# The phylogenomic landscape of the genus *Serratia*

David J. Williams[1,2], Patrick A. D. Grimont[3], Adrián Cazares[2,4], Francine Grimont[3], Elisabeth Ageron[3], Kerry A. Pettigrew[5], Daniel Cazares[2], Elisabeth Njamkepo[6], François-Xavier Weill[6], Eva Heinz[2,7], Matthew T. G. Holden[5], Nicholas R. Thomson[2,8]* and Sarah J. Coulthurst[1]*.

[1]Division of Molecular Microbiology, School of Life Sciences, University of Dundee, Dundee, UK. [2]Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK. [3]Unité Biodiversité des Bactéries Pathogènes Emergentes, INSERM Unité 389, Institut Pasteur, Paris, France. [4]European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, UK. [5]School of Medicine, University of St Andrews, St Andrews, UK. [6]Institut Pasteur, Université de Paris, Unité des Bactéries Pathogènes Entériques, Paris, France. [7]Departments of Vector Biology and Clinical Sciences, Liverpool School of Tropical Medicine, Liverpool, UK. [8]Department of Infection Biology, London School of Hygiene and Tropical Medicine, London, UK.

*Correspondence may be addressed to Nicholas Thomson (nrt@sanger.ac.uk) or Sarah Coulthurst (s.j.coulthurst@dundee.ac.uk)

1

1   **Abstract**

2   The genus *Serratia* has been studied for over a century and includes clinically-important and diverse

3   environmental members. Despite this, there is a paucity of genomic information across the genus and a

4   robust whole genome-based phylogenetic framework is lacking. Here, we have assembled and analysed

5   a representative set of 664 genomes from across the genus, including 215 historic isolates originally

6   used in defining the genus. Phylogenomic analysis of the genus reveals a clearly-defined population

7   structure which displays deep divisions and aligns with ecological niche, as well as striking congruence

8   between historical biochemical phenotyping data and contemporary genomics data. We show that

9   *Serratia* is a diverse genus which displays striking plasticity and ability to adapt to its environment,

10   including a highly-varied portfolio of plasmids, and provide evidence of different patterns of gene flow

11   across the genus. This work provides an essential platform for understanding the emergence of clinical

12   and other lineages of *Serratia*.

13    The genus *Serratia* was originally described in Italy in the early 19[th] century, following the observation

14    of a blood-like red discoloration appearing on polenta from organic growth[1]. It has since become clear

15    that *Serratia* species are ubiquitous, free-living, motile Gram-negative proteobacteria, traditionally

16    considered members of the *Enterobacteriaceae*. The genus *Serratia* represents a broad and diverse

17    genus of more than ten species, delineated by DNA-DNA hybridization and characterized by extensive

18    physiological and biochemical tests[2–12]. Despite being a diverse genus, much of the contemporary

19    research and understanding of *Serratia* has focused on the type species, *Serratia marcescens*. *S.*

20    *marcescens* has served as a model system for studying key bacterial traits, including protein secretion

21    systems[13] and motility[14], but it also represents an important opportunistic human pathogen[15–17] for which

22    there has been a dramatic rise in the incidence of multi-drug resistance and reported cases of problematic

23    nosocomial infections[18].

24    Other members of this genus include *S. rubidaea* and *S. liquefaciens,* which have also been reported to

25    cause hospital-acquired infections, albeit less frequently[17,19–21]. In addition to infection of human hosts,

26    members of multiple *Serratia* species represent insect pathogens or are otherwise associated with

27    insects. *Serratia entomophila* has been used as a biocontrol agent in New Zealand to predate upon the

28    pasture pest, *Costelytra zealandica*[12,22–25], and *S. proteamaculans*[26] and *S. marcescens*[13,27–29] have also

29    been shown to be insect pathogens. In contrast, *S. ficaria* is associated with the pollination and

30    oviposition cycle between figs and fig wasps, respectively[6]. In addition, underlining the ubiquitous

31    nature of this genus, *Serratia* species can be found in a multitude of environmental niches[7,9,29–36],

32    including frequent isolation from aqueous environments[17,21].

33    Given its importance to human health, it is perhaps unsurprising that the majority of genomic

34    information available for *Serratia* originate from clinically-derived *S. marcescens.* Recently a number

35    of *S. marcescens* sequences have also been included within large scale metagenomic studies from pre-

36    term neonates[90] or nosocomial environments[91]. However, for these, as for all the sequences from

37    clinically isolated strains, there is a critical lack of a robust phylogenetic framework for the *Serratia*

38    genus within which the *S. marcescens* sequences can be placed. In an attempt to understand the many

39    facets and functions of the different species within the *Serratia* genus, and, importantly, to understand

40    the context within which *S. marcescens* is becoming more widespread as a problematic opportunistic

41    pathogen, we aimed to assemble a balanced genomic dataset that reflected the entire *Serratia* genus.

42    We supplemented existing publicly-available, published *Serratia* genome sequences by sequencing

43    *Serratia* isolates that were non-clinical in origin, and mainly belonging to *Serratia* species other than

44    *S. marcescens*. We included the historic collection of Patrick Grimont, located at the Institut Pasteur, a

45    collection that includes the original strains used to define biochemically and phenotypically the vast

46    majority of the known *Serratia* species[3,4,6,9–12] These strains, kept in cold storage for between 20 to 40

47    years, represented a unique resource, providing the opportunity to compare historical biochemical,

48    phenotypic and DNA-DNA hybridisation data with contemporary genomics data. In this study, we bring

49 together all of the previous molecular and biochemical knowledge from this important genus and place

50 it in a genomics framework. This not only explains why previous definitions of this genus were robust,

51 or not, it also highlights important differences in diversity, plasticity and niche adaptation of the species

52 within it.

53

54 **Results**

55 **Deep divisions demarcate phylogroups within *Serratia***

56 Here we sequenced and analysed a collection of 256 novel *Serratia* genome sequences and combined

57 these data with 408 published genomes. Our total set of 664 genome sequences included those of 215

58 isolates from the original genus-defining Grimont collection sequenced here, 205 multidrug-resistant *S.*

59 *marcescens* isolates from UK hospitals isolated between 2001 and 2011[18], and an additional 41, more

60 diverse, *Serratia* isolates from UK hospitals, sequenced here for comparison with the latter collection[18].

61 We inferred the genus phylogeny from the whole genome data using a core-gene alignment-based

62 approach. It is evident from Figure 1 that there are deep divisions within this phylogeny that correlate

63 with both the current genus taxonomy and with species-level grouping calculated using genome-wide

64 average nucleotide identity (ANI; clustered using a cutoff of 95 percent; Fig. 1, Supplementary Fig. 2).

65 To gain a finer-scaled view, we used hierarchical Bayesian clustering (FastBaps) to four levels in order

66 to further subdivide the phylogeny and the species-level groups. In total, we identified 7, 16, 23 and 27

67 clusters across the four levels, respectively (Supplementary Fig. 3). FastBaps level 1 clusters comprise

68 monophyletic clades reflecting individual or multiple ANI groupings within the genus, consistent with

69 speciation or species complexes[37]. The second and third levels reveal the presence of several

70 subdivisions within some of the species-level phylogroups (Fig. 1, Supplementary Fig. 3). Hereafter,

71 we refer to the clusters set out by FastBaps level 3 as Lineages 1-23 (L1-23; Fig. 1).

72 Interestingly, within the tree there are two examples of singleton genomes occupying both a single ANI

73 species-level phylogroup and representing a discrete FastBaps lineage (L16 and L23). Although both

74 are situated within well-defined species, these two singletons are hereafter referred to as "*S.*

75 *marcescens*-like" and "*S. rubidaea*-like", for L16 and L23, respectively. Our phylogeny also resolves

76 previous taxonomic discrepancies. Here, based on the core-gene phylogeny, the *liquefaciens* complex,

77 made up of *S. liquefaciens*, *S. grimesii*, and *S. proteamaculans,* is monophyletic (Fig. 1). Previous work

78 had suggested *S. proteamaculans* be resolved into both *S. proteamaculans sensu stricto* and a

79 subspecies, termed *S. proteamaculans* subs. *quinovora*[10,21]. However, a species level distinction

80 between these two taxa, rather than a sub-species one, was subsequently proposed[30]. We observe that

81 genomes labelled as *S. proteamaculans* and *S. proteamaculans* subs. *quinovora* form two separate ANI

82 phylogroups in a monophyletic branch made up of L6-8 (Fig. 1, Supplementary Fig. 2). This supports

4

83    the presence of two distinct species-level groups, which we refer to as *S. proteamaculans* and *S.*

84    *quinivorans* in accordance with the latter work[30]. Furthermore, only a single genome links L7 and L8

85    into a single ANI phylogroup within *S. quinivorans* (Supplementary Fig. 4), which may suggest that a

86    further sub-species separation within *S. quinivorans* is appropriate.

87    **Concordance between historical biochemical phenotyping and metabolomic predictions**

88    The current taxonomic structure of the *Serratia* genus, summarised by Grimont and Grimont, 2005 [92],

89    is based on 41 phenotypic and biochemical tests used to differentiate between different *Serratia* species

90    or species-complexes. Many of the *Serratia* isolates originally used to define the genus taxonomy were

91    sequenced here (Fig. 1; Supplementary Table 1), presenting a unique opportunity to reconcile this

92    historical biochemical metadata with genomic predictions. First, we calculated the genus pan-genome

93    using a population structure-aware approach[38]. The pan-genome comprised 47,743 discrete gene groups

94    (Fig. 2a), of which 2252 were present in at least 99 percent of all genomes in the dataset (which would

95    be defined as a "traditional" core genome), however only 1655 of these were present in at least 95

96    percent of genomes within each FastBaps lineage, and therefore are core to all lineages (Fig. 2b). These

97    1655 genes are hereafter defined as the "genus-core". From the 47,743 genes of the pan-genome, we

98    predicted the metabolic potential of the genus. We identified 641 different complete metabolic

99    pathways using Pathway tools[39] (Fig. 3a), of which 260 were core to the genus, being present in all

100    known lineages (L1-23) (Fig. 3a, Supplementary Fig. 9).

101    Of the 41 metabolic tests used to define species within the genus, some were also used to define biotypes

102    within a species. It can be seen that the fine-scale delimitation of biotypes within the *liquefaciens*

103    complex corresponds with the phylogenetic structure we observe (Supplementary Fig. 10). Similarly,

104    ten of the 41 metabolic tests were used previously to split *marcescens* into ten biotypes, which reflected

105    differences in niche occupancy[21]. It is clear that these biotypes are also robust markers of phylogenetic

106    subdivisions within this important species (Fig. 3c, 4c).

107    *In silico* pathway predictions were used to identify the genes and/or pathways linked to the biotype tests

108    for *S. marcescens*, where strain biotype metadata was available. Four of these ten tests (growth on *m*-

109    erythritol, trigonelline, 3-hydroxybenzoate and lactose) were not investigated because there were no

110    corresponding pathway assignments in the predicted metabolic network. Next, the phylogenetic

111    distribution of these metabolic genes was plotted across the phylogeny and extrapolated back to the

112    most basal internal node differentiating two biotypes, based on gene/pathway presence or absence (Fig.

113    3b). Where we had no genome representative of a particular biotype in which to identify the genes for

114    the cognate pathway, these were predicted using the results of the *in silico* metabolic prediction for

115    known pathways (Fig. 3c). Although there are some discrepancies between pathway presence/absence

116    and historical phenotypes for different biotypes, Fig. 3b shows a robust linkage between inferred

117    phylogeny and biotype. Across the species, these data show that L15 corresponds with the pigmented

118    biogroups A2, A6 and A1, whilst L14, L12 and L9 correspond with the nonpigmented biogroup A4 and

119    biotypes TCT and A5. Furthermore, adding the source of sample isolation shows that niche occupancy

120    also aligns well with the population structure and biotype data (Fig. 3c). In particular, strains

121    representing original biotypes described in the 1980s that were associated with hospitalised patients

122    (TCT, TC, TT, A5, A8) are situated in the same phylogenetic position as contemporary clinically-

123    isolated *S. marcescens*, implying important adaptations that can be linked to risk of disease or greater

124    fitness in hospital environments.

125    **Codon usage redundancy may facilitate a GC shift within *Serratia***

126    Changes in GC content of coding sequences over time have been hypothesised to reflect subtle

127    differences in mutational bias as a consequence of long-term niche adaptation or different lifestyles[40].

