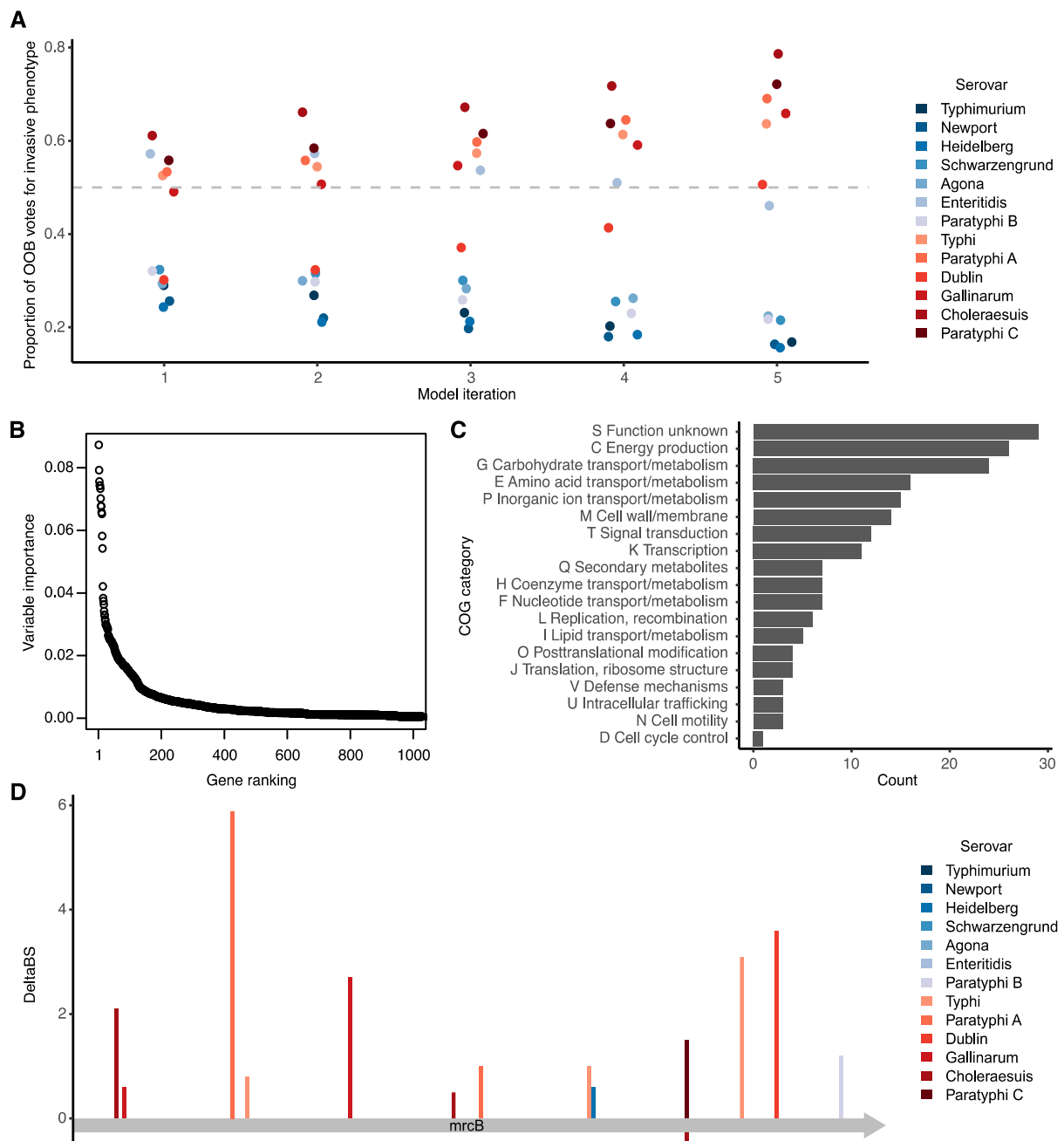


**Fig 1 | Overview of the approach employed in this study**

For each genome, the functional significance of sequence variation within protein coding

genes is quantified using the DeltaBS metric. Following scoring, a bootstrap sampling of

genomes are used to train each decision tree. For each node in the tree, a random subset of

genes are sampled, and the most informative gene from this set is chosen to split the data.

For each node in the tree, the predictive utility of the selected gene (variable importance) is

tested by calculating how well the gene separates the samples according to phenotype.


We then employed random forests to identify the genes which were most informative of

phenotype when viewed collectively. Random forests work by building an ensemble of

decision trees designed to predict a characteristic of the samples (Breiman 2001), in this

case adaptation to an extraintestinal, or invasive, niche. For each node in the decision tree,

the best gene of a random sampling from the training gene set is selected according to its

ability to separate a randomly selected subset of samples by phenotype based on DeltaBS

values. The process of building a random forest produces measures of variable importance

that can be used to assess the relative utility of different genes in classification of *Salmonella*

strains based on lifestyle.

6

140 **_A small subset of genes are strongly predictive of invasiveness in Salmonella_**

141 To obtain an indication of the proportion of the genome that shows patterns of unusual

142 sequence variation associated with an invasive phenotype, we trained a random forest

143 model on a set of 6,438 orthologous genes. Accuracy of the model was assessed using out-

144 of-bag accuracy. This out-of-bag (OOB) measure of accuracy gives us an indication of how

145 well each decision tree in the forest performs at predicting phenotype in a serovar it has

146 never encountered before, using information on DeltaBS differences collected from other

147 serovars. Next, we performed iterative feature selection to improve the performance of the

148 model. This process involved repeated rounds of selecting the top 50% of predictors and re-

149 training the model, until the model achieved perfect OOB predictive performance on the

150 training dataset (Fig 2A). When the full set of filtered orthologous genes was used to build a

151 model, a subset of genes ranked much higher than the others in variable importance (VI)

152 (Fig 2B). We then saw a tailing off of VI, resulting in 4,721 orthologous groups either not

153 being used in the model, or not improving classification accuracy (as indicated by VI ≤ 0).

154 The final model used 196 of the original 6,438 genes for prediction (Supplemental Table S2).

155 This model additionally achieved perfect classification accuracy on an independent set of

156 genomes of the same serovars as our training data (Supplemental Fig S1).

7

157

**Fig 2 | A subset of *Salmonella* genes are strongly indicative of invasive potential**

A: Out-of-bag votes for phenotype of each serovar cast by each model. Model 1 is the model

built using all predictor variables, then each successive model was built using sparsity

pruning from the previous model's predictor variables. Model 5 is the final model with 100%

accuracy. Out-of-bag votes include only those votes cast by trees that were not trained on a

given sample. The dashed grey line indicates the voting threshold to classify an isolate as

invasive. Invasive serovars are coloured in red and gastrointestinal serovars are coloured in

blue.

