

# The genus *Serratia* revisited by genomics

Received: 1 January 2022

Accepted: 24 August 2022

Published online: 03 September 2022

 Check for updates

David J. Williams <sup>1,2</sup>, Patrick A. D. Grimont <sup>3</sup>, Adrián Cazares <sup>2,4</sup>, Francine Grimont<sup>3</sup>, Elisabeth Ageron<sup>3,9</sup>, Kerry A. Pettigrew <sup>5</sup>, Daniel Cazares<sup>2</sup>, Elisabeth Njamkepo<sup>6</sup>, François-Xavier Weill <sup>6</sup>, Eva Heinz <sup>2,7</sup>, Matthew T. G. Holden <sup>5</sup>, Nicholas R. Thomson <sup>2,8</sup>  & Sarah J. Coulthurst <sup>1</sup> 


The genus *Serratia* has been studied for over a century and includes clinically-important and diverse environmental members. Despite this, there is a paucity of genomic information across the genus and a robust whole genome-based phylogenetic framework is lacking. Here, we have assembled and analysed a representative set of 664 genomes from across the genus, including 215 historic isolates originally used in defining the genus. Phylogenomic analysis of the genus reveals a clearly-defined population structure which displays deep divisions and aligns with ecological niche, as well as striking congruence between historical biochemical phenotyping data and contemporary genomics data. We highlight the genomic, phenotypic and plasmid diversity of *Serratia*, and provide evidence of different patterns of gene flow across the genus. Our work provides a framework for understanding the emergence of clinical and other lineages of *Serratia*.

The genus *Serratia* was originally described in Italy in the early 19th century, following the observation of a blood-like red discoloration appearing on polenta from organic growth<sup>1</sup>. It has since become clear that *Serratia* species are ubiquitous, free-living, motile Gram-negative proteobacteria, traditionally considered members of the *Enterobacteriaceae*. The genus *Serratia* represents a broad and diverse genus of more than ten species, delineated by DNA-DNA hybridisation and characterised by extensive physiological and biochemical tests<sup>2–12</sup>. Despite being a diverse genus, much of the contemporary research and understanding of *Serratia* has focused on the type species, *Serratia marcescens*. *S. marcescens* has served as a model system for studying key bacterial traits, including protein secretion systems<sup>13</sup> and motility<sup>14</sup>, but it also represents an important opportunistic human pathogen<sup>15–17</sup> for which there has been a dramatic rise in the incidence of multi-drug resistance and reported cases of problematic nosocomial infections<sup>18</sup>.

Other members of this genus include *S. rubidaea* and *S. liquefaciens*, which have also been reported to cause hospital-acquired infections, albeit less frequently<sup>17,19–21</sup>. In addition to infection of human

hosts, members of multiple *Serratia* species represent insect pathogens or are otherwise associated with insects. *Serratia entomophila* has been used as a biocontrol agent in New Zealand to predate upon the pasture pest, *Costelytra zealandica*<sup>12,22–25</sup>, and *S. proteamaculans*<sup>26</sup> and *S. marcescens*<sup>13,27–29</sup> have also been shown to be insect pathogens. In contrast, *S. ficaria* is associated with the pollination and oviposition cycle between figs and fig wasps, respectively<sup>6</sup>. In addition, underlining the ubiquitous nature of this genus, *Serratia* species can be found in a multitude of environmental niches<sup>7,9,29–36</sup>, including frequent isolation from aqueous environments<sup>17,21</sup>.