128    Given that there are clear differences in the lifestyles and niches between *Serratia* species and intra-

129    species lineages, we investigated the distribution of GC content across the genus. We observe that

130    *Serratia* is broadly divided into two phylogenetically-coherent groups based on whole-genome GC

131    content: *marcescens*, *entomophila*, *ficaria* and *rubidaea* show a GC content of ~59% (59.0-59.9%),

132    whilst *odorifera*, *fonticola*, *plymuthica*, *liquefaciens*, *proteamaculans*, *quinivorans* and *grimesii* have a

133    GC content ranging from 52.7 to 56.1% (Fig. 4a). The singleton "*S. rubidaea*-like" and "*S. marcescens*-

134    like" genomes have an average GC content of 57.7% and 58.9%, respectively, consistent with their

135    positions in the tree adjacent to *rubidaea* or *marcescens*. Additionally, we observed no difference in

136    overall GC content between coding and non-coding regions (Fig. 4b).

137    To understand how this GC pattern impacts on protein coding, we investigated the variation in GC

138    content over the three codon positions for all lineages using the genus-core set of 1655 genes (Fig. 4c;

139    Fig. 2b), and separately for all other genes, designated "non-core". We observed no obvious difference

140    in GC content between genes that were core or non-core at all three codon positions, termed as GC1-3

141    (Fig. 4c). The GC content at GC2 is essentially fixed across the genus (Fig. 4c), whilst GC1 shows a

142    slight skew across the genus, varying by approximately 1%. Codon position GC3 showed a clear bias

143    for A/T-ending codons in low GC species and G/C-ending codons in high GC species, as expected[41]

144    (Fig. 4d). Hence, the difference in average G/C across the genus is largely explained by variation in

145    codon position GC3. For example, GC3 in *S. grimesii* is 20% lower than in *S. marcescens* L9 (Fig. 4c).

146    Taken together, the variations in metabolic capability and GC content between both species and niche-

147    adapted lineages are indicative of long-term niche adaptation within evolutionary timescales.

148    **Pan-genome analysis highlights lineage-specific gene gain and loss as well as intra-genus gene**

149    **flow**

150    The results so far suggest that the pan-genome of *Serratia* lineages is phylogenetically constrained, yet

151    members of *Enterobacteriaceae* are known to have a highly plastic gene content through horizontal

gene transfer (HGT). To investigate this further, we sought to understand genus-wide species plasticity. Plasticity can be estimated by comparing the pan-genome size and complexity against the size of the genus-core gene set (Fig. 2a, b). Given the uneven sampling of some taxa we performed a population structure-aware analysis of the pan-genome, as noted above, in order to define the "genus-core" genome. We overlaid this classification system[38] onto intersections of multi- and single-lineage core genomes (Fig. 2b). Genes were defined as core to a lineage if a gene was present in at least 95 percent of the genomes in each lineage, and the union of all lineage-core genes was defined as the genus-core, consisting of 1655 genes (Fig. 2a,b). This analysis showed lineage- and species-level core gene gain and loss, which are markedly larger in terms of the number of genes when looking at lineages that have a very small sample size. For example, *S. odorifera* (L21; two genomes) and *S. marcescens*-like (L16; one genome), have 1725 and 345 genes core only to those specific lineages (Fig. 2b).

Significant variations in the pan-genome between different *Serratia* species were evident. For example, whilst *S. entomophila* and *S. fonticola* display similar core gene branch lengths, indicative of similar evolutionary timescales, *S. entomophila* has a closed pan genome whilst *S. fonticola* has an open pan-genome (Fig. 2c). The difference in the size of the accessory genome between these two species is 6395 genes, with *S. entomophila* and *S. fonticola* and having accessory genome sizes of 2764 and 9159 genes, respectively (Supplementary Table 2). In contrast, *S. ficaria,* which has similar internal branch lengths, has a more open pan-genome curve (Fig. 2c), and an accessory genome of 3490 genes, despite being represented by fewer genomes in the analysis (Supplementary Table 2), suggesting that different *Serratia* species have varying propensities to gene gain and loss.

Evidence of core gene gain and loss possibly reflective of speciation or niche adaptation can be seen when examining this data. For example, 99 genes are found core to all three lineages in *S. entomophila*, and 41 genes are found core to the entire *S. liquefaciens* complex, which comprises *S. liquefaciens, grimesii, proteamaculans* and *quinivorans* (Fig. 2b). Within the pan-genome we identified lineage- and species-exclusive gene sets, as well as those whose genes are also present at intermediate or rare frequencies across the genus (Fig. 2b). For example, of the 99 *S. entomophila* species-core genes, 35 genes were found across the rest of the genus (Supplementary Fig. 5), shared between both high and low GC members. In contrast, in the 41 genes core to the *S. liquefaciens* complex, very few are found outside the complex, and where they are, they are predominantly present in *S. ficaria* and *S. plymuthica* (Supplementary Fig. 6). The sharing of genes across the genus, implying potential gene flow, raises questions about whether GC3 has been ameliorated to reflect the GC3 trend in a potential recipient genome. Of the 35 genes from the high GC species *S. entomophila* that are found across the low GC species *S. liquefaciens* complex, *S. plymuthica* and *S. fonticola*, the GC3 values of these genes are lower than when found in *S. entomophila* (Supplementary Fig. 11). Similarly, *S. liquefaciens* complex core genes which are also found in *S. ficaria* and *S. plymuthica*, both species with higher GC than members of the *liquefaciens* complex, appear to have ameliorated GC3 (Supplementary Fig. 12).

188     In an attempt to understand the mechanisms by which genes are gained and lost, we focused initially

189     on *S. marcescens*. We investigated the genetic context of the metabolic gene loci associated with

190     different biotypes. In doing so, we identified a hypervariable locus analogous to the plasticity zone seen

191     in *Yersinia*[42]. Variation in this locus explained some of the biochemical differences seen within *S.*

192     *marcescens*. This plasticity zone was located between two sets of tRNAs: one encoding tRNA-Pro$_{ggg}$,

193     the other encoding tRNA-Ser$_{tga}$ and tRNA-Thr$_{tgt}$. It encoded the genes required for gentisate degradation

194     (*nag* gene cassette) and/or tetrathionate reduction (*ttr* gene cassette), present in the same order and

195     orientation across the species, located alongside three sets of genes that were variably present across

196     the *S. marcescens* phylogeny (Fig. 5). These three sets comprised: (1) four genes including one encoding

197     a cyclic AMP (cAMP) phosphodiesterase; (2) an acyltransferase; and (3) a two-gene toxin cassette. A

198     gene predicted to encode a DNA damage inducible protein I (*dinI*) was always present, downstream of

199     the *ttr*/*nag*/cAMP genes and upstream of the acyltransferase gene. In a small number of instances,

200     frameshifts have truncated or split coding genes in this region. Additionally, prophage sequences can

201     also be found flanking these variable sets of genes in some genomes (Fig. 5). Interestingly, in L13 and

202     L9, when the *nag* genes are present, an additional gene, encoding a gene with predicted 3-

203     chlorobenzoate degradation activity, is present 3' of the other genes in the cassette (Fig. 5).

204     Further evidence of gene flow can be seen in *S. marcescens* (Supplementary Figs. 7, 8). Certain genes

205     core to *S. marcescens* L10, L11 and L14 were also found in members of of other lineages, including *S.*

206     *marcescens* L15 and L9, and *S. proteamaculans* L6 (Supplementary Fig. 8). On closer inspection, the

207     genes shared with *S. marcescens* L15 and *S. proteamaculans* L6 comprise a Type VI Secretion System

208     (T6SS). Whilst polyphyletic across *S. marcescens*, this T6SS is syntenic when found in *S. marcescens*

209     but is encoded in a different region of the chromosome when present in *S. proteamaculans*. In both

210     cases, this T6SS is encoded adjacent to a tRNA, and also an integrase in *S. proteamaculans*, potentially

211     suggestive of horizontal transfer across the genus from *marcescens*. There are also 42 genes core to the

212     clinically-associated *S. marcescens* L9, for which 37 are also found polyphyletically across the rest of

213     *S. marcescens* (Supplementary Fig. 8). Many of these genes are predicted to be components of fimbrial

214     usher systems (Supplementary Fig. 8).

215     **Contribution of plasmids to gene content and flow varies across the *Serratia* genus**

216     To understand the potential contribution of plasmids to the plasticity seen in this genus, we searched

217     for plasmid contigs in our genus-wide dataset. This uncovered 409 putative plasmids in 228 genomes

218     and 9 species, 301 (73%) of them present in *S. marcescens* (Fig. 6; Supplementary Table 3). The

219     collection of identified plasmids displays a wide range of sizes (~1-310 kb) and GC content (~30-66%),

220     indicating diversity. However, the distribution of these traits varied amongst *Serratia* species

221     (Supplementary Figs. 13-15). For example, plasmids identified in *S. marcescens* and *liquefaciens* show

222     a markedly broader range of size and GC content compared with those detected in *S. entomophila* and

223    *quinivorans.* Seventy out of a total of 113 predicted plasmid replicons were found within *S. marcescens*

224    L9 and L12, which are the 'clinical' lineages in which 97% and 81% percent of the isolates,

225    respectively, are known to be human- or clinically-associated. In terms of mobility, 296 (72%) of the

226    plasmids were predicted to be conjugative or mobilizable (Fig. 6; Supplementary Table 3), highlighting

227    their potential role in HGT. Consistent with this notion, the predicted host range for this collection of

228    plasmids ranges from single genus to multi-phyla, with the most heterogeneous host range profile

229    observed for plasmids found *S. marcescens* (Fig. 6).

230    A network visualization of the all-versus-all Mash distances[43] calculated for the *Serratia* plasmids was

231    used to explore their diversity. The resulting network comprises 113 clusters, of which 53 (47%)

232    correspond to singletons, illustrating the diversity of *Serratia* plasmids (Fig. 6). Differences in plasmid

233    abundance between clusters were evident from the network, as four top clusters included 36% of the

234    plasmids identified in *Serratia* genomes. Overall, the plasmids clustering was concordant with their size

235    and GC content but also with the host species (Fig. 6c, Supplementary Fig. 15), suggesting limited

236    between-species plasmid transfer within *Serratia.* Nevertheless, some multi-species clusters were

237    identified, perhaps hinting at recent plasmid acquisition events. A cluster formed by plasmids of four

238    non-*marcescens* species was the largest in the network. This cluster mainly consists of large MOBP

239    conjugative plasmids related to the amber disease associated plasmid (pADAP), which is required for

240    virulence of *S. entomophila* and *S. proteamaculans* in the larvae of the grass grub *Costelytra*

241    *zealandica*[12,22]. Interestingly whilst pathogenic potential in *Costelytra zealandica* is a defining trait of

242    *S. entomophila*, the presence of a pADAP-related plasmid was not universal or a defining trait for either

243    *S. entomophila* or *S. proteamaculans,* being found in members of *S. entomophila*, *S. quinivorans*, and

244    *S. proteamaculans,* and a single *S. liquefaciens* genome.

245    Notably, the predicted host range of the plasmids brings an additional perspective on their potential

246    dynamics within the genus. Most plasmids identified in *S. marcescens* appear to be restricted to this

247    species within *Serratia.* Yet many of them have a predicted host range that goes beyond the taxonomic

248    rank of family, implying transfer outside the genus, including two clusters of small ColRNAI plasmids

249    predicted to cross multiple phyla (Fig. 6). In contrast, the largest plasmid cluster (pADAP-like),

250    featuring multiple non-*marcescens* species, seems to be restricted to the *Serratia* genus. Altogether, this

251    picture may suggest that the ecological niche of *S. marcescens* has favoured plasmid exchange with

252    diverse hosts outside the genus but has also promoted plasmid containment within the species in

253    *Serratia*. The diversity of plasmids identified in *S. marcescens* and their predicted host range thus

254    implies a major role for this species in the gene flow outside the genus and to a lesser but relevant extent

255    within it.

256    **A genomics perspective on a historical phenotype**

257    A famous characteristic often popularly associated with *Serratia* spp. is the production of the red

258  pigment prodigiosin[17]. However, in fact, prodigiosin production has only been observed in *S.*
259  *marcescens* biogroups A1, A2 and A6, some *rubidaea* and some *plymuthica* isolates[21]. The *pig* gene
260  cluster comprises fourteen genes (*pigA-pigN*) required for the production of prodigiosin[44]. Searching
261  across the genus for *pig* gene cluster loci and flanking regions showed that, consistent with the earlier
262  biotyping observations, the *pig* cluster is only encoded in certain *S. marcescens*, *S. rubidaea* and *S.*
263  *plymuthica* genomes (Figs. 3, 7) which are associated with biotypes or biogroups known to be
264  pigmented. In each case, the cluster presents exactly the same contiguous order of genes (*pigA-pigN*),
265  however, notably, it is found in separate genomic loci in each of the three different species (Fig. 7).
266  Representative *pig* gene clusters from each species share ~77-80% identity at the nucleotide level (Fig.
267  7b), which is similar to the shared nucleotide percentage identity between these species at fully syntenic
268  regions in the chromosome. This indicates that the *pig* gene cluster has been acquired horizontally on
269  at least three separate occasions.