166    B: Of all genes used in the original training dataset, a small minority are given high

167    importance in identifying invasive strains. Variable importance is shown for the top 1000

168    genes used in the original training set. Variable importance was measured as average

169    decrease in Gini index in a random forest model trained on all orthologous groups that met

170    the inclusion criteria (N = 6,438).

171    C: Functional categories associated with the top predictive genes.

172    D: Mutations in *mrcB* (penicillin-binding protein 1b), one of the top three predictors.

173    Mutations in different strains are colour-coded, with bars in red indicating a mutation in an

174    extraintestinal strain and bars in blue indicating a mutation in a gastrointestinal strain. An

175    estimate of the effect of the mutation on protein function (DeltaBS) is shown on the y-axis,

176    with positive values indicating higher chance of a mutation being deleterious to protein

177    function. The x-axis represents the length of the protein.


178    ***Predictive genes are typically degraded or absent in invasive isolates***

179    We anticipated that the majority of informative genes identified in our study would be genes

180    that showed functional degradation in invasive isolates but not in gastrointestinal isolates. Of

181    the top predictors in our study (N = 196), 154 showed significantly greater mutational burden

182    in extraintestinal strains compared to gastrointestinal strains (Mann-Whitney U test, adjusted

183    *P*-value < 0.05), compared to 9 genes that showed significantly greater mutational burden in

184    gastrointestinal strains. Of the genes that were more conserved in invasive isolates, one was

185    the aldo-keto reductase *yakC*, which was deleted or truncated in all but one gastrointestinal

186    strain and intact in all invasive strains. Another was the chaperone protein *yajL*, which

187    appears to be important for oxidative stress tolerance (Kthiri et al. 2010; Le et al. 2012).

188

189    Among the top predictors were several sets of genes belonging to the same operon (S2

190    Table). Examples included the *ttr*, *cbi* and *pdu* operons, which are all required for the

191    anaerobic metabolism of 1,2-propanediol (Roth et al. 1996). These operons have previously

192    been identified as key degraded pathways in invasive isolates (Thomson et al. 2008; Nuccio

9

193    and Bäumler 2014; Langridge et al. 2015), and indicate the agreement of this method with

194    other studies linking loss of gene function to host niche. Overall, a large proportion of the

195    identified genes were involved in metabolism (Fig 2C), consistent with the findings of similar

196    studies (Nuccio and Bäumler 2014; Langridge et al. 2015). Other major categories affected

197    include proteins involved in cell wall and membrane function, perhaps suggesting changes

198    affecting recognition by the host immune system, and signal transduction, suggesting some

199    degree of consistent regulatory rewiring during adaptation to an extraintestinal niche.

200    ***Sequence changes in key indicator genes involve independent mutations in each***

201    ***serovar, contributing to similar functional outcomes***

202    When examining individual genes that showed differences in mutational burden between

203    invasive and gastrointestinal isolates, we found that most of these mutations had occurred

204    independently, and had occurred at different sites in the protein. While the majority of genes

205    identified appeared to be cases of gene degradation in invasive lineages, some genes

206    showed more subtle signs of mutational burden, restricted to nonsynonymous changes of

207    modest predicted functional impact. An example of this, Fig 2D, illustrates mutation

208    accumulation in one of the top candidate genes, *mrcB*, encoding penicillin-binding protein 1b

209    (PBP1b). Not only does *mrcB* carry more mutations in invasive serovars compared to

210    gastrointestinal serovars, the mutations have occurred independently in different positions

211    within the protein. Penicillin-binding proteins are the major target of β-lactam antibiotics and

212    are important for synthesis and maturation of peptidoglycan (Typas et al. 2011). PBP1b in

213    particular extends and crosslinks peptidoglycan chains during cell division. While PBP1b is

214    not essential, it has been shown to be synthetically lethal with PBP1a and is important for

215    competitive survival of extended stationary phase, osmotic stress (Pepper et al. 2006), and

216    — in *Salmonella* Typhi — growth in the presence of bile (Langridge et al. 2009). Bile is an

217    important environmental challenge for *Salmonella*, particularly for extraintestinal serovars

218    which colonize the gall bladder (Crawford et al. 2010). While there are more mutations in

219    invasive than in gastrointestinal serovars, the mutations that occur in this protein are all

220 amino acid substitutions of modest predicted impact. This suggests that sequence changes

221 could result in a modification of protein function, rather than a loss, consistent with the

222 importance of PBP1b for the survival of *S.* Typhi during a typical infection cycle (Langridge et

223 al. 2009).

### *S. Dublin and S. Enteritidis serovars are more difficult to classify than others*

225 To anticipate the performance of our random forest model on new data we computed out-of-

226 bag (OOB) error. Because random forests train each decision tree on a random subset of

227 the training data, OOB error can be computed by testing the performance of these trees on

228 data they have not been trained on, providing inbuilt cross-validation (Breiman 2001). In our

229 case, perfect OOB classifications were only achieved by the fifth iteration of the model. The

230 need for iterative improvement of the model came from difficulty in correctly classifying the

231 reference strains for serovars Enteritidis and Dublin. This is reflective of their relatively

232 recent divergence and niche adaptation compared to other serovars in the study. *S.*

233 Gallinarum was classified much more readily than *S.* Entereitidis and *S.* Dublin, despite

234 being closely related to both serovars, perhaps due to its host restriction.

235

236 *S.* Enteritidis was initially mis-classified as invasive, indicating that it shares genomic trends

237 with invasive lineages. Genomic analyses have indicated that the ancestor of *S.* Enteritidis

238 previously possessed intact pathogenicity islands (SPI-6 and SPI-19), each encoding a type

239 six secretion system (Langridge et al. 2015; Blondel et al. 2009). These loci have been

240 implicated in host-adaptation and survival during extraintestinal infection (Blondel et al. 2013;

241 Mulder et al. 2012), and it has been speculated based on their loss and other evidence that

242 classical *S.* Enteritidis has been adapting towards greater host generalism with respect to its

243 ancestral state (Langridge et al. 2015). This could explain the greater number of disrupted

244 and deleted genes relative to other gastrointestinal serovars used in this study, and the

245 difficulty in classifying it correctly. Conversely, *S.* Dublin was initially mis-classified as

246 gastrointestinal. In previous studies *S.* Dublin has been shown to possess fewer

247    pseudogenes than related invasive isolates (Nuccio and Bäumler 2014; Langridge et al.

248    2015), suggesting a lower degree of host adaptation than other invasive isolates. Indeed, *S.*

249    Dublin is more promiscuous in its host range, primarily infecting cattle (Kingsley and Bäumler

250    2000) while still causing sporadic human disease (Harvey et al. 2017). It seems likely that a

251    subset of informative genes identified in early iterations of the model may have been

252    indicators of host restriction or generalism rather than broad extraintestinal adaptation.