Given its importance to human health, it is perhaps unsurprising that the majority of genomic information available for *Serratia* originate from clinically-derived *S. marcescens*. Recently a number of *S. marcescens* sequences have also been included within large-scale metagenomic studies from pre-term neonates<sup>37</sup> or nosocomial environments<sup>38</sup>. However, for these, as for all the sequences from clinically isolated strains, there is a critical lack of a robust phylogenetic framework for the *Serratia* genus within which the *S. marcescens* sequences can be placed. The historic collection of Patrick Grimont,

<sup>1</sup>Division of Molecular Microbiology, School of Life Sciences, University of Dundee, Dundee, UK. <sup>2</sup>Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK. <sup>3</sup>Unité Biodiversité des Bactéries Pathogènes Emergentes, INSERM Unité 389, Institut Pasteur, Paris, France. <sup>4</sup>European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, UK. <sup>5</sup>School of Medicine, University of St Andrews, St Andrews, UK. <sup>6</sup>Institut Pasteur, Université de Paris, Unité des Bactéries Pathogènes Entériques, Paris, France. <sup>7</sup>Departments of Vector Biology and Clinical Sciences, Liverpool School of Tropical Medicine, Liverpool, UK. <sup>8</sup>Department of Infection Biology, London School of Hygiene and Tropical Medicine, London, UK. <sup>9</sup>Present address: Université Paris Cité, INSERM UMR-S1151, CNRS UMR-S8253, Institut Necker Enfants Malades, Paris, France.  e-mail: [nrt@sanger.ac.uk](mailto:nrt@sanger.ac.uk); [s.j.coulthurst@dundee.ac.uk](mailto:s.j.coulthurst@dundee.ac.uk)

located at the Institut Pasteur, includes the original strains used to define biochemically and phenotypically the vast majority of the known *Serratia* species<sup>3,4,6,9–12</sup>. These strains, kept in cold storage for between 20 and 40 years, represented a unique resource, providing the opportunity to compare historical biochemical, phenotypic and DNA-DNA hybridisation data with contemporary genomics data, across the entire *Serratia* genus. In an attempt to understand the many facets and functions of the different species within the *Serratia* genus, and, importantly, to understand the context within which *S. marcescens* is becoming more widespread as a problematic opportunistic pathogen, we aimed to assemble and analyse a balanced genomic dataset that reflected the entire *Serratia* genus.

In this work, we supplement existing publicly available, published *Serratia* genome sequences by sequencing over 200 *Serratia* isolates that are non-clinical in origin, and mainly belonging to species other than *S. marcescens*, from the historic Grimont collection. This allows us to bring together previous molecular and biochemical knowledge from this important genus and place it in a contemporary genomics framework. Our findings not only explain why previous definitions of the genus *Serratia* were robust, or not, they also highlight important differences in diversity, plasticity and niche adaptation of the species within it.

## Results

### Deep divisions demarcate phylogroups within *Serratia*

Here we sequenced and analysed a collection of 256 novel *Serratia* genome sequences and combined these data with 408 published genomes. Our total set of 664 genome sequences included those of 215 isolates from the original genus-defining Grimont collection sequenced here, 205 multidrug-resistant *S. marcescens* isolates from UK hospitals isolated between 2001 and 2011<sup>18</sup>, and an additional 41, more diverse, *Serratia* isolates from UK hospitals, sequenced here for comparison with the latter collection<sup>18</sup>.

We inferred the genus phylogeny from the whole genome data using a core-gene alignment-based approach (Fig. 1, Supplementary Fig. 1). It is evident from Fig. 1 that there are deep divisions within this phylogeny that correlate with both the current genus taxonomy and with species-level grouping calculated using genome-wide average nucleotide identity (ANI; clustered using a cutoff of 95 percent; Fig. 1, Supplementary Fig. 2). To gain a finer-scaled view, we used hierarchical Bayesian clustering (FastBaps) to four levels in order to further subdivide the phylogeny and the species-level groups. In total, we identified 7, 16, 23 and 27 clusters across the four levels, respectively (Supplementary Fig. 3). FastBaps level 1 clusters comprise monophyletic clades reflecting individual or multiple ANI groupings within the genus, consistent with speciation or species complexes<sup>39</sup>. The second and third levels reveal the presence of several subdivisions within some of the species-level phylogroups (Fig. 1, Supplementary Fig. 3). Hereafter, we refer to the clusters set out by FastBaps level 3 as Lineages 1–23 (L1–23; Fig. 1).