## Discussion

272  With advancements in technology, the methods used to delineate and decipher prokaryotic species
273  boundaries have changed over time, as researchers attempt to resolve the shortcomings of earlier
274  approaches and build upon the understanding of biology at any given point in time. This study has, in
275  part, investigated the relationship by which species level boundaries have been determined within a
276  genus, namely phenotypic characterisation and whole genome sequencing. It also highlights how, in
277  order to make appropriate conclusions from these approaches, the currently available data requires to
278  be constantly filtered, checked and reviewed.

279  Following its original identification in the early 19th century, the nomenclature and number of species
280  within *Serratia* underwent several iterations as additional strains with similar, yet distinct phenotypes
281  were identified and added to an expanding membership of the genus[17,21]. Then in the 1970s and 80s,
282  comprehensive biochemical and phenotypic characterisation along with the use of DNA-DNA
283  hybridisation, allowed the genus to be defined as a collection of ten clearly defined species. Since the
284  advent of the genomic era, despite the potential of genomic approaches to resolve fine-scaled
285  differences between taxa, no similar-scale work within *Serratia* has been attempted, nor do we have a
286  robust phylogenetic framework against which we are able to recognize novel *Serratia* spp or emerging
287  lineages. Such a framework is also required to resolve confusion over existing species. For example,
288  strain DSM 21420, a nematode-associated strain proposed to belong to *S. nematodiphila*[45], sits within
289  the broadly non-clinical *S. marcescens* L15, suggesting that it does not in fact represent a separate
290  species. Conversely, the identification of singleton ANI phylogroups and FastBaps lineages (*S.*
291  *marcescens*-like L16 and *S. rubidaea*-like L23) highlights that there is likely further species diversity
292  to be discovered. This may be partly due to geography and lack of sampling: the strain that occupies

10

293    L16 (MSU97) was sourced from a plant in the Carrao River in Venezuela[46], a region which is not highly
294    sampled.

295    It is interesting to consider how the computational approaches used here to classify and describe the
296    genus parallel the original biotyping. In the earlier studies, *in vitro* DNA-DNA hybridisation was used
297    to assess genomic relatedness between novel *Serratia* strains[5,11], an approach for which ANI is in many
298    ways an *in silico* proxy, whilst the connection between *in silico* prediction of metabolic potential and
299    the lab-based tests detecting the corresponding metabolic pathway in the original biochemical-based
300    biotyping is obvious. Furthermore, in some cases these biotypes highlight further clusters within
301    lineages that match branching within the phylogeny. For example, biotypes C1c, EB and RB, and
302    biotypes A1b, A1a, A6 and A2 are all monophyletic within *S. proteamaculans* L6 and *S. marcescens*
303    L15, respectively (Figs. 3b, 5). This highlights just how accurate the original biochemical-based typing
304    was for defining species.

305    This accuracy is particularly striking given that we have observed that presence or absence of metabolic
306    pathways (corresponding to the historic biotyping tests used) can be due to repeated gene gain or loss
307    in the same locus over short evolutionary distances. For example, the genes required for the degradation
308    of gentisate and the reduction of tetrathionate are gained and lost within and between lineages in *S.*
309    *marcescens,* in the same locus and also in the same conserved order (Fig. 5). This would explain why
310    the original phenogrouped biotypes based on biochemical typing had "variable" results for certain
311    metabolic tests, such as gentisate degradation being observed to be variable in the clinical biotypes A8a
312    and A8c[92]. This locus-specific pathway gain and loss in historic isolates is also seen in more
313    contemporary strains (Fig. 5). The maintenance of this plasticity zone suggests that there are transient
314    and frequently re-occurring environmental selective pressures where the benefit and cost of these
315    pathways is great enough to provide selection both for and against them. In other words, the data suggest
316    that both the loss and re-acquisition of these elements is of benefit to *S. marcescens* at various times.

317    It is also noteworthy that the environment from which strains were isolated across our assembled dataset
318    tends to match the environments and niches with which each biotype was historically associated[21]. Of
319    particular interest is the observation that the predominantly hospital-associated biotypes of *S.*
320    *marcescens* that were defined in the 1980s (A5, A8, TCT) sit within L9 and L12 defined in the current
321    study. These lineages are mainly comprised of recently-sequenced genomes from hospital settings,
322    including a large collection of clinically-derived *S. marcescens* isolates from the UK that represent the
323    recent emergence of hospital-adapted clones exhibiting recent acquisition of MDR phenotypes[18]. The
324    fact that these lineages of clinically-associated *S. marcescens* were identified back in the 1970s and 80s
325    shows that the original biochemical characterisation of *Serratia* captured the emergence of *S.*
326    *marcescens* lineages that have subsequently been reported to be associated with human disease many
327    times in recent years[18,19,47–52]. The apparent specialisation of *S. marcescens* L9 to be a clinically-adapted

11

328  pathogen is further highlighted by plasmid replicon identification and the types of lineage-specific core
329  genes observed. The identification of numerous plasmid replicons in these lineages (L9, L12 and L14)
330  as opposed to the rest of the genus is perhaps unsurprising, given that most known plasmids are
331  associated with multi-drug resistance and hospital environments. Fimbrial genes are well-known
332  pathogenicity factors and multiple different fimbrial genes are found to be core to L9 but accessory to
333  multiple other *S. marcescens* lineages. This potential gene flow from L9 across the rest of *S. marcescens*
334  may be one reason why isolates from more "environmental" *S. marcescens* lineages are still isolated
335  from nosocomial settings. In these other lineages, *S. marcescens* is still an opportunistic pathogen, with
336  nosocomial isolates being genetically similar to strains that have colonised or infected plants, insects or
337  other environments. Indeed, bee-associated *S. marcescens* cause infections in bees in a similar manner
338  to how *S. marcescens* can cause bloodstream infections in preterm neonates[29]. Taking the historic
339  biotyping data along with the population structure defined here, the combined data suggest that *S.*
340  *marcescens* is highly plastic in its nature yet can also become specialised in a particular niche.

341  Speciation and niche specialisation events or processes are seen across the phylogeny, as highlighted
342  by the long branch lengths between divisions, separations in GC content, variation in metabolic
343  potential, and enrichment for certain isolation source sites in different lineages. These divisions likely
344  represent ancient speciation events that have occurred as *Serratia* has spread to be ubiquitous
345  worldwide. As mentioned above, changes in GC content can be a response to long-term niche
346  adaptation, however there is no commonly held theory or understanding of the possible reasons that
347  underpin this. One possible factor that may have influenced the variation in GC content observed across
348  *Serratia* is a difference in ideal growth temperature: higher GC *Serratia* species tend to be able to grow
349  better at higher temperatures than lower GC *Serratia* species[21]. Another possibility is that the observed
350  GC-dependent change in codon usage, which does not alter protein sequence or function, is indicative
351  of a shift to an optimal set of codons for each particular *Serratia* species, although the evolutionary
352  pressure that would drive such a shift is not clear. Importantly, however, this division in GC content
353  does not seem to be a barrier for gene flow in *Serratia*, since genes core to the high GC species *S.*
354  *entomophila* can also be found in polyphyletic and variable patterns across the genus, including in low
355  GC *Serratia* species (Supplementary Fig. 11). However, it is formally possible that these genes could
356  have been horizontally acquired from non-*Serratia* sources.

357  This study also provides definitive genomic evidence to explain the variation in a classical *Serratia*
358  phenotype, namely the production of the red pigment prodigiosin (Fig. 7). The high level of synteny
359  within the *pig* gene cluster together with the absence of homology in the flanking regions indicates that
360  the ability to produce prodigiosin has been acquired on at least three separate occasions within *Serratia*,
361  namely in subsets of *S. marcescens*, *S. plymuthica* and *S. rubidaea*. Given the relatively low shared
362  nucleotide identity, it is unclear when and how these genes were incorporated into the chromosome,

363    and whether each event reflects gene flow within the genus or separate acquisition from an external

364    source. This genomic evidence of separate acquisition of the *pig* clusters matches the historical metadata

365    noting that *S. marcescens*, *S. plymuthica,* and *S. rubidaea* all variably produce a red pigment[21].

366    Prodigiosin has been reported to display many functions, including anti-protozoal, anti-fungal, anti-

367    bacterial, immunosuppressive, and anti-cancer activity[44]. The biological advantage for these individual

368    *Serratia* species, or subsets thereof, to be able to produce prodigiosin is unclear, however it could reflect

369    a degree of convergent evolution within *Serratia*, or perhaps the varied potential functions of

370    prodigiosin may provide different fitness benefits to different species. Further evidence for convergent

371    evolution in the genus is provided by the observation that members of both *S. proteamaculans* and *S.*

372    *entomophila* carry pADAP, which is required for the pathogenesis of grass grub larvae[22,53].

373    In conclusion, we have demonstrated the power of combining phenotypic metadata with a

374    comprehensive and balanced genomics-based phylogeny to define an important and diverse bacterial

375    genus, its plasticity and its niche adaptation. The dataset and phylogeny that we present here will

376    provide a vital platform for future work, including in the tracking of further emergence of pathogenic

377    *Serratia* or changes in the portfolio of anti-microbial resistance genes or pathogenicity factors.

## Material and Methods

### Bacterial strains

Bacterial isolates sequenced in this study are listed in Supplementary Table 1, along with relevant metadata and summaries of sequencing and assembly statistics.

### Bacterial culture and resuscitation, genomic DNA isolation and sequencing

279 isolates in the Institut Pasteur collection were successfully resuscitated from agar stabs kept in cold storage for ~20 years. Isolates were resuscitated in the original agar stabs with 2-3 ml of Tryptic Soy Broth and incubated static and upright at 30°C for up to three days, or until clear signs of growth were visible, followed by sub-culture on solid LB media. In rare cases of mixed colony morphology, or abnormal looking colonies, a number of colonies were selected and streaked through two to three times. In such cases, the *Serratia* were identified, where possible, by red pigmentation and/or a strong potato-like odour. In cases of mixed pigmentation, a representative colony of each type of pigment type (or lack of pigment) were taken forward. DNA extraction was carried out using the Maxwell 16 Cell DNA purification kit (Promega) on the automated Maxwell 16 MDx instrument (Promega), according to the manufacturer's instructions. 400 µl of mid-log culture (grown at 30°C in LB), sub-cultured from a liquid overnight culture, was used for DNA extraction. DNA samples were sequenced using the Illumina HiSeq X10 platform (Illumina, Inc) at the Wellcome Sanger Institute. DNA fragments of approximately 450 bp were produced from 0.5 µg DNA for Illumina library creation and were sequenced on a 150 bp paired-end run.

42 isolates from UK hospitals were received from frozen stocks, freshly streaked plates, or in bead suspensions, and were grown on solid media to ensure uniform single colonies. As with isolates from the Pasteur collection, samples from mid-log cultures were used for DNA extraction. DNA samples were sequenced using short-read technology only, or a hybrid approach of both long-read and short-read technology, as detailed in Supplementary Table 2. For short-read sequencing, DNA was extracted using a DNeasy extraction kit (Qiagen). DNA quality was assessed using a Qubit 3.0 (Invitrogen) and Bioanalyzer (Agilent), then subsequently diluted to a concentration of 0.4 ng/µl. DNA library preparation was performed using the Illumina Nextera protocol and PCR clean up was performed using AMPure beads (Beckman). Multiplexed samples were then run on the MiSeq (Illumina). Adaptor sequences were automatically trimmed by the MiSeq platform and then raw reads were downloaded from basespace in FASTQ format. For long-read sequencing, high molecular weight DNA was isolated using the MasterPure DNA Purification kit (Epicentre, no. MC85200). Sequencing was performed using the PacBio Sequel (Pacific Biosciences) or MinION (Oxford Nanopore Technologies) sequencing platforms. For PacBio sequencing, 10 µg DNA was sequenced using polymerase version P6 and C4 sequencing chemistry reagents. For MinION sequencing, 5 µg DNA in 35 µl nuclease-free water for each sample was sequenced using the SQK-LSK108 kit using a FLO-MIN106 flow cell. DNA ends

14

413  were repaired and dA-tailed using NEBNext End Repair/dA-tailing module, following by ligation of

414  barcodes. DNA concentration and clean up steps were performed using AMPureXP beads (New

415  England Biolabs). 12 samples (from 12 isolates) were multiplexed on a single MinION run. Basecalling

416  and demultiplexing was performed by Albacore v2. In all cases, kits were used according to the

417  manufacturers' instructions.