253    ***Patterns of gene degradation identified in established invasive lineages are present in***

254    ***novel lineages of S. Typhimurium and S. Enteritidis associated with systemic***

255    ***infection***

256    In recent years there have been reports of novel *S.* Typhimurium and *S.* Enteritidis lineages

257    associated with invasive disease in sub-Saharan Africa (Kingsley et al. 2009; Okoro et al.

258    2012; Feasey et al. 2016) in populations with a high prevalence of immunosuppressive

259    illness such as HIV, malaria, and malnutrition (Uche et al. 2017). These lineages contribute

260    to a staggering burden of invasive non-typhoidal salmonella (iNTS) disease, which is

261    responsible for an estimated 3.4 million cases and circa 680,000 deaths annually (Ao et al.

262    2015). Based on epidemiological analysis, high-throughput metabolic screening of selected

263    strains, and analysis of pseudogenes it has been suggested that these lineages may be

264    rapidly adapting to cause invasive disease in the human niche created by widespread

265    immunosuppressive illness (Kingsley et al. 2009; Feasey et al. 2012; Okoro et al. 2012,

266    2015; Feasey et al. 2016).

267

268    Two iNTS-associated lineages have recently been described within serovar Enteritidis

269    (Feasey et al. 2016), geographically restricted to West Africa and Central/East Africa,

270    respectively. Initial observations have demonstrated that a representative isolate of the

271    Central/East African clade has a reduced capacity to respire in the presence of metabolites

272    requiring cobalamin for their metabolism and has lost the ability to colonize a chick infection

273    model (Feasey et al. 2016), suggesting adaptation to a new host niche. Similarly, two iNTS

274  disease associated lineages have been described in serovar Typhimurium (Okoro et al.

275  2012), both members of sequence type 313 (ST313), generally referred to as Lineage I and

276  II in the literature. Lineage II appears to have largely replaced Lineage I since 2004, and it

277  has been suggested this is due to Lineage II possessing a gene encoding chloramphenicol

278  resistance (Okoro et al. 2012). Laboratory characterization of Lineage II strains has shown

279  that they are not host-restricted (Parsons et al. 2013; Ramachandran et al. 2017), but do

280  appear to possess characteristics suggestive of adaptation to an invasive lifestyle

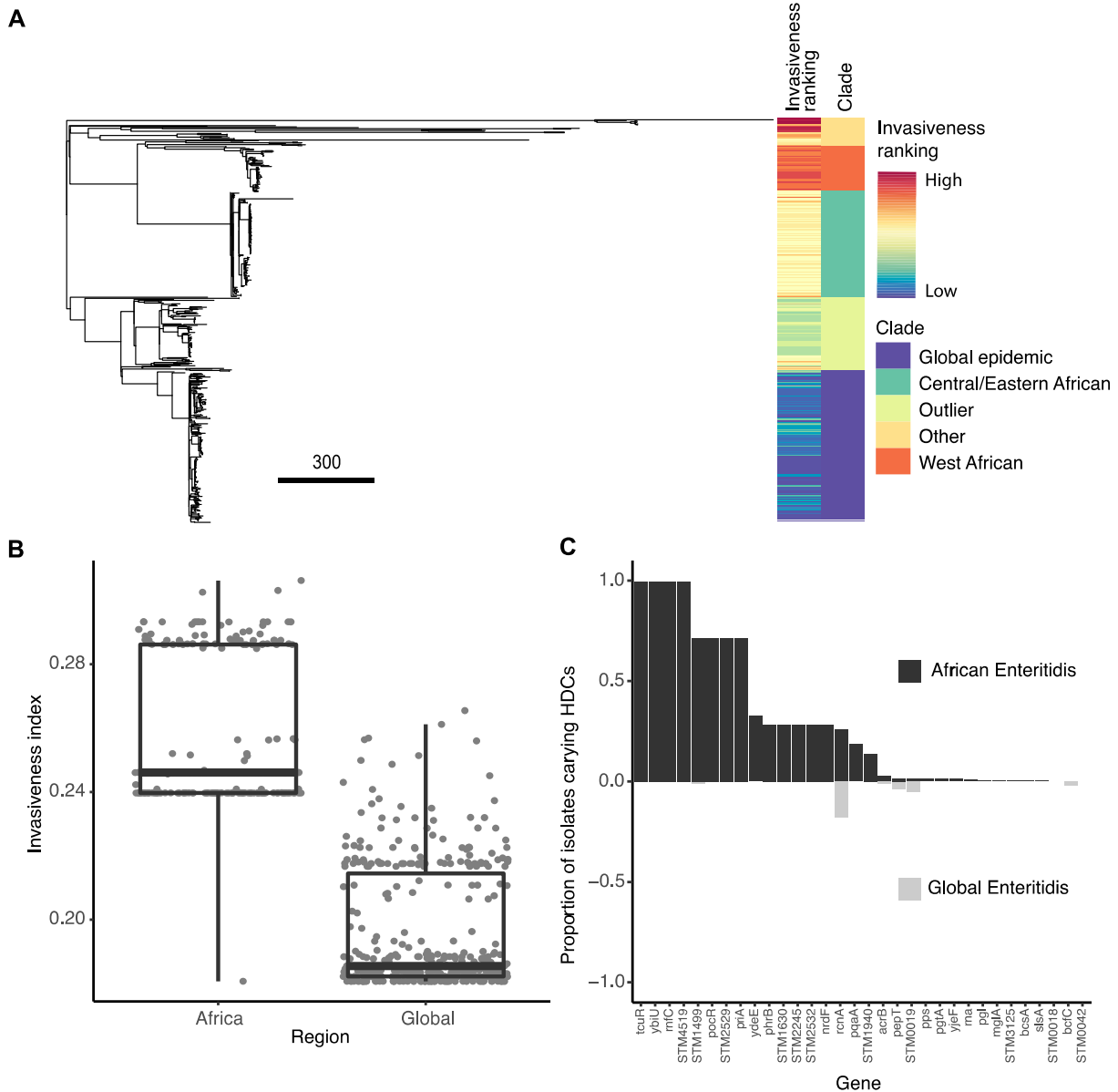281  (Ramachandran et al. 2015; Carden et al. 2015; Singletary et al. 2016; Carden et al. 2017).

282

283  Given the evidence of adaptation to an invasive niche in these lineages, we asked if

284  genomics signatures of extraintestinal adaptation we had detected previously could be

285  detected in iNTS disease associated lineages. To this end, we applied our predictive model

286  trained on well-characterized extraintestinal strains to calculate an invasiveness index, the

287  fraction of decision trees in the random forest voting for an invasive phenotype. First, we

288  compared isolates from African iNTS-associated clades of *S.* Enteritidis (N=233) to a global

289  collection of isolates generally associated with intestinal infection (N=100) (Feasey et al.