Interestingly, within the tree, there are two examples of singleton genomes occupying both a single ANI species-level phylogroup and representing a discrete FastBaps lineage (L16 and L23). Although both are situated within well-defined species, these two singletons are hereafter referred to as ‘*S. marcescens*-like’ and ‘*S. rubidaea*-like’, for L16 and L23, respectively. Our phylogeny also resolves previous taxonomic discrepancies. Here, based on the core-gene phylogeny, the *liquefaciens* complex, historically made up of *S. liquefaciens*, *S. grimesii*, *S. proteamaculans* and *S. proteamaculans* subs. *quinovora* is monophyletic (Fig. 1). Previous work had suggested *S. proteamaculans* be resolved into both *S. proteamaculans sensu stricto* and a sub-species, termed *S. proteamaculans* subs. *quinovora*<sup>10,21</sup>. However, a species-level distinction between these two taxa, rather than a

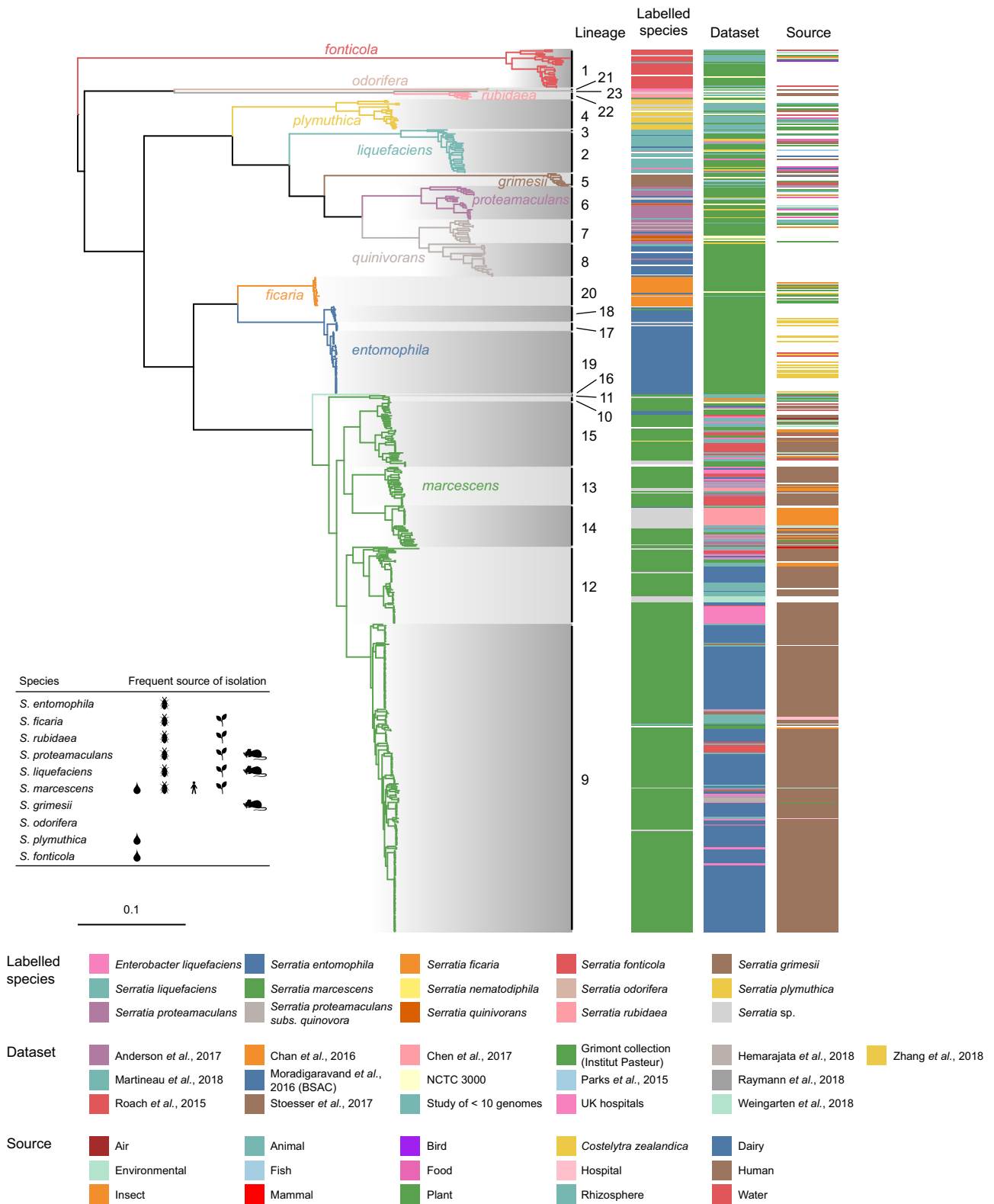
sub-species one, was subsequently proposed<sup>30</sup>. We observe that genomes labelled as *S. proteamaculans* and *S. proteamaculans* subs. *quinovora* form two separate ANI phylogroups in a monophyletic branch made up of L6–8 (Fig. 1, Supplementary Fig. 2). This supports the presence of two distinct species-level groups, which we refer to as *S. proteamaculans* and *S. quinivorans* in accordance with the latter work<sup>30</sup>, thereby extending the membership of the *liquefaciens* complex to four species, including *S. quinivorans*. Furthermore, only a single genome links L7 and L8 into a single ANI phylogroup within *S. quinivorans* (Supplementary Fig. 4), which may suggest that a further sub-species separation within *S. quinivorans* is appropriate.