**Sequence data quality control**

419  Read sets obtained from all samples were compared to the MiniKraken database by Kraken v0.10[54],

420  and then corrected using Bracken[55] which assigns reads to a specific reference sequence, species or

421  genus. If reads were not able to be assigned to a taxonomic class, they were classed as 'unclassified'.

422  Any read sets that belonged to genera other than *Serratia* were discarded from any further analysis,

423  along with any assemblies obtained from those read sets.

424  Any read sets with more than an estimated five percent of heterozygous SNPs across the whole genome

425  were removed from further analysis, in addition to any assemblies obtained from those read sets.

426  Heterozygous SNPs were calculated using a software pipeline from the pathogen informatics team at

427  the Wellcome Sanger Institute. Specifically, read sets from each *Serratia* sample were aligned to an

428  appropriate reference for that sample, given the taxonomic profile from the Kraken and Bracken output.

429  Reads were aligned to the reference using bwa v0.7.17[56], and parsed using samtools v0.1.19[57] and

430  bcftools v0.1.19[57]. Reads were considered as heterozygous if there were at least two variants at the same

431  base, both supported by a number of reads that was fewer than 90 percent of the total reads mapped to

432  that site. Read coverage to each strand was considered independently. The minimum total coverage

433  required was 4x, and the minimum total coverage for each strand was 2x. Calculated heterozygous SNP

434  coverage was then predicted by scaling the number of observed heterozygous SNPs against the

435  proportion of the reference that was covered by read mapping.

436  Eight genome sequences from the Pasteur collection dataset and one from the UK hospitals set were

437  removed due to the above criteria. In addition, a number of the isolates resuscitated from the Pasteur

438  collection were duplicate samples of the same strain. After inspection of preliminary phylogenetic trees

439  from core-gene alignments (see below), a further 56 genomes were removed from the Pasteur collection

440  dataset due to being duplicates of the same-named strain.

**Publicly available genome sequences**

442  Previously-published, publicly-available assembled genome sequences were downloaded from the

443  NCBI GenBank database (https://ftp.ncbi.nlm.nih.gov/genomes/genbank/) as of 19/03/2019. Genomes

444  were downloaded if the species was attributed to any of the following: *Serratia sp.*, *odorifera*, *rubidaea*,

445  *plymuthica*, *liquefaciens*, *grimesii*, *oryzae*, *proteamaculans*, *quinivorans*, *nematodiphila*, *ficaria*,

446  *entomophila* or *marcescens.* Assemblies smaller than 4.5 Mbp or larger than 6.5 Mbp were removed

447     from the analysis, along with any assemblies comprised of more than 250 contigs. Quast v4.6.0[58] was

448     used to extract statistics for genomes and genomic assemblies, specifically whole genome GC content,

449     number of contigs and assembly size. Initial phylogenetic trees with additional non-*Serratia* reference

450     sequences (*Yersinia enterocolitica*, *Rahnella aquatilis* and *Dickeya solani*) were computed, and

451     genomes detemined by visual inspection as being non-*Serratia* or close to non-*Serratia* members of

452     *Enterobacteriacaeae* were removed from any subsequent analysis. Ten genomes were excluded on this

453     basis, including several so-called *Serratia sp.* and *Serratia oryzae*.

454     **Genome assembly and annotation**

455     The assembly method used for genome assembly and annotation for each genome are detailed in

456     Supplementary Table 1. For samples sequenced using short-read only data, genomes were assembled

457     in two different ways depending on their origin. Isolates in the Institut Pasteur collection were

458     assembled through assembly pipelines at the Wellcome Sanger Institute. For each sample, sequence

459     reads were used to create multiple assemblies using VelvetOptimiser v2.2.5

460     (https://github.com/tseemann/VelvetOptimiser) and Velvet v1.2[59]. An assembly improvement step was

461     applied to the assembly with the best N50 and contigs were scaffolded using SSPACE[60] and sequence

462     gaps filled using GapFiller[61]. For isolates from UK hospitals that were only sequenced by short-read

463     technology, these short reads were assembled using SPAdes v3.6.1[62], using default settings.

464     For hybrid short- and long-read assemblies of selected isolates from UK hospitals, genomes were

465     assembled using Unicycler v0.4.7[63]. Long-read-only assemblies from MinION or PacBio long reads

466     were generated first, using Canu v1.6[64], with the expected genome size set as 5.4 Mbps, the minimum

467     read length and overlap length set to 100 bp, and "corOutCoverage" set to 1000. Long-read assemblies

468     were then used as input to Unicycler, using the --existing_long_read_assembly flag. Sets of paired-end

469     Illumina reads were then used as input to Unicycler alongside this long-read assembly and also the long

470     reads. The "--mode" flag was set to "normal". In the event that Unicycler was not able to produce

471     circularised assemblies, Circlator v1.5.5[65] was used to circularise assemblies.

472     Assembled genomes were then annotated using Prokka v.1.13.3[66].

473     **Pan-genome analysis**

474     Pan-genomes were calculated from 664 *Serratia* sequences using Panaroo v1.2.3[67], with Prokka-

475     annotated genomes as input. For initial protein clustering, a protein similarity threshold was set at 95

476     percent (0.95). The subsequent clustering of these groups into protein families was performed using a

477     threshold of 70 percent identity (0.7). The "--clean-mode" flag was set to "moderate". A core-gene

478     alignment was created using the "-a" flag, specifying mafft as the aligner using the "--aligner" flag,

479     with core genes specified by being present in at least 95 percent of genomes (631/664). Pan-genome

480     gene accumulation curves were generated using the *specaccum* function from the R package Vegan

16

481    v2.5.7[68], with 100 random permutations.

482    Population structure-aware classification of genes across the genus was performed upon the gene

483    presence/absence matrix created by Panaroo through the use of the twilight analysis package[38]. Groups

484    were defined by the lineages set by the third level of Fastbaps clustering (see below), and singleton

485    lineages were included in the analysis (""--min_size 1"). The core and rare thresholds were set at 0.95

486    and 0.15, respectively.

487    Preliminary core-gene alignments using the pan-genome software Roary v3.12.0[69], including all

488    downloaded genomes from the NCBI GenBank datbase, duplicate genomes from the Pasteur collection

489    and non-*Serratia* Enterobacteriacaeae members, were computed for initial tree-drawing to remove

490    contaminants and assess whether duplicate strains (from data supplied in strain name information, for

491    example, labels on agar stabs from strains in the Pasteur collection) were found in the same position in

492    the tree. Non-*Serratia Enterobacteriacaeae* were also used to determine the location of the root for all

493    visualisations of the *Serratia* genus phylogenetic tree.

**494    Clustering, phylogroup determination, core-gene alignment filtering and phylogenetic tree**
**495    construction**

496    For the *Serratia* phylogeny, a concatenated core-gene alignment from 2252 genes (2,820,212 bp in

497    length) from Panaroo v1.2.3[67] (as described above) was filtered to remove monomorphic sites that were

498    exclusively A, T, G or C using SNP-sites v2.5.1[70]. The resulting alignment was 398,551 bp in length.

499    IQtree v.1.6.10[71] was then used for maximum-likelihood tree construction using 1000 ultrafast

500    bootstraps[72] using the TIM2e+ASC+R4 model chosen using modelfinder[73]. Both the ultrafast bootstraps

501    and modelfinder were implemented in IQtree. The *Serratia* phylogenetic tree was rooted at the position

502    of a *Yersinia enterocolitica* outgroup root after analysis of preliminary trees based on exclusively

503    polymorphic variant sites (filtered using SNP-sites v2.4.1) from preliminary core-gene alignments

504    (determined using Roary v.3.12.0 as described above). Trees were constructed using modelfinder

505    implemented in IQtree v1.6.10, followed by tree construction using IQtree v.1.6.10.

506    Whole-genome assemblies were compared in a pairwise manner using fastANI v1.3[74], and phylogroups

507    determined through clustering these comparisons using a cutoff of 95% average nucleotide identity

508    (ANI). Genomic assemblies were then clustered base on this cutoff value, using the script

509    fastANI_to_clusters.py which uses the networkx package (https://networkx.github.io/), and visualised

510    using Cytoscape v3.7.1[75]. The phylogeny was partitioned into lineages defined through hierarchical

511    bayesian clustering using Fastbaps v1.0.4[76]. Fastbaps was used to cluster the phylogeny over four levels,

512    with the third levels selected for lineage designation. The SNP sites-filtered core-gene alignment was

513    used as input to Fastbaps, alongside the rooted phylogenetic tree to provide a guide for the hierarchical

514    partitioning.

17

**Functional and metabolic pathway analysis**

*In silico* reconstruction of metabolic pathways was performed using Pathway tools v23.5[39] , using a multi-processing wrapper tool mpwt (https://github.com/AuReMe/mpwt)[77]. In order to arrange input data into the appropriate format, and subsequently parse the output, a collection of Python and R scripts were written (https://github.com/djw533/pathwaytools_gff2gbk). Further specific information about how to run this can be found in the readme hosted at the github repository. In brief: Representative protein sequences for each of the 47,743 protein family groups identified in the pan-genome analysis were extracted from the pan-genome graph-associated data using Cytoscape v3.7.1, and functionally annotated using EggNOG-mapper v1.0.3[78], using the following flags "-m diamond -d none --tax_scope auto --go_evidence non-electronic --target_orthologs all --seed_ortholog_evalue 0.001 -- seed_ortholog_score 60 --query-cover 20 --subject-cover 0 --override". Using the EggNOG annotations from representative protein sequences, annotated genomes (as .gff files) were updated with the Enzyme Commision (EC) numbers, Gene Ontology (GO) terms and predicted function for each protein family group from the pan-genome analysis, using the script gffs2gbk.py in pathwaytools_gff2gbk. This script also appropriately organises the input data required for mpwt given a file listing the taxon IDs for each genome. Pathway tools was then run by running the multi-processing wrapper mpwt was then run with the "--patho" and "--taxon_id" flags, whilst providing the file containing taxon ids linked to each genome. The *in silico*-reconstructed metabolic pathways for all genomes were then collated using compare_pgdbs.R in pathwaytools_gff2gbk, and downstream analysis conducted in R, as shown in https://github.com/djw533/Serratia_genus_paper/figure_scripts.

**Plasmid replicon identification**

Plasmid sequences were identified in the collection of *Serratia* genome assemblies with the MOB-recon tool using the MOB-suite v3.0.3 databases and default settings[79]. Characterisation of the identified plasmids, including predicted transferability of the plasmid, was performed with MOB-typer from the MOB-suite package. Charts illustrating plasmid counts and features were generated in R using ggplot2[80]. K-mer-based sketches of the plasmid sequences (s=1000, k=21) were generated with the mash v2.3 sketch algorithm[43]. Pairwise mutation distances between sketches were estimated using mash dist with a distance threshold of 0.05 and otherwise default settings. The resulting all-pairs distance matrix was used for graph-based clustering of the plasmid sequences in Cytoscape v3.8.2[75] using the "connected components cluster" algorithm from the clusterMaker2 v2 app[81].

**GC content analysis**

Whole-genome GC content calculated using Quast v4.6.0. GC content for each gene, and the average GC value for codon positions 1, 2 and 3 for each gene was then calculated using the script GC_from_panaroo_gene_alignments.py, which uses the gene_data.csv file created from Panaroo

549    (detailed above). Intragenic nucleotide sequence was extracted for all protein-encoding sequences using

550    gals_parser_with_fasta.py with the "-t nuc" flag. Intergenic GC values were then calculated by using

551    Bedtools[82] complement from Bedtools v2.29.0 to identify the inverse of all coding regions (i.e all

552    intergenic regions). Bedtools getfasta from Bedtools v2.29.0 was then used to extract the intergenic

553    regions as nucleotide sequence. Average GC values for the total intergenic and intragenic regions were

554    then calculated using get_gc_content.py.