290  2016).

291

292  Our model gave iNTS-associated *S.* Enteritidis strains a higher invasiveness index than the

293  globally distributed isolates (Fig 3A,B, Supplemental Table S3), indicating the presence of

294  genetic changes paralleling those that have occurred in extraintestinal serovars of

295  *Salmonella*. Similar gene signatures were only rarely observed in the global epidemic clade

296  (Fig 3C). These findings are consistent with the metabolic changes observed by Feasey et

297  al. (2016) in the Central/Eastern African clade compared to the global epidemic clade. In

298  particular we found signs of gene sequence variation uncharacteristic of gastrointestinal

299  *Salmonella* across a number of key genomic indicators, including *tcuR, ttrA, pocR, pduW,*

300  *eutH,* SEN2509 (a putative anaerobic dimethylsulfoxide reductase) and SEN3188 (a putative

301  tartrate dehydratase subunit), all in pathways previously identified by Nuccio and Bäumler

302    (2014) as being involved in the utilization of host-derived nutrients in the inflamed gut

303    environment. This indicates that our model is able to identify early signatures of adaptation,

304    even in these recently emerged strains that still retain some capacity to cause enterocolitis

305    (Feasey et al. 2016).



306

**Fig 3 | Voting of the model on African iNTS and global gastrointestinal isolates**

308    A: Maximum likelihood phylogeny of all *S*. Enteritidis isolates included in the study,

309    annotated with invasiveness ranking and clade.

310    B: Invasiveness indices for African and non-African clades of *Salmonella*. Lower and upper

311    boundaries of the boxplots correspond to the 25th and 75th quantiles.

312    C: The proportion of isolates from each tested dataset carrying a hypothetically disrupted

313    coding sequence (HDC, as defined by a DeltaBS>3 relative to the reference serovar). Genes

314    are ordered by the amount of degradation observed in African clades. African strains are

315    shown in the positive y-axis in darker grey, global strains are shown in the negative y-axis in

316    lighter grey.

317

318    To confirm this, we performed an additional comparison of *S.* Typhimurium ST313 isolates

319    (N=208), to global isolates from other STs, predominantly ST19, associated with

320    gastroenteritis (N=51) (Okoro et al. 2015; Ashton et al. 2017). Similarly to iNTS associated

321    *S.* Enteritidis isolates, *S.* Typhimurium ST313 isolates has a higher invasiveness index than

322    isolates from other STs (Supplemental Fig S2, Supplemental Table S4). Within ST313,

323    Lineage II scored higher than Lineage I, possibly suggesting differential adaptation to the

324    extraintestinal niche. We found that there were in fact more degraded genes unique to

325    Lineage I than Lineage II, but that these genes were assigned less weight in the model, so

326    did not impact score as strongly (Supplemental Fig S2 & S3). Interestingly, ST313 has

327    recently been shown not to be entirely restricted to Africa, with isolation reported in Brazil

328    (Almeida et al. 2017) and the UK (Ashton et al. 2017). We included a collection of UK ST313

329    strains (Ashton et al. 2017) in our analysis, and found that their invasiveness index tended to

330    be elevated compared to non-ST313 salmonellae, and intermediate between Lineage I and

331    II, suggesting that some of the changes we are detecting are ancestral to ST313 as a whole

332    (Supplemental Fig S3).

333

334    To test whether we could detect a recent case of accelerated adaptation over the course of a

335    single infection, we scored the invasiveness index of a collection of hypermutator *S.*

336    Enteritidis isolates collected over a ten year period that were adapting to chronic systemic

337    infection of an immunocompromised patient (Klemm et al. 2016). We found a significant

338    positive correlation between invasiveness index and duration of carriage (r=0.96, n=6,

339    *P*=0.002, Supplemental Fig S4).

**Discussion**

340

341   Parallel evolution appears to be common in niche adaptation, which allows us to identify

342   genes that are important for survival in different environments. Parallelism has been

343   observed across vastly different time scales in adapting pathogens. Parallel evolution in the

344   distantly related genuses *Salmonella* and *Yersinia* during adaptation to invasive infection of

345   the human host has lead to independent losses of the *ttr*, *cbi* and *pdu* genes, important for

346   anaerobic metabolism during intestinal infection (McNally et al. 2016). Within genuses,

347   parallelism has been observed when distinct lineages acquire similar virulence factors

348   leading to similar phenotypes, as with *Yersinia pseudotuberculosis* and *enterocolitica*

349   (Reuter et al. 2014), or the repeated emergence of the *Shigella* phenotype within the

350   *Escherichia* (The et al. 2016). Even on the scale of a single human lifetime, parallel

351   adaptation has been observed in *Pseudomonas aeruginosa* lineages adapting to infection of

352   the lungs of children with cystic fibrosis (Marvig et al. 2015), or a hypermutator strain of

353   *Salmonella* adapting to an immunocompromised host (Klemm et al. 2016). With pathogen

354   sequencing for disease surveillance becoming increasingly routine (Quick et al. 2016;

355   Aanensen et al. 2016; Schürch and Schaik 2017), we have the opportunity to search for

356   signals of parallel evolution as new pathogens emerge, or old pathogens expand into new

357   niches.

358

359   Here, we have developed an approach for automatically learning which genes contribute to

360   this parallel adaptation. Leveraging the DeltaBS functional variant scoring approach we

361   developed previously (Wheeler et al. 2016) allowed us to construct scores which integrate

362   independent mutations and indels that impact gene function. Using these scores, we were

363   able to construct a classifier model which is able to separate *Salmonella* serovars adapted to

364   an extraintestinal niche from gastrointestinal strains. Importantly, the random forest classifier

365   that we used produces interpretable lists of genes involved in this adaptation, which agree

366   with results in the literature attained through manual curation of pseudogenes. Additionally,

367    we have shown that this classifier is able to identify nascent signatures of adaptation in

368    strains of *Salmonella* which have been evolving in response to large populations of

369    immunocompromised patients in resource-poor nations.

370

371    Other automated approaches to detecting adaptation have been developed which search for

372    SNPs (Lippert et al. 2011) or words (Lees et al. 2016; Earle et al. 2016) associated with

373    phenotype. These approaches, termed microbial genome-wide association studies

374    (GWASs), have used techniques adapted from human GWASs, but better cater to

375    methodological issues that arise due to the differences between human and bacterial

376    inheritance patterns. Major differences impacting analyses are stronger linkage

377    disequilibrium (LD) between genetic variants in bacterial genomes, greater population

378    stratification, and often stronger selection for traits (Chen and Shapiro 2015). Greater LD

379    and population stratification often result in traits being linked closely with particular lineages,

380    and a large number of variants unique to a lineage being spuriously associated with

381    phenotype. Correction for population stratification allows greater discrimination of true and

382    false positive associations, but results in a substantial loss of power to detect true positives

383    (Chen and Shapiro 2015), particularly in phenotypes that are highly polygenic and are not

384    under strong positive selection (Power et al. 2017). This can be corrected by increasing the

385    sample size of the study, but increasing sample size can make measurement of complex

386    phenotypes infeasible (Dutilh et al. 2013).