### Concordance between historical biochemical phenotyping and metabolomic predictions

The current taxonomic structure of the *Serratia* genus, summarised by Grimont and Grimont, 2005<sup>40</sup>, is based on 41 phenotypic and biochemical tests used to differentiate between different *Serratia* species or species-complexes. Many of the *Serratia* isolates originally used to define the genus taxonomy were sequenced here (Fig. 1; Supplementary Data 1), presenting a unique opportunity to reconcile this historical biochemical metadata with genomic predictions. First, we calculated the genus pan-genome using a population structure-aware approach<sup>41</sup>. The pan-genome comprised 47,743 discrete gene groups (Fig. 2a), of which 2252 were present in at least 99 percent of all genomes in the dataset (which would be defined as a traditional core genome), however, only 1655 of these were present in at least 95 percent of genomes within each FastBaps lineage, and therefore are core to all lineages (Fig. 2b). These 1655 genes are hereafter defined as the genus-core. From the 47,743 genes of the pan-genome, we predicted the metabolic potential of the genus. We identified 641 different complete metabolic pathways using Pathway tools<sup>42</sup> (Fig. 3a), of which 260 were core to the genus, being present in all known lineages (L1–23) (Fig. 3a, Supplementary Fig. 9).

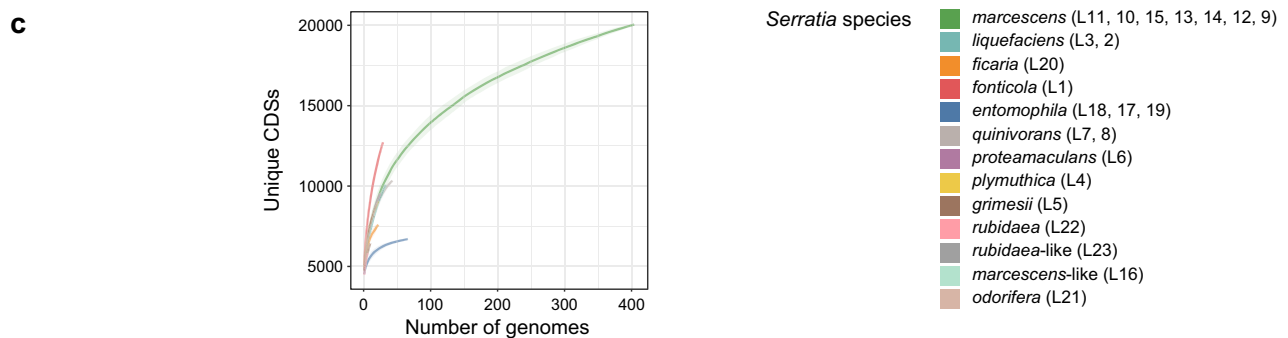
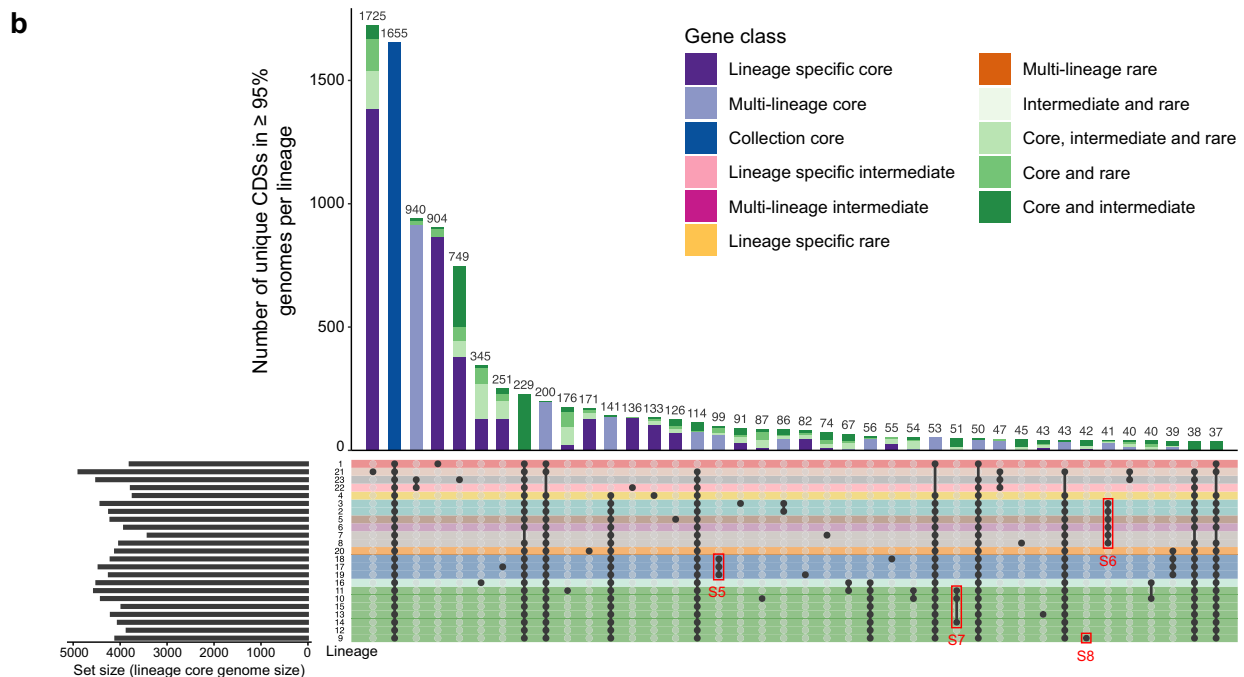
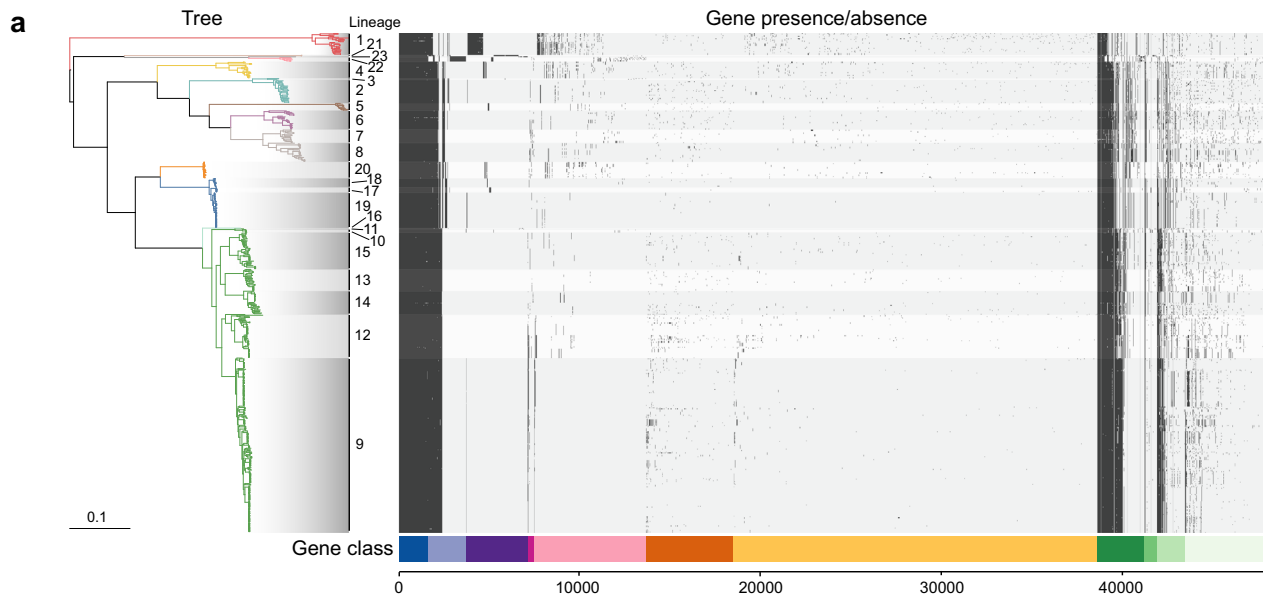
Of the 41 metabolic tests used to define species within the genus, some were also used to define biotypes within a species. It can be seen that the fine-scale delimitation of biotypes within the *liquefaciens* complex generally agrees well with the phylogenetic structure we observe (Supplementary Fig. 10). Similarly, ten of the 41 metabolic tests were used previously to split *marcescens* into ten biotypes, which reflected differences in niche occupancy<sup>21</sup>. It is clear that these biotypes are also robust markers of phylogenetic subdivisions within this important species (Fig. 3c).