555    **Retrieval of specific gene clusters**

556    Gene clusters containing co-localised *pig* genes (*pigA-M*) were identified using Hamburger

557    (github.com/djw533/hamburger), which uses protein HMM profiles for each target gene in the gene

558    cluster. User-set parameters define the minimum number of HMMsearch[83] hits required to report the

559    presence of each system in a genome, in addition to the maximum number of non-hit genes that are

560    permitted between two hit genes in a contiguous set of genes. Gene clusters were reported as prodigiosin

561    clusters for loci encoding at least nine genes containing Pfam domains characteristic of 11 of the 14 *pig*

562    genes with no more than five non-model genes between any "hits". Extracted genomic sequences were

563    then compared using blast+ v2.2.31[84] and genoplotR v0.8.11[85]. Blastn was used with the flags "--task

564    Blastn --perc_identity 20 --evalue 10000". Functions created to use these can be found in

565    micro.gen.extra on https://github.com/djw533/micro.gen.extra.

566    Gene clusters around other genes of interest, such as the plasticity zone in *S. marcescens* located

567    between tRNA-Pro$_{ggg}$ and tRNA-Ser$_{tga}$, were extracted using the script pull_out_around_point.py, and,

568    if in the unwanted orientation, flipped using gff_reverse.py.

569    **Phage prediction**

570    Phage regions were predicted using Phaster[86] on the webserver (https://phaster.ca/), using default

571    settings.

572    **Data visualisation**

573    Phylogenetic trees were visualised using the R package ggtree v2.4.2[88]. Synteny of regions of bacterial

574    genomes extracted by Hamburger were visualised using the R package genoplotR v0.8.11[85]. Genetic

575    organisation of genes were plotted using the R package gggenes v0.4.1 (https://wilkox.org/gggenes/).

576    Other plots were created using the R package ggplot2 v3.3.5[80]. As mentioned above, networks were

577    viewed using Cytoscape v3.7.1[75]. Sets were visualised as Upset plots using UpsetR v1.4.0[89].

578

579    **Code availability statement**

580    All custom scripts for which github repositories are not specified above can be found in

19

581    https://github.com/djw533/Serratia_genus_paper/, along with all Rscripts used to plot figures. Rscripts

582    make use of the tidyverse[87] collection of packages. R version 4.0.3 was used for all analysis and plotting.

583    Other packages can be found at https://github.com/djw533/hamburger,

584    https://github.com/djw533/micro.gen.extra, and https://github.com/djw533/pathwaytools_gff2gbk.

585

**Data availability statement**

587    The whole genome sequences generated during the current study are in the process of deposition and

588    will be made fully available. The read sets for the majority of the sequences (from the Institut Pasteur

589    isolates) are already available in the ENA (https://www.ebi.ac.uk/ena/browser/home), with project

590    number PRJEB24638. Other whole genome sequences analysed during the study are available from

591    NCBI GenBank (https://ftp.ncbi.nlm.nih.gov/genomes/genbank/), with the accession numbers for the

592    individual sequences given in Supplementary Table 1.

593

**References**

595  1.    Merlino, C. P. Bartolomeo Bizio's letter to the most eminent priest, Angelo Bellani, concerning
596      the phenomenon of the red-colored polenta [translated from the Italian]. *Journal of Bacteriology*
597      (1924).
598  2.    Grimont, P. A. D. & Dulong de Rosnay, H. L. C. Numerical Study of 60 Strains of *Serratia*.
599      *Journal of General Microbiology* **72**, 259-268 (1972).
600  3.    Grimont, P. A. D., Grimont, F. & Dulong de Rosnay, H. L. C. Taxonomy of the genus Serratia.
601      *Journal of General Microbiology* **98**, 39–66 (1977).
602  4.    Grimont, F., Grimont, P. A. D. & Dulong de Rosnay, H. L. C. Characterization of *Serratia*
603      *marcescens, S. liquefaciens, S. plymuthica* and *S. marinorubra* by Electrophoresis of their
604      Proteinases. *Journal of General Microbiology* **99**, 301-310 (1977).
605  5.    Grimont, P. A. D. *et al.* Deoxyribonucleic Acid Relatedness Between *Serratia plymuthica* and
606      Other *Serratia* Species, with a Description of *Serratia odorifera* sp. nov. (Type Strain: ICPB
607      3995). *International Journal of Systemic Bacteriology* **28**, 453-463(1978).
608  6.    Grimont, P. A. D., Grimont, F. & Starr, M. P. *Serratia ficaria* sp. nov., a bacterial species
609      associated with Smyrna figs and the fig wasp *Blastophaga psenes*. *Current Microbiology* **2**,
610      277–282 (1979).
611  7.    Gavini, F. *et al. Serratia fonticola*, a New Species from Water *International Journal of Systemic*
612      *Bacteriology* **29**, 92-101 (1979).
613  8.    Holmes, B. Proposal to Conserve the Specific Epithet *liquefaciens* Over the Specific Epithet
614      *proteamaculans* in the Name of the Organism Currently Known as *Serratia liquefaciens* (Grimes
615      and Hennerty 1931) Bascomb et al. 1971. Request for an Opinion. *International Journal of*
616      *Systemic Bacteriology* **30**, 220-222 (1980).
617  9.    Grimont, P. A. D., Grimont, F. & Starr, M. P. *Serratia* species isolated from plants. *Current*
618      *Microbiology* **5**, 317–322 (1981).
619  10.   Grimont, P. A. D., Grimont, F. & Irino, K. Biochemical characterization of *Serratia liquefaciens*
620      *sensu stricto*, *Serratia proteamaculans*, and *Serratia grimesii* sp. nov. *Current Microbiology* **7**,
621      69–74 (1982).
622  11.   Grimont, P. A. D., Irino, K. & Grimont, F. The *Serratia liquefaciens-S. proteamaculans-S.*
623      *grimesii* complex: DNA relatedness. *Current Microbiology* **7**, 63–67 (1982).
624  12.   Grimont, P. A. D., Jackson, T. A., Ageron, E. & Noonan, M. J. *Serratia entomophila* sp. nov.

Associated with Amber Disease in the New Zealand Grass Grub *Costelytra zealandica*. *International Journal of Systematic Bacteriology* **38**, 1–6 (1988).

13. Murdoch, S. L. *et al.* The opportunistic pathogen *Serratia marcescens* utilizes Type VI secretion to target bacterial competitors. *Journal of Bacteriology* **193**, 6057–69 (2011).

14. Williamson, N. R., Fineran, P. C., Ogawa, W., Woodley, L. R. & Salmond, G. P. C. Integrated regulation involving quorum sensing, a two-component system, a GGDEF/EAL domain protein and a post-transcriptional regulator controls swarming and RhlA-dependent surfactant biosynthesis in *Serratia*. *Environmental Microbiology* **10**, 1202–1217 (2008).

15. Kurz, C. L. *et al.* Virulence factors of the human opportunistic pathogen *Serratia marcescens* identified by *in vivo* screening. *The EMBO journal* **22**, 1451–60 (2003).

16. Khanna, A., Khanna, M. & Aggarwal, A. *Serratia marcescens* - a rare opportunistic nosocomial pathogen and measures to limit its spread in hospitalized patients. *Journal of Clinical and Diagnostic Research* **7**, 243–6 (2013).

17. Mahlen, S. D. *Serratia* infections: from military experiments to current practice. *Clinical Microbiology Reviews* **24**, 755–91 (2011).

18. Moradigaravand, D., Boinett, C. J., Martin, V., Peacock, S. J. & Parkhill, J. Recent independent emergence of multiple multidrug-resistant *Serratia marcescens* clones within the United Kingdom and Ireland. *Genome Research* **26**, 1101–1109 (2016).

19. Karkey, A. *et al.* Outbreaks of *Serratia marcescens* and *Serratia rubidaea* bacteremia in a central Kathmandu hospital following the 2015 earthquakes. *Transactions of The Royal Society of Tropical Medicine and Hygiene* **112**, 467–472 (2018).

20. Dubouix, A. *et al.* Epidemiological investigation of a *Serratia liquefaciens* outbreak in a neurosurgery department. *Journal of Hospital Infection* **60**, 8–13 (2005).

21. Grimont, F. & Grimont, P. A. D. The Genus *Serratia*. in *Prokaryotes* (eds. Martin Dworkin, Stanley Falkow, Eugene Rosenberg, Karl-Heinz Schleifer & Erko Stackebrandt) 219–244 (Springer-Verlag, 2006).

22. Hurst, M. R. H., Glare, T. R., Jackson, T. A. & Ronson, C. W. Plasmid-Located Pathogenicity Determinants of *Serratia entomophila*, the Causal Agent of Amber Disease of Grass Grub, Show Similarity to the Insecticidal Toxins of *Photorhabdus luminescens*. *Journal of Bacteriology* **182**, 5127-5138 (2000).

23. Hurst, M. R. H., Glare, T. R. & Jackson, T. A. Cloning *Serratia entomophila* antifeeding genes - a putative defective prophage active against the grass grub *Costelytra zealandica*. *Journal of Bacteriology* **186**, 5116–28 (2004).

24. Nuñez-Valdez, M. E. *et al.* Identification of a putative Mexican strain of *Serratia entomophila* pathogenic against root-damaging larvae of *Scarabaeidae* (Coleoptera). *Applied and Environmental Microbiology* **74**, 802–10 (2008).

25. Rodríguez-Segura, Z., Chen, J., Villalobos, F. J., Gill, S. & Nuñez-Valdez, M. E. The lipopolysaccharide biosynthesis core of the Mexican pathogenic strain *Serratia entomophila* is associated with toxicity to larvae of *Phyllophaga blanchardi*. *Journal of Invertebrate Pathology* **110**, 24–32 (2012).

26. Hurst, M. R. H. *et al.* *Serratia proteamaculans* Strain AGR96X Encodes an Antifeeding Prophage (Tailocin) with Activity against Grass Grub (*Costelytra giveni*) and Manuka Beetle (*Pyronota* Species) Larvae. *Applied and Environmental Microbiology* **84**, (2018).

27. Flyg, C., Kenne, K. & Boman, H. G. Insect Pathogenic Properties of *Serratia marcescens*: Phage-resistant Mutants with a Decreased Resistance to Cecropia Immunity and a Decreased Virulence to Drosophila. *Microbiology* **120**, 173–181 (1980).

28. Ishii, K., Adachi, T., Hara, T., Hamamoto, H. & Sekimizu, K. Identification of a *Serratia marcescens* virulence factor that promotes hemolymph bleeding in the silkworm, *Bombyx mori*. *Journal of Invertebrate Pathology* **117**, 61–67 (2014).

29. Raymann, K., Coon, K. L., Shaffer, Z., Salisbury, S. & Moran, N. A. Pathogenicity of *Serratia marcescens* strains in honey bees. *mBio* **9**, e01649-18 (2018).

30. Ashelford, K. E., Fry, J. C., Bailey, M. J. & Day, M. J. Characterization of *Serratia* isolates from soil, ecological implications and transfer of *Serratia proteamaculans* subsp. *quinovora* Grimont *et al.* 1983 to Serratia quinivorans corrig., sp. nov. *International Journal of Systematic and*

679       *Evolutionary Microbiology* **52**, 2281–2289 (2002).

680   31.  Lim, Y.-L. L. *et al.* Complete genome sequence of *Serratia fonticola* DSM 4576T, a potential
681       plant growth promoting bacterium. *Journal of Biotechnology* **214**, 43–44 (2015).

682   32.  Abebe-Akele, F. *et al.* Genome sequence and comparative analysis of a putative
683       entomopathogenic *Serratia* isolated from *Caenorhabditis briggsae*. *BMC Genomics* **16**, 531
684       (2015).

685   33.  Petersen, L. M. & Tisa, L. S. Friend or foe? A review of the mechanisms that drive *Serratia*
686       towards diverse lifestyles. *Canadian Journal of Microbiology* **59**, 627–640 (2013).

687   34.  Cheng, T. H. *et al.* Genome Sequence of *Serratia marcescens* subsp. *sakuensis* Strain K27, a
688       Marine Bacterium Isolated from Sponge (*Haliclona amboinensis*). *Genome Announcements* **6**,
689       e00022-18 (2018).

690   35.  Matilla, M. A., Udaondo, Z. & Salmond, G. P. C. Genome Sequence of the Oocydin A-
691       Producing Rhizobacterium *Serratia plymuthica* 4Rx5. *Microbiology Resource Announcements*
692       **7**, e00997-18 (2018).

693   36.  Chen, S., Blom, J. & Walker, E. D. Genomic, Physiologic, and Symbiotic Characterization of
694       *Serratia marcescens* Strains Isolated from the Mosquito *Anopheles stephensi*. *Frontiers in*
695       *Microbiology* **8**, 1483 (2017).