387

388    DeltaBS differs from current approaches by allowing the estimation of the combined effects

389    of variants, both common and rare, on gene function. The weighting scheme can also

390    combine data on gene presence/absence, indels and SNPs into a single metric. It

391    significantly reduces the number of association tests that need to be performed to

392    comprehensively capture much of the genetic diversity in a species, increasing power to

393    detect associations, and reducing the requirement for such large sample sizes. The

394    approach also aids in identifying genetic variants that are most likely to have a phenotypic

395   effect within LD blocks. The DeltaBS variant scoring approach can be readily applied to large

396   datasets, and could be employed in a linear mixed model (LMM) based association testing

397   framework (Lippert et al. 2011), or used in a hybrid LMM-random forest based approach

398   (Stephan et al. 2015) to preserve the ability of the metric to detect epistasis between genes

399   (Wei et al. 2014).


**Methods**

400


***Genome data and identification of orthologs***

401

402   Genomes for 13 *Salmonella enterica* serovars were retrieved from the NCBI database

403   (accessions and serovar information can be found in S1 Table). The serovars were divided

404   into gastrointestinal and extraintestinal serovars according to the classifications made by

405   Nuccio and Bäumler (2014). Ortholog calls were also taken from the Supplementary Material

406   of Nuccio and Bäumler (2014).


***Measuring the divergence of genes from predicted sequence constraints***

407

408   Profile hidden Markov models (HMMs) for Gammaproteobacterial proteins were retrieved

409   from the eggNOG database (Huerta-Cepas et al. 2016). We chose this source of HMMs

410   because it is publically available, allowing for better reproduction of analyses, and we feel it

411   provides a good balance between collecting enough sequence diversity to capture typical

412   patterns of sequence variation in a protein, without sacrificing sensitivity in the detection of

413   deleterious mutations, as we have observed with Pfam HMMs (Wheeler et al. 2016). Each

414   protein sequence was searched against the HMM database using hmmsearch from the

415   HMMER3.0 package (http://hmmer.org). The top scoring model corresponding to each

416   protein was used for analysis (N = 8,060 groups). Orthologous groups (OGs) with no

417   corresponding eggNOG HMM, or more than one top model hit were excluded from further

418   analysis (N = 1,524). If most genes in an OG had a significant hit (E-value<0.0001) to the

419   same eggNOG model, any genes within this OG that did not were assigned a score of zero,

18

420     reflecting a loss of the function of that protein. These cases typically reflected a truncation

421     that had occurred early in the protein sequence. Additionally, genes with no variation in

422     bitscore for the match between protein sequences and their respective eggNOG HMM

423     across isolates were excluded (N = 188). After this filtering process, 6,439 orthologous

424     groups remained for analysis. Residue-specific DeltaBS (as in Fig 2D) was calculated by

425     aligning orthologous sequences, choosing a reference sequence (from *S.* Typhimurium), and

426     substituting each variant match state and any accompanying insertions into the reference

427     sequence and calculating the difference in bitscore caused by the substitution.


428     ***Training a random forest classifier***

429     The R package "randomForest" (Liaw and Wiener 2002) was used to build random forest

430     classifiers using a variety of parameters to assess which were best for accuracy. Prediction

431     accuracy, as measured by out-of-bag (OOB) error rate, stabilised at 1000 trees, so we chose

432     this as a parameter for optimising the number of genes sampled per node (mtry). mtry

433     values of 1, $p/10$, $p/5$, $p/3$, $p/2$ and $p$ (where $p$ = the number of predictors) were tested, and

434     we found that at mtry=$p/10$, the number of genes that were either not incorporated into trees,

435     or did not improve the homogeneity of daughter nodes when they were incorporated into

436     trees (as measured by mean decrease in Gini index, (Breiman et al. 1984)) stabilised at

437     ~92%.

438

439     To improve the performance of the model, we performed five model building and sparsity

440     pruning cycles. For the first cycle, we built a random forest model using all genes that met

441     the inclusion criteria, and performed sparsity pruning by eliminating all variables that had a

442     mean Gini index (variable importance) of zero or lower (meaning the gene was either not

443     included in the model or did not improve model accuracy when it was). Four successive

444     rounds of model building and sparsity pruning involved building a new model with the pruned

445     dataset, then pruning the genes with the lowest 50% of variable importances. The resulting

446     model had 100% out-of-bag classification accuracy. We also tested the accuracy of the full

447 model on a collection of alternative strains related to the training dataset (see Table S1).

448 Orthologs to the top genes identified by our model were identified using phmmer from the

449 HMMER3.0 package (http://hmmer.org).


450 ***Invasive non-typhoidal Salmonella analysis***

451 Read data from Feasey et al. (2016) and Klemm et al. (2016) was mapped to the reference

452 genome *S.* Enteritidis P125109. Reads from Okoro et al. (2015) and Ashton et al. (2017)

453 were mapped to the reference genome *S.* Typhimurium LT2. For samples in the Okoro

454 study, if an isolate was sequenced using multiple runs, the most recent run was chosen for

455 analysis. All reads were mapped using BWA mem (Li and Durbin 2009) and regions near

456 indels were realigned using GATK (McKenna et al. 2010). Picard

457 (http://broadinstitute.github.io/picard) was used to identify and flag optical duplicates

458 generated during library preparation. SNPs and indels were called using samtools v1.2

459 mpileup (Li 2011), and were filtered to exclude those variants with coverage <10 or quality

460 <30. For tree building, a pseudogenome was constructed by substituting high confidence

461 (coverage >4, quality >50) variant sites in the reference genome, and masking any sites with

462 low confidence with an "N". Insertions relative to the reference genome were ignored, and

463 deletions were filled with an "N". Pseudogenome alignments were then used as input to

464 produce trees using Gubbins (Croucher et al. 2015) to exclude recombination events, and

465 RAxML v8.2.8 (Stamatakis 2014) to build maximum likelihood trees using a GTR + Gamma

466 model.

467

468 Sequences for the 196 genes of interest used in the random forest model were retrieved for

469 each isolate and translated. These were then scored using their respective profile HMMs.

470 Score data was collated, and any missing values were marked as 'NA' and imputed using

471 the na.roughfix function from the randomForest R package (Liaw and Wiener 2002). This is

472 a different approach used to that of the training dataset, due to the potentially lower quality of

473 the sequenced genomes leading to gene absence due to low coverage rather than true

474 deletion or severe truncation. The relationship between invasiveness ranking and phylogeny

475 were visualised using Phandango (Hadfield et al. 2017).