In silico pathway predictions were used to identify the genes and/or pathways linked to the biotype tests for *S. marcescens*, where strain biotype metadata was available. Four of these ten tests (growth on *m*-erythritol, trigonelline, 3-hydroxybenzoate and lactose) were not investigated because there were no corresponding pathway assignments in the predicted metabolic network. Next, the phylogenetic distribution of these metabolic genes was plotted across the phylogeny and extrapolated back to the most basal internal node differentiating two biotypes, based on gene/pathway presence or absence (Fig. 3b). Where we had no genome representative of a particular biotype in which to identify the genes for the cognate pathway, these were predicted using the results of the in silico metabolic prediction for known pathways (Fig. 3c). Although there are some discrepancies between pathway presence/absence and historical phenotypes for different biotypes, Fig. 3b shows a robust linkage between inferred phylogeny and biotype. Across the species, these data show that L15 corresponds with the pigmented biogroups A2, A6 and A1, whilst L14, L12 and L9 correspond with the nonpigmented biogroup A4 (L14), biotype TCT (L12), and biotypes A5, TT, TC, biogroup A8 (L9). Furthermore, adding the source of sample isolation shows that niche occupancy also aligns well with the population structure and biotype



**Fig. 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. The inset illustrates the frequent sources of isolation of each species<sup>21</sup>, with symbols representing water, insects, human, plants and small mammals. The phylogenetic tree is reproduced in Supplementary Fig. 1 with the addition of the outgroup root, the country of isolation for each strain, and bootstrap values.



data (Fig. 3c). In particular, strains representing original biotypes described in the 1980s that were associated with hospitalised patients (TCT, TC, TT, A5, A8) are situated in the same phylogenetic position as contemporary clinically-isolated *S. marcescens*, implying important adaptations that can be linked to risk of disease or greater fitness in hospital environments.

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

Changes in GC content of coding sequences over time have been hypothesised to reflect subtle differences in mutational bias as a consequence of long-term niche adaptation or different lifestyles<sup>43</sup>. Given that there are clear differences in the lifestyles and niches

**Fig. 2 | The pan-genome of *Serratia*.** **a** Presence/absence matrix of the 47,743 genes in the *Serratia* pan-genome, generated using Panaroo and overlaid with shading according to lineage, alongside the maximum-likelihood tree in Fig. 1. The presence/absence matrix is ordered by gene class as defined by Twilight. Gene class is first defined within each lineage by calculating whether genes are core (in  $\geq 95\%$  of strains in each lineage), intermediate (in  $>15\%$  and  $\leq 95\%$  of strains), or rare (in  $\leq 15\%$  of strains). Classification of each gene group per lineage is then compared between lineages. Gene groups core to all lineages are collection core, gene groups core to only certain lineages are multi-lineage core, and genes core to only a single lineage are lineage-specific core. Individual genes found at intermediate or rare occurrence in all, multiple, or single lineages are classified similarly, as intermediate or rare genes. These three classes are indicated by colour: core, blue shades; intermediate, pink shades; and rare, orange shades. Genes which are in one classification (core, intermediate, rare) in a particular lineage but in another classification in a separate

lineage are termed hybrid classes (green shades). **b** UpSetR plot showing the 40 largest intersections of lineage-specific core genomes (genes present in  $\geq 95\%$  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 (here L22 and L23) have been included. Rows in the presence/absence matrix correspond to each lineage and are coloured according to *Serratia* species defined by fastANI. Red boxes indicate intersections of genes represented in Supplementary Figs. 5–8. **c** Estimated pan-genome accumulation curves for each *Serratia* phylogroup. Shaded region represents standard deviation. Throughout, species are coloured according to the key in **c**. Source data are provided as a Source Data file.

between *Serratia* species and intra-species lineages, we investigated the distribution of GC content across the genus. We observe that *Serratia* is broadly divided into two phylogenetically-coherent groups based on whole-genome GC content: *marcescens*, *entomophila*, *ficaria* and *rubidaea* show a GC content of  $\sim 59\%$  (59.0–59.9%), whilst *odorifera*, *fonticola*, *plymuthica*, *liquefaciens*, *proteamaculans*, *quinivorans* and *grimesii* have a GC content ranging from 52.7 to 56.1% (Fig. 4a). The singleton *S. rubidaea*-like and *S. marcescens*-like genomes have an average GC content of 57.7% and 58.9%, respectively, consistent with their positions in the tree adjacent to *rubidaea* or *marcescens*. Additionally, we observed that whilst there is a difference in the GC content of coding and non-coding regions of the same genome, the difference in GC content between species is equivalent in both coding and non-coding regions (Fig. 4b).

To understand how this GC pattern impacts protein coding, we investigated the variation in GC content over the three codon positions for all lineages using the genus-core set of 1655 genes (Fig. 4c; Fig. 2b), and separately for all other genes, designated ‘non-core’. We observed no obvious difference in GC content between genes that were core or non-core, at all three codon positions, termed as GC1-3 (Fig. 4c). The GC content at GC2 is essentially fixed