696   37.  Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High
697       throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat.*
698       *Commun.* **9**, (2018).

699   38.  Horesh, G. *et al.* Different evolutionary trends form the twilight zone of the bacterial pan-
700       genome. *Microbial Genomics* **7**, 000670 (2021).

701   39.  Karp, P. D. *et al.* Pathway Tools version 23.0 update: software for pathway/genome informatics
702       and systems biology. *Briefings in Bioinformatics* **22**, 109 (2021).

703   40.  Foerstner, K. U., von Mering, C., Hooper, S. D. & Bork, P. Environments shape the nucleotide
704       composition of genomes. *EMBO Reports* **6**, 1208–1213 (2005).

705   41.  Palidwor, G. A., Perkins, T. J. & Xia, X. A General Model of Codon Bias Due to GC Mutational
706       Bias. *PLOS ONE* **5**, e13431 (2010).

707   42.  Reuter, S. *et al.* Parallel independent evolution of pathogenicity within the genus *Yersinia*.
708       *Proceedings of the National Academy of Sciences of the United States of America* **111**, 6768–
709       6773 (2014).

710   43.  Ondov, B. D. *et al.* Mash: Fast genome and metagenome distance estimation using MinHash.
711       *Genome Biology* **17**, 132 (2016).

712   44.  Harris, A. K. P. *et al.* The *Serratia* gene cluster encoding biosynthesis of the red antibiotic,
713       prodigiosin, shows species- and strain-dependent genome context variation. *Microbiology* **150**,
714       3547–3560 (2004).

715   45.  Kwak, Y., Khan, A. R. & Shin, J.-H. Genome sequence of *Serratia nematodiphila* DSM 21420T,
716       a symbiotic bacterium from entomopathogenic nematode. *Journal of Biotechnology* **193**, 1–2
717       (2015).

718   46.  Matilla, M. A., Udaondo, Z., Krell, T. & Salmond, G. P. C. Genome Sequence of *Serratia*
719       *marcescens* MSU97, a Plant-Associated Bacterium That Makes Multiple Antibiotics. *Genome*
720       *Announcements* **5**, (2017).

721   47.  Cristina, M. L., Sartini, M. & Spagnolo, A. M. *Serratia marcescens* infections in neonatal
722       intensive care units (NICUs). *International Journal of Environmental Research and Public*
723       *Health* **16**, (2019).

724   48.  Daoudi, A., Benaoui, F., el Idrissi Slitine, N., Soraa, N. & Rabou Maoulainine, F. M. An
725       Outbreak of Serratia marcescens in a Moroccan Neonatal Intensive Care Unit . *Advances in*
726       *Medicine* **2018**, 1–4 (2018).

727   49.  Moles, L. *et al. Serratia marcescens* colonization in preterm neonates during their neonatal
728       intensive care unit stay. *Antimicrobial Resistance and Infection Control* **8**, 135 (2019).

729   50.  Martineau, C. *et al. Serratia marcescens* outbreak in a neonatal intensive care unit: New insights
730       from next-generation sequencing applications. *Journal of Clinical Microbiology* **56**, (2018).

731   51.  Escribano, E. *et al.* Influence of a *Serratia marcescens* outbreak on the gut microbiota
732       establishment process in low-weight preterm neonates. *PLOS ONE* **14**, e0216581 (2019).

52.  Montagnani, C. *et al.* *Serratia marcescens* outbreak in a neonatal intensive care unit: Crucial role of implementing hand hygiene among external consultants. *BMC Infectious Diseases* **15**, 11 (2015).

53.  Hurst, M. R. H., Becher, S. A. & O'Callaghan, M. Nucleotide sequence of the *Serratia entomophila* plasmid pADAP and the *Serratia proteamaculans* pU143 plasmid virulence associated region. *Plasmid* **65**, 32–41 (2011).

54.  Wood, D. E. & Salzberg, S. L. Kraken: Ultrafast metagenomic sequence classification using exact alignments. *Genome Biology* **15**, 1–12 (2014).

55.  Lu, J., Breitwieser, F. P., Thielen, P. & Salzberg, S. L. Bracken: Estimating species abundance in metagenomics data. *PeerJ Computer Science* **2017**, e104 (2017).

56.  Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).

57.  Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).

58.  Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* **29**, 1072–1075 (2013).

59.  Zerbino, D. R. & Birney, E. Velvet: Algorithms for de novo short read assembly using de Bruijn graphs. *Genome Research* **18**, 821–829 (2008).

60.  Boetzer, M., Henkel, C. v., Jansen, H. J., Butler, D. & Pirovano, W. Scaffolding pre-assembled contigs using SSPACE. *Bioinformatics* **27**, 578–579 (2011).

61.  Boetzer, M. & Pirovano, W. Toward almost closed genomes with GapFiller. *Genome Biology* **13**, 1–9 (2012).

62.  Bankevich, A. *et al.* SPAdes: A new genome assembly algorithm and its applications to single-cell sequencing. *Journal of Computational Biology* **19**, 455–477 (2012).

63.  Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial genome assemblies from short and long sequencing reads. *PLOS Computational Biology* **13**, e1005595 (2017).

64.  Koren, S. *et al.* Canu: Scalable and accurate long-read assembly via adaptive κ-mer weighting and repeat separation. *Genome Research* **27**, 722–736 (2017).

65.  Hunt, M. *et al.* Circlator: automated circularization of genome assemblies using long sequencing reads. *Genome Biology* **16**, 294 (2015).

66.  Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).

67.  Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo pipeline. *Genome Biology* **21**, 1–21 (2020).

68.  Dixon, P. VEGAN, a package of R functions for community ecology. *Journal of Vegetation Science* **14**, 927–930 (2003).

69.  Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**, 3691–3693 (2015).

70.  Page, A. J. *et al.* SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments. *Microbial Genomics* **2**, (2016).

71.  Nguyen, L. T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular Biology and Evolution* **32**, 268 (2015).

72.  Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2: Improving the Ultrafast Bootstrap Approximation. *Molecular Biology and Evolution* **35**, 518–522 (2018).

73.  Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermiin, L. S. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature Methods* **14**, 587–589 (2017).

74.  Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nature Communications 2018 9:1* **9**, 1–8 (2018).

75.  Shannon, P. *et al.* Cytoscape: A software Environment for integrated models of biomolecular interaction networks. *Genome Research* **13**, 2498–2504 (2003).

76. Tonkin-Hill, G., Lees, J. A., Bentley, S. D., Frost, S. D. W. & Corander, J. Fast hierarchical Bayesian analysis of population structure. *Nucleic Acids Research* **47**, 5539 (2019).

77. Belcour, A. *et al.* Inferring Biochemical Reactions and Metabolite Structures to Understand Metabolic Pathway Drift. *iScience* **23**, (2020).

78. Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Research* **47**, D309–D314 (2019).

79. Robertson, J., Bessonov, K., Schonfeld, J. & Nash, J. H. E. Universal whole-sequence-based plasmid typing and its utility to prediction of host range and epidemiological surveillance. *Microbial Genomics* **6**, 1–12 (2020).

80. Wickham, Hadley. *Ggplot2 : elegant graphics for data analysis*. (Springer, 2009).

81. Morris, J. H. *et al.* ClusterMaker: A multi-algorithm clustering plugin for Cytoscape. *BMC Bioinformatics* **12**, 1–14 (2011).

82. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).

83. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Computational Biology* **7**, e1002195 (2011).

84. Camacho, C. *et al.* BLAST+: Architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).

85. Guy, L., Kultima, J. R., Andersson, S. G. E. & Quackenbush, J. GenoPlotR: comparative gene and genome visualization in R. *Bioinformatics* **27**, 2334–2335 (2011).

86. Arndt, D. *et al.* PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic Acids Research* **44**, W16 (2016).

87. Wickham, H. *et al.* Welcome to the Tidyverse. *Journal of Open Source Software* **4**, 1686 (2019).

88. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T.-Y. ggtree : an r package for visualization and annotation of phylogenetic trees with their covariates and other associated data. *Methods in Ecology and Evolution* **8**, 28–36 (2017).

89. Conway, J. R., Lex, A. & Gehlenborg, N. UpSetR: an R package for the visualization of intersecting sets and their properties. *Bioinformatics* **33**, 2938–2940 (2017).

90. Ward, D. V. *et al.* Metagenomic Sequencing with Strain-Level Resolution Implicates Uropathogenic *E. coli* in Necrotizing Enterocolitis and Mortality in Preterm Infants. *Cell Reports* **14**, 2912-24 (2016).

91. Roach, D. J. *et al.* A Year of Infection in the Intensive Care Unit: Prospective Whole Genome Sequencing of Bacterial Clinical Isolates Reveals Cryptic Transmissions and Novel Microbiota. *PLoS Genetics* **11**, e1005413 (2015).

92. Grimont, F. & Grimont, P. A. D. Genus XXXIV, *Serratia*. In *Bergey's Manual of Systematic Bacteriology, Volume 2 Part B* (eds. George Garrity, Don Brenner, Nole Kreig & James Staley) 799-810 (Springer, 2005).

**Acknowledgements**

24

842

**Author Contributions**

844    D.J.W., N.R.T. and S.J.C. conceived the study; D.J.W. performed the bioinformatics analyses, with

845    contributions from A.C.L. and D.C.L; P.A.D.G., F.G. and E.A. performed identification and

846    biochemical characterisation of *Serratia* isolates in the Institut Pasteur collection; D.J.W., K.P., E.N.

847    and F.X.W. contributed to isolate resuscitation and sequencing; D.J.W., A.J.C., E.H., M.T.G.H., N.R.T.

848    and S.J.C analysed and interpreted results; D.J.W., N.R.T. and S.J.C. wrote the paper with input from

849    the other authors.

850

**Competing Financial Interests**

852    The authors declare no competing financial interests.

**Figure 1: Phylogeny of the genus *Serratia*.** Maximum-likelihood phylogenetic tree constructed from polymorphic sites of a core-gene alignment comprised of 2252 genes from 664 *Serratia* genomes, comprising

408 genomes from publicly available databases, and 256 sequenced in this study. Tree constructed with 1000 ultrafast bootstraps. The core-gene alignment was produced from a Panaroo pan-genome analysis run with "--clean_mode moderate" and the protein family threshold set to 70% shared sequence identity. Branches are coloured according to phylogroups defined by clustering assemblies at 95% ANI. Clades are shaded according to lineage, calculated through hierarchical bayesian clustering to three levels using FastBaps. "Labelled species" refers to the labelled name of species on the provided *Serratia* strain sample, or species name associated with published *Serratia* genome sequences in the NCBI GenBank database.
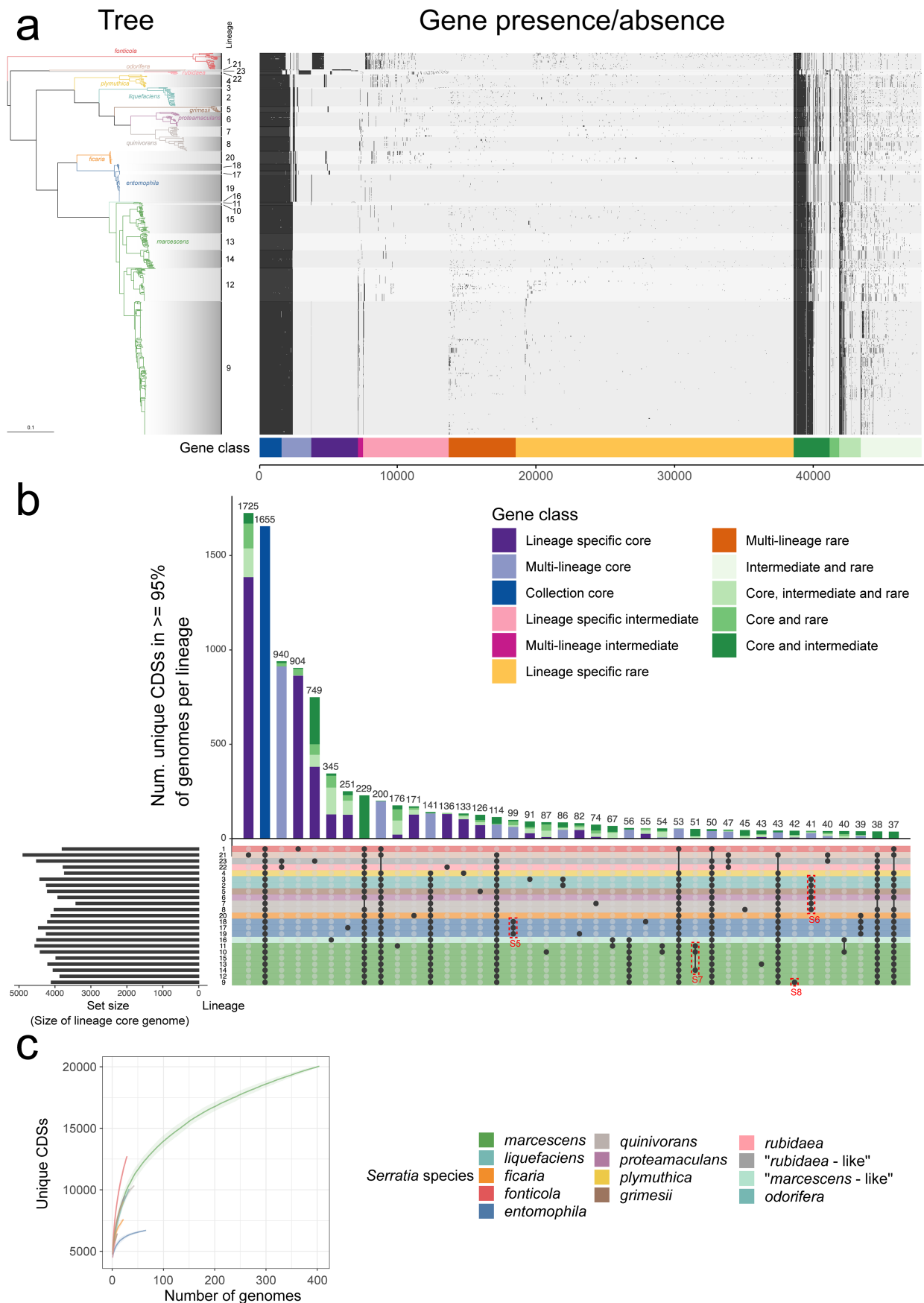
**Figure 2: The pan-genome of *Serratia*.** (a) Presence/absence matrix of the 46,588 genes in the *Serratia* pangenome, generated using Panaroo and overlaid with shading according to lineage, alongside the maximum-likelihood tree based on the core-gene alignment shown in Fig. 1. The presence/absence matrix is ordered by gene class as defined by Twilight. (b) UpSetR plot showing the 50 largest intersections of lineage-specific core