476 **Data access**

477 All genome sequence data are publically available, and accessions are provided in the

478 appropriate Supplemental Tables. Code and data for reproducing this analysis, performing

479 an equivalent analysis using new data, and assessing the invasiveness index of other

480 *Salmonella* strains is publically available at github.com/UCanCompBio/invasive_salmonella.

493 **References**

494 Aanensen DM, Feil EJ, Holden MTG, Dordel J, Yeats CA, Fedosejev A, Goater R, Castillo-Ramírez S,

495     Corander J, Colijn C, et al. 2016. Whole-Genome Sequencing for Routine Pathogen

496     Surveillance in Public Health: a Population Snapshot of Invasive Staphylococcus aureus in

497     Europe. *MBio* **7**. doi: 10.1128/mBio.00444-16.

498    Alam MT, Petit RA 3rd, Crispell EK, Thornton TA, Conneely KN, Jiang Y, Satola SW, Read TD. 2014.

499        Dissecting vancomycin-intermediate resistance in staphylococcus aureus using genome-wide

500        association. *Genome Biol Evol* **6**: 1174–1185.

501    Almeida F, Seribelli AA, da Silva P, Medeiros MIC, Dos Prazeres Rodrigues D, Moreira CG, Allard

502        MW, Falcão JP. 2017. Multilocus sequence typing of Salmonella Typhimurium reveals the

503        presence of the highly invasive ST313 in Brazil. *Infect Genet Evol* **51**: 41–44.

504    Ao TT, Feasey NA, Gordon MA, Heddy KH, Angulo FJ, Crump JA. 2015. Global Burden of Invasive

505        Nontyphoidal Salmonella Disease, 2010[1]. *Emerging Infectious Disease journal* **21**: 941.

506    Ashton PM, Owen SV, Kaindama L, Rowe WPM, Lane C, Larkin L, Nair S, Jenkins C, de Pinna E,

507        Feasey N, et al. 2017. Salmonella enterica Serovar Typhimurium ST313 Responsible For

508        Gastroenteritis In The UK Are Genetically Distinct From Isolates Causing Bloodstream Infections

509        In Africa. *bioRxiv* 139576. doi: 10.1101/139576.

510    Bäumler A, Fang FC. 2013. Host specificity of bacterial pathogens. *Cold Spring Harb Perspect Med* **3**:

511        a010041.

512    Bayjanov JR, Molenaar D, Tzeneva V, Siezen RJ, van Hijum SAFT. 2012. PhenoLink--a web-tool for

513        linking phenotype to ~omics data for bacteria: application to gene-trait matching for Lactobacillus

514        plantarum strains. *BMC Genomics* **13**: 170.

515    Blondel CJ, Jiménez JC, Contreras I, Santiviago CA. 2009. Comparative genomic analysis uncovers

516        3 novel loci encoding type six secretion systems differentially distributed in Salmonella

517        serotypes. *BMC Genomics* **10**: 354.

518    Blondel CJ, Jiménez JC, Leiva LE, Alvarez SA, Pinto BI, Contreras F, Pezoa D, Santiviago CA,

519        Contreras I. 2013. The type VI secretion system encoded in Salmonella pathogenicity island 19

520        is required for Salmonella enterica serotype Gallinarum survival within infected macrophages.

521        *Infect Immun* **81**: 1207–1220.

522    Breiman L. 2001. Random Forests. *Mach Learn* **45**: 5–32.

523    Breiman L, Friedman J, Stone CJ, Olshen RA. 1984. *Classification and Regression Trees*. Chapman

524        and Hall/CRC.

525    Carden SE, Walker GT, Honeycutt J, Lugo K, Pham T, Jacobson A, Bouley D, Idoyaga J, Tsolis RM,

526        Monack D. 2017. Pseudogenization of the Secreted Effector Gene sseI Confers Rapid Systemic

527        Dissemination of S. Typhimurium ST313 within Migratory Dendritic Cells. *Cell Host Microbe* **21**:

528       182–194.

529       Carden S, Okoro C, Dougan G, Monack D. 2015. Non-typhoidal Salmonella Typhimurium ST313

530              isolates that cause bacteremia in humans stimulate less inflammasome activation than ST19

531              isolates associated with gastroenteritis. *Pathog Dis* **73**. doi: 10.1093/femspd/ftu023.

532       Chen PE, Shapiro BJ. 2015. The advent of genome-wide association studies for bacteria. *Curr Opin*

533              *Microbiol* **25**: 17–24.

534       Crawford RW, Rosales-Reyes R, Ramírez-Aguilar M de la L, Chapa-Azuela O, Alpuche-Aranda C,

535              Gunn JS. 2010. Gallstones play a significant role in Salmonella spp. gallbladder colonization and

536              carriage. *Proc Natl Acad Sci U S A* **107**: 4353–4358.

537       Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, Parkhill J, Harris SR. 2015.

538              Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome

539              sequences using Gubbins. *Nucleic Acids Res* **43**: e15.

540       Dutilh BE, Backus L, Edwards RA, Wels M, Bayjanov JR, van Hijum SAFT. 2013. Explaining microbial

541              phenotypes on a genomic scale: GWAS for microbes. *Brief Funct Genomics* **12**: 366–380.

542       Earle SG, Wu C-H, Charlesworth J, Stoesser N, Gordon NC, Walker TM, Spencer CCA, Iqbal Z,

543              Clifton DA, Hopkins KL, et al. 2016. Identifying lineage effects when controlling for population

544              structure improves power in bacterial association studies. *Nat Microbiol* **1**: 16041.

545       Fauci AS, Morens DM. 2012. The perpetual challenge of infectious diseases. *N Engl J Med* **366**: 454–

546              461.

547       Feasey NA, Dougan G, Kingsley RA, Heyderman RS, Gordon MA. 2012. Invasive non-typhoidal

548              salmonella disease: an emerging and neglected tropical disease in Africa. *Lancet* **379**: 2489–

549              2499.

550       Feasey NA, Hadfield J, Keddy KH, Dallman TJ, Jacobs J, Deng X, Wigley P, Barquist Barquist L,

551              Langridge GC, Feltwell T, et al. 2016. Distinct Salmonella Enteritidis lineages associated with

552              enterocolitis in high-income settings and invasive disease in low-income settings. *Nat Genet* **48**:

553              1211–1217.

554       Frank SA, Schmid-Hempel P. 2008. Mechanisms of pathogenesis and the evolution of parasite

555              virulence. *J Evol Biol* **21**: 396–404.