genomes (genes present in ≥ 95% of strains in each lineage). Lineages with membership to each intersection are shown by the presence of a black dot in the presence/absence matrix underneath the stacked bar plot. Stacked bar plots representing the number of genes in each intersection are coloured according to the gene classes assigned by Twilight, where singleton lineages have been included (in this case, lineages 22 and 23 are singletons). Rows in the presence/absence matrix correspond to each lineage and are coloured according to *Serratia* species defined by fastANI. Dashed red boxes indicate intersections of genes represented in Supplementary Figures S5-S8. (c) Estimated pan-genome accumulation curves for each *Serratia* phylogroup. Shaded region represents standard deviation.
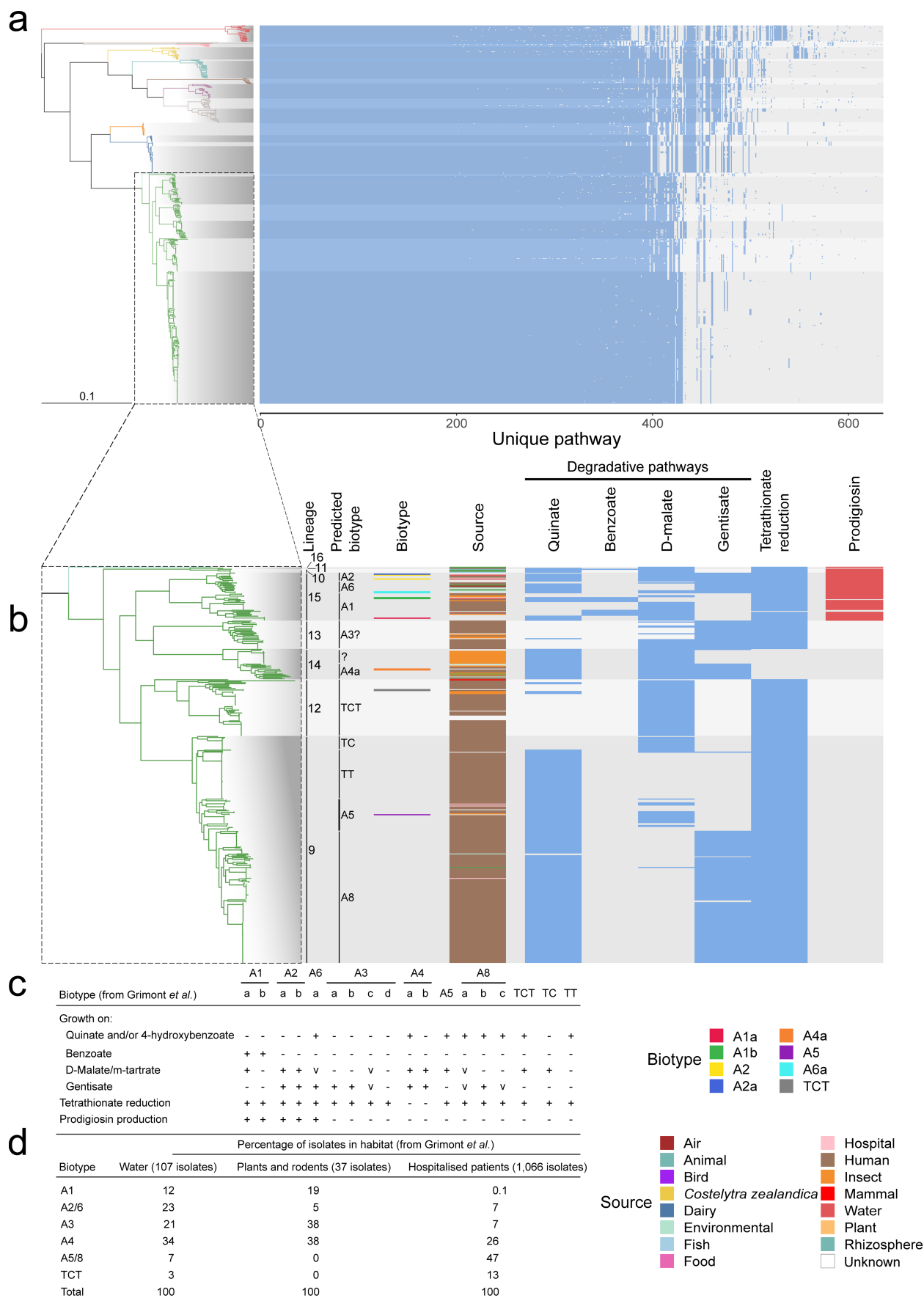
**Figure 3: Predicted metabolic pathways in *Serratia* and correspondence with historical biotyping.** (a) Predicted metabolic pathways across *Serratia*, predicted using Pathway Tools following re-annotation of assemblies using Interproscan/EggNOG-based functional annotation of representative sequences of protein groups defined by Panaroo. (b) Presence/absence of selected complete metabolic pathways across *Serratia*

30

*marcescens*. Pathways were selected according to a subset of the biochemical tests originally used to group *Serratia* isolates into Biotypes (c). (d) Table of habitat source for different *S. marcescens* biotypes, reproduced from Grimont *et al*., 2006.
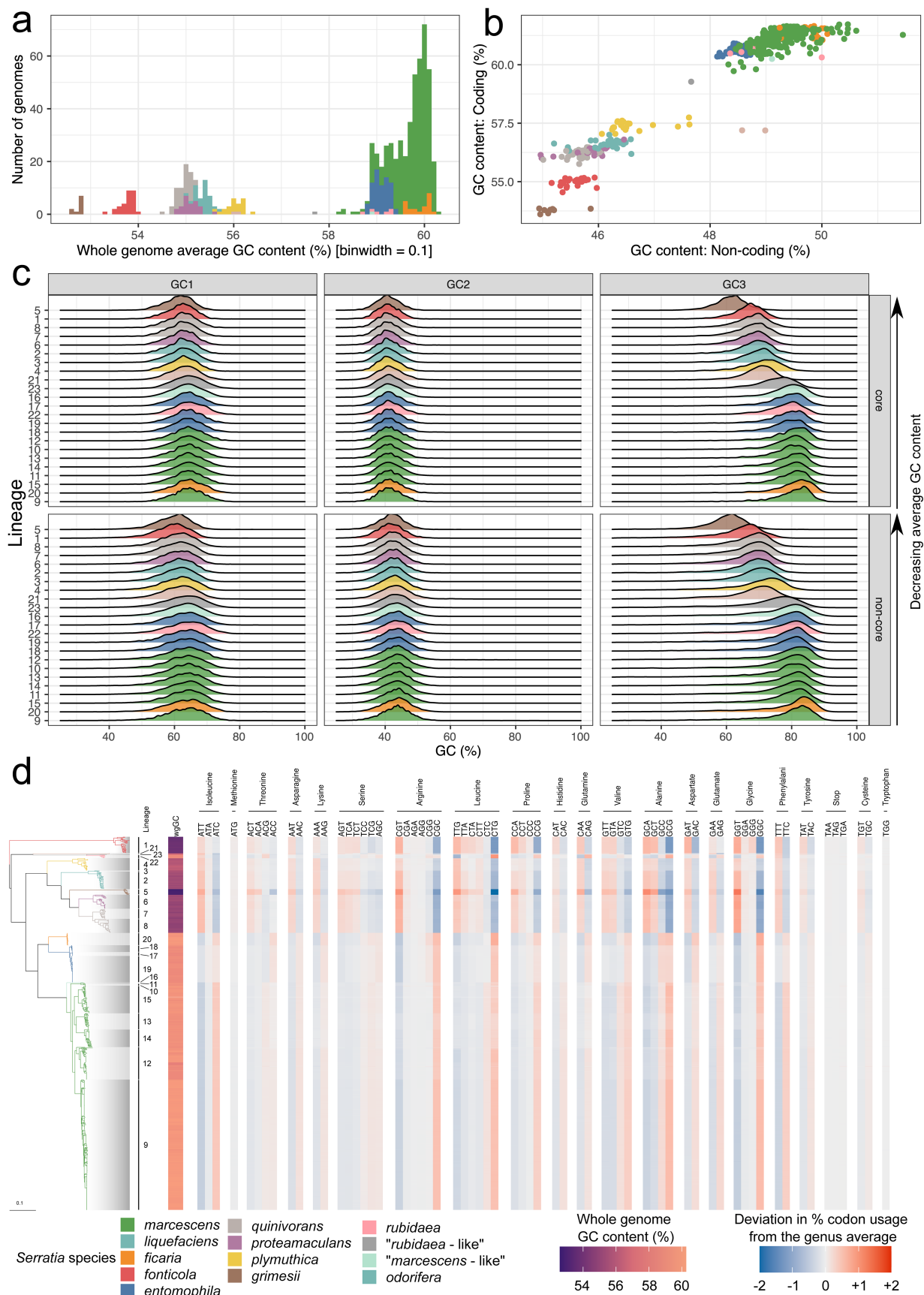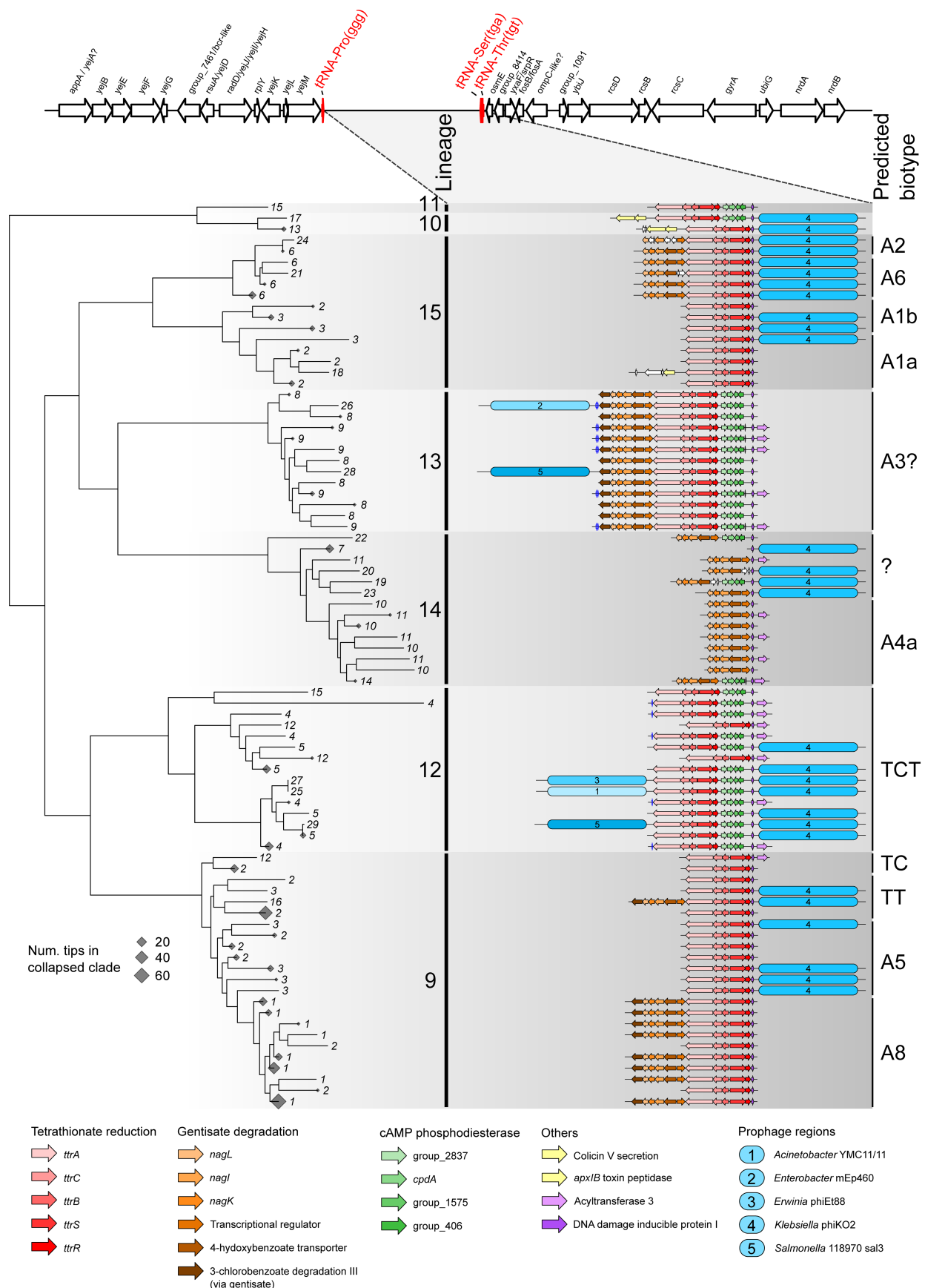
**Figure 4:** *Serratia* **is split by GC content.** (a) Histogram of GC content (average over whole genome) across *Serratia*. (b) Comparison of GC content in all coding regions vs. that of all intergenic regions. (c) Distribution of GC content in codon positions 1, 2 and 3 in all genus-core (core) and non-genus-core (non-core) genes across each lineage. Data is normalised according to gene length. Ridgeplots are coloured according to *Serratia* species/phylogroup. Lineages are ordered from top to bottom according to average GC content across the

whole genome. (d) Codon usage (CU) within the genus-core genome. Blue to red colour represents deviation from the average CU across the entire genus for each codon, with this genus-average CU calculated from a per-lineage mean CU value to account for the different numbers of sequences in each lineage. The whole genome GC content is also shown in the left-most column.

**Figure 5: A tRNA-associated hypervariable region ('plasticity zone') is a hot-spot for horizontal transfer of gene cassettes for metabolic pathways used for biotyping within *S. marcescens*.** The gene arrangement between the conserved tRNA-Pro$_{ggg}$ and tRNA-Ser$_{tga}$ in *S. marcescens* is plotted against a maximum-likelihood sub-phylogeny from the tree shown in Fig. 1. Clades for which all descending tips represent strains that have an identical set of genes in the locus depicted are collapsed, and denoted by a diamond shape within the tree.

The size of the diamond represents the number of tips in each collapsed clade. Tips lacking a completely assembled gene locus between tRNA-Pro$_{ggg}$ and tRNA-Ser$_{tga}$ have been pruned from the tree. Genes are coloured according to their role, or in the absence of any predicted function, named according to the group number assigned by Panaroo in the pan-genome (Fig. 2). Prophage regions and the closest related prophage sequence determined by PHASTER are indicated.
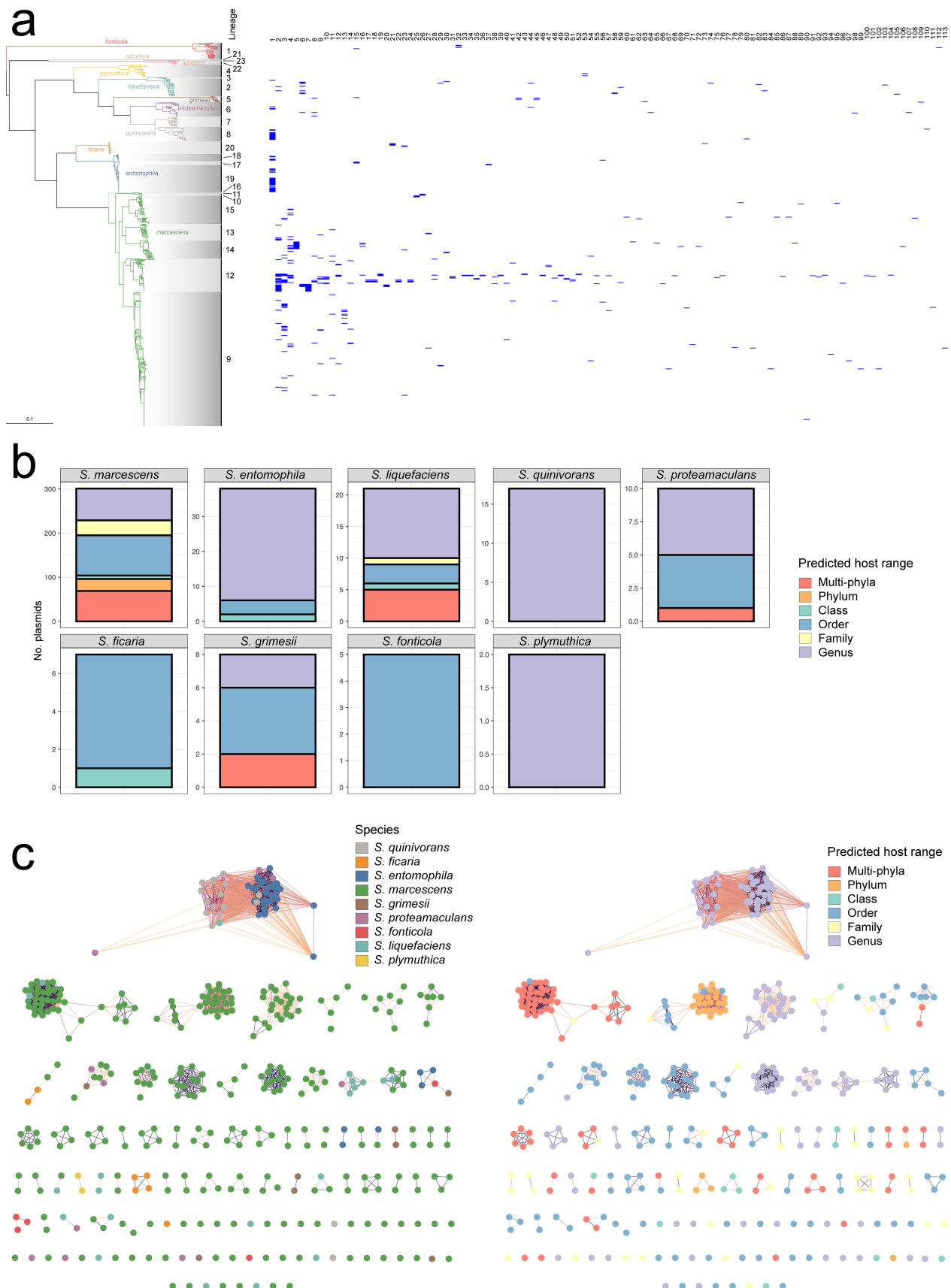
**Figure 6: Predicted plasmids across *Serratia*.** (a) Distribution of the 131 plasmid clusters identified against the phylogeny of *Serratia* shown in Fig. 1. (b) Number and predicted host range of plasmids identified in *Serratia* genomes. (c) Diversity of *Serratia* plasmids according to species (left) and predicted host range (right). Within each panel, the order of clusters (from left-right in descending rows) is the same order as presented in the heatmap in panel a (left-right).
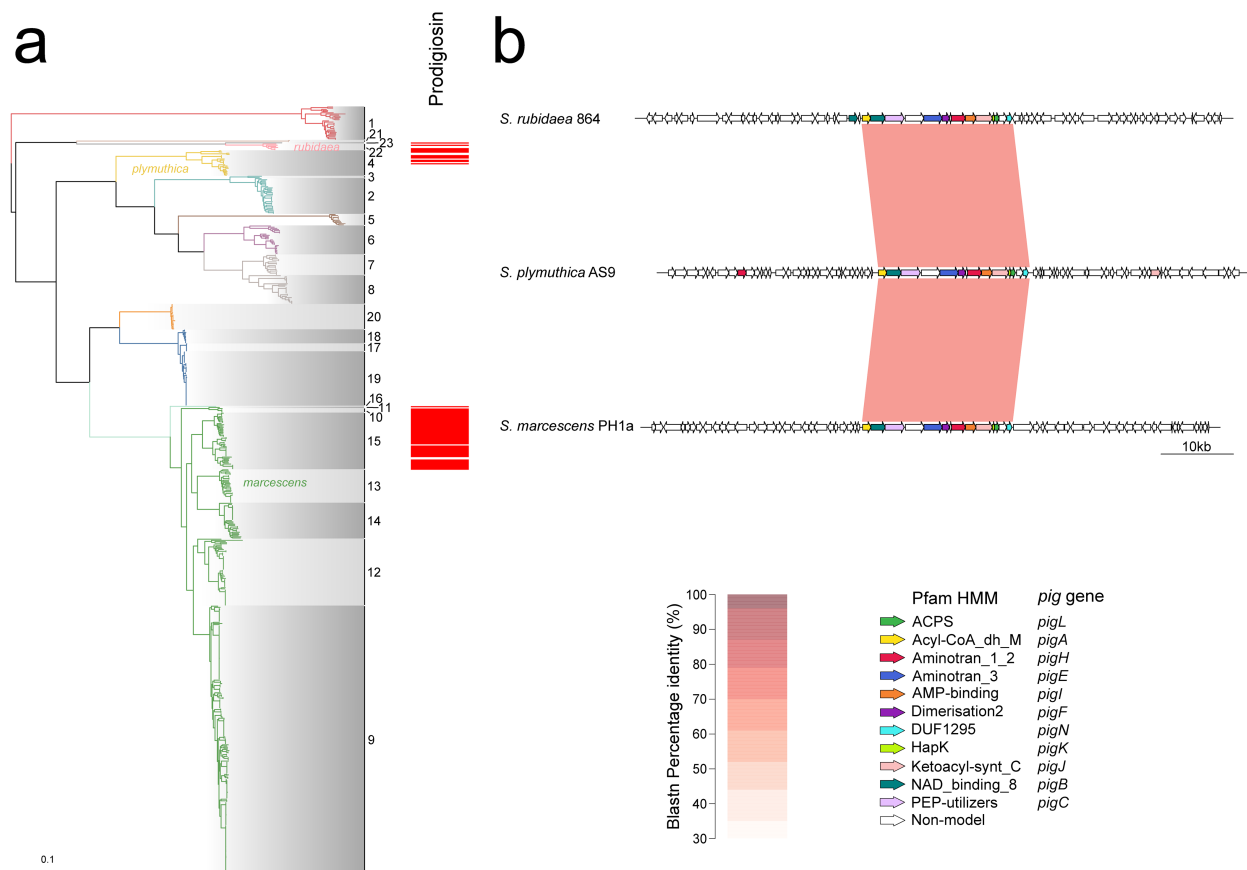
**Figure 7: The prodigiosin gene cluster is present variably across *Serratia* and in different genomic loci.** (a) Prodigiosin (*pig*) gene clusters identified using Hamburger are plotted against the maximum-likelihood phylogeny of *Serratia* shown in Fig. 1. (b) Pairwise blastn comparison of *pig* loci (core *pig* cluster +/- 30 kb) from representative members of the three species containing *pig* genes, extracted using Hamburger.