556       Hadfield J, Croucher NJ, Goater RJ, Abudahab K, Aanensen DM, Harris SR. 2017. Phandango: an

557              interactive viewer for bacterial population genomics. *Bioinformatics*. doi:

558     10.1093/bioinformatics/btx610.

559     Harvey RR, Friedman CR, Crim SM, Judd M, Barrett KA, Tolar B, Folster JP, Griffin PM, Brown AC.

560     2017. Epidemiology of Salmonella enterica Serotype Dublin Infections among Humans, United

561     States, 1968–2013. *Emerging Infectious Disease journal* **23**: 1493.

562     Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC, Rattei T, Mende DR,

563     Sunagawa S, Kuhn M, et al. 2016. eggNOG 4.5: a hierarchical orthology framework with

564     improved functional annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic Acids*

565     *Res* **44**: D286–93.

566     Kingsley RA, Bäumler AJ. 2000. Host adaptation and the emergence of infectious disease: the

567     Salmonella paradigm. *Mol Microbiol* **36**: 1006–1014.

568     Kingsley RA, Kay S, Connor T, Barquist L, Sait L, Holt KE, Sivaraman K, Wileman T, Goulding D,

569     Clare S, et al. 2013. Genome and transcriptome adaptation accompanying emergence of the

570     definitive type 2 host-restricted Salmonella enterica serovar Typhimurium pathovar. *MBio* **4**:

571     e00565–13.

572     Kingsley RA, Msefula CL, Thomson NR, Kariuki S, Holt KE, Gordon MA, Harris D, Clarke L,

573     Whitehead S, Sangal V, et al. 2009. Epidemic multiple drug resistant Salmonella Typhimurium

574     causing invasive disease in sub-Saharan Africa have a distinct genotype. *Genome Res* **19**:

575     2279–2287.

576     Klemm EJ, Gkrania-Klotsas E, Hadfield J, Forbester JL, Harris SR, Hale C, Heath JN, Wileman T,

577     Clare S, Kane L, et al. 2016. Emergence of host-adapted Salmonella Enteritidis through rapid

578     evolution in an immunocompromised host. *Nat Microbiol* **1**: 15023.

579     Kthiri F, Gautier V, Le H-T, Prère M-F, Fayet O, Malki A, Landoulsi A, Richarme G. 2010.

580     Translational defects in a mutant deficient in YajL, the bacterial homolog of the parkinsonism-

581     associated protein DJ-1. *J Bacteriol* **192**: 6302–6306.

582     Kuo C-H, Ochman H. 2010. The extinction dynamics of bacterial pseudogenes. *PLoS Genet* **6**.

583     http://dx.doi.org/10.1371/journal.pgen.1001050.

584     Laabei M, Recker M, Rudkin JK, Aldeljawi M, Gulay Z, Sloan TJ, Williams P, Endres JL, Bayles KW,

585     Fey PD, et al. 2014. Predicting the virulence of MRSA from its genome sequence. *Genome Res*

586     **24**: 839–849.

587     Langridge GC, Fookes M, Connor TR, Feltwell T, Feasey N, Parsons BN, Seth-Smith HMB, Barquist

588     L, Stedman A, Humphrey T, et al. 2015. Patterns of genome evolution that have accompanied

589     host adaptation in Salmonella. *Proc Natl Acad Sci U S A* **112**: 863–868.

590  Langridge GC, Phan M-D, Turner DJ, Perkins TT, Parts L, Haase J, Charles I, Maskell DJ, Peters SE,

591     Dougan G, et al. 2009. Simultaneous assay of every Salmonella Typhi gene using one million

592     transposon mutants. *Genome Res* **19**: 2308–2316.

593  Lees JA, Vehkala M, Välimäki N, Harris SR, Chewapreecha C, Croucher NJ, Marttinen P, Davies MR,

594     Steer AC, Tong SYC, et al. 2016. Sequence element enrichment analysis to determine the

595     genetic basis of bacterial phenotypes. *Nat Commun* **7**: 12797.

596  Le H-T, Gautier V, Kthiri F, Malki A, Messaoudi N, Mihoub M, Landoulsi A, An YJ, Cha S-S, Richarme

597     G. 2012. YajL, prokaryotic homolog of parkinsonism-associated protein DJ-1, functions as a

598     covalent chaperone for thiol proteome. *J Biol Chem* **287**: 5861–5870.

599  Lerat E, Ochman H. 2005. Recognizing the pseudogenes in bacterial genomes. *Nucleic Acids Res*

600     **33**: 3125–3132.

601  Liaw A, Wiener M. 2002. Classification and regression by randomForest. *R news* **2**: 18–22.

602  Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping and

603     population genetical parameter estimation from sequencing data. *Bioinformatics* **27**: 2987–2993.

604  Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.

605     *Bioinformatics* **25**: 1754–1760.

606  Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D. 2011. FaST linear mixed

607     models for genome-wide association studies. *Nat Methods* **8**: 833–835.

608  Loman NJ, Pallen MJ. 2015. Twenty years of bacterial genome sequencing. *Nat Rev Microbiol* **13**:

609     787–794.

610  Marvig RL, Sommer LM, Molin S, Johansen HK. 2015. Convergent evolution and adaptation of

611     Pseudomonas aeruginosa within patients with cystic fibrosis. *Nat Genet* **47**: 57–64.

612  McClelland M, Sanderson KE, Clifton SW, Latreille P, Porwollik S, Sabo A, Meyer R, Bieri T, Ozersky

613     P, McLellan M, et al. 2004. Comparison of genome degradation in Paratyphi A and Typhi,

614     human-restricted serovars of Salmonella enterica that cause typhoid. *Nat Genet* **36**: 1268–1274.

615  McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D,

616     Gabriel S, Daly M, et al. 2010. The Genome Analysis Toolkit: a MapReduce framework for

617     analyzing next-generation DNA sequencing data. *Genome Res* **20**: 1297–1303.

618  McNally A, Thomson NR, Reuter S, Wren BW. 2016. "Add, stir and reduce": Yersinia spp. as model

619      bacteria for pathogen evolution. *Nat Rev Microbiol* **14**: 177–190.

620  Merhej V, Georgiades K, Raoult D. 2013. Postgenomic analysis of bacterial pathogens repertoire

621      reveals genome reduction rather than virulence factors. *Brief Funct Genomics* **12**: 291–304.

622  Mulder DT, Cooper CA, Coombes BK. 2012. Type VI secretion system-associated gene clusters

623      contribute to pathogenesis of Salmonella enterica serovar Typhimurium. *Infect Immun* **80**: 1996–

624      2007.

625  Nuccio S-P, Bäumler AJ. 2014. Comparative Analysis of Salmonella Genomes Identifies a Metabolic

626      Network for Escalating Growth in the Inflamed Gut. *MBio* **5**: e00929–14–e00929–14.

627  Okoro CK, Barquist L, Connor TR, Harris SR, Clare S, Stevens MP, Arends MJ, Hale C, Kane L,

628      Pickard DJ, et al. 2015. Signatures of Adaptation in Human Invasive Salmonella Typhimurium

629      ST313 Populations from Sub-Saharan Africa. *PLoS Negl Trop Dis* **9**: e0003611.

630  Okoro CK, Kingsley RA, Connor TR, Harris SR, Parry CM, Al-Mashhadani MN, Kariuki S, Msefula CL,

631      Gordon MA, de Pinna E, et al. 2012. Intracontinental spread of human invasive Salmonella

632      Typhimurium pathovariants in sub-Saharan Africa. *Nat Genet* **44**: 1215–1221.

633  Pallen MJ, Wren BW. 2007. Bacterial pathogenomics. *Nature* **449**: 835–842.

634  Pappu V, Pardalos PM. 2014. High-Dimensional Data Classification. In *Clusters, Orders, and Trees:*

635      *Methods and Applications* (eds. F. Aleskerov, B. Goldengorin, and P.M. Pardalos), *Springer*

636      *Optimization and Its Applications*, pp. 119–150, Springer New York.

637  Parkhill J, Dougan G, James KD, Thomson NR, Pickard D, Wain J, Churcher C, Mungall KL, Bentley

638      SD, Holden MT, et al. 2001. Complete genome sequence of a multiple drug resistant Salmonella

639      enterica serovar Typhi CT18. *Nature* **413**: 848–852.

640  Parsons BN, Humphrey S, Salisbury AM, Mikoleit J, Hinton JCD, Gordon MA, Wigley P. 2013.

641      Invasive non-typhoidal Salmonella typhimurium ST313 are not host-restricted and have an

642      invasive phenotype in experimentally infected chickens. *PLoS Negl Trop Dis* **7**: e2487.

643  Pepper ED, Farrell MJ, Finkel SE. 2006. Role of penicillin-binding protein 1b in competitive stationary-

644      phase survival of Escherichia coli. *FEMS Microbiol Lett* **263**: 61–67.

645  Power RA, Parkhill J, de Oliveira T. 2017. Microbial genome-wide association studies: lessons from

646      human GWAS. *Nat Rev Genet* **18**: 41–50.

647  Quick J, Loman NJ, Duraffour S, Simpson JT, Severi E, Cowley L, Bore JA, Koundouno R, Dudas G,

648    Mikhail A, et al. 2016. Real-time, portable genome sequencing for Ebola surveillance. *Nature*
649    **530**: 228–232.

650    Rabsch W, Andrews HL, Kingsley RA, Prager R, Tschäpe H, Adams LG, Bäumler AJ. 2002.
651    Salmonella enterica serotype Typhimurium and its host-adapted variants. *Infect Immun* **70**:
652    2249–2255.

653    Ramachandran G, Panda A, Higginson EE, Ateh E, Lipsky MM, Sen S, Matson CA, Permala-Booth J,
654    DeTolla LJ, Tennant SM. 2017. Virulence of invasive Salmonella Typhimurium ST313 in animal
655    models of infection. *PLoS Negl Trop Dis* **11**: e0005697.

656    Ramachandran G, Perkins DJ, Schmidlein PJ, Tulapurkar ME, Tennant SM. 2015. Invasive
657    Salmonella Typhimurium ST313 with naturally attenuated flagellin elicits reduced inflammation
658    and replicates within macrophages. *PLoS Negl Trop Dis* **9**: e3394.

659    Reuter S, Connor TR, Barquist L, Walker D, Feltwell T, Harris SR, Fookes M, Hall ME, Petty NK,
660    Fuchs TM, et al. 2014. Parallel independent evolution of pathogenicity within the genus Yersinia.
661    *Proc Natl Acad Sci U S A* **111**: 6768–6773.

662    Roth JR, Lawrence JG, Bobik TA. 1996. Cobalamin (coenzyme B12): synthesis and biological
663    significance. *Annu Rev Microbiol* **50**: 137–181.

664    Schürch AC, Schaik W. 2017. Challenges and opportunities for whole-genome sequencing--based
665    surveillance of antibiotic resistance. *Ann N Y Acad Sci* **1388**: 108–120.

666    Singletary LA, Karlinsey JE, Libby SJ, Mooney JP, Lokken KL, Tsolis RM, Byndloss MX, Hirao LA,
667    Gaulke CA, Crawford RW, et al. 2016. Loss of Multicellular Behavior in Epidemic African
668    Nontyphoidal Salmonella enterica Serovar Typhimurium ST313 Strain D23580. *MBio* **7**. doi:
669    10.1128/mBio.02265-15.

670    Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
671    phylogenies. *Bioinformatics* **30**: 1312–1313.

672    Stephan J, Stegle O, Beyer A. 2015. A random forest approach to capture genetic effects in the
673    presence of population structure. *Nat Commun* **6**: 7432.

674    The HC, Thanh DP, Holt KE, Thomson NR, Baker S. 2016. The genomic signatures of Shigella
675    evolution, adaptation and geographical spread. *Nat Rev Microbiol*.
676    http://dx.doi.org/10.1038/nrmicro.2016.10.

677    Thomson NR, Clayton DJ, Windhorst D, Vernikos G, Davidson S, Churcher C, Quail MA, Stevens M,

678    Jones MA, Watson M, et al. 2008. Comparative genome analysis of Salmonella Enteritidis PT4

679    and Salmonella Gallinarum 287/91 provides insights into evolutionary and host adaptation

680    pathways. *Genome Res* **18**: 1624–1637.

681    Touw WG, Bayjanov JR, Overmars L, Backus L, Boekhorst J, Wels M, van Hijum SAFT. 2013. Data

682    mining in the Life Sciences with Random Forest: a walk in the park or lost in the jungle? *Brief*

683    *Bioinform* **14**: 315–326.

684    Typas A, Banzhaf M, Gross CA, Vollmer W. 2011. From the regulation of peptidoglycan synthesis to

685    bacterial growth and morphology. *Nat Rev Microbiol* **10**: 123–136.

686    Uche IV, MacLennan CA, Saul A. 2017. A Systematic Review of the Incidence, Risk Factors and

687    Case Fatality Rates of Invasive Nontyphoidal Salmonella (iNTS) Disease in Africa (1966 to

688    2014). *PLoS Negl Trop Dis* **11**: e0005118.

689    Wei W-H, Hemani G, Haley CS. 2014. Detecting epistasis in human complex traits. *Nat Rev Genet*

690    **15**: 722–733.

691    Wheeler NE, Barquist L, Kingsley RA, Gardner PP. 2016. A profile-based method for identifying

692    functional divergence of orthologous genes in bacterial genomes. *Bioinformatics* **32**: 3566–3574.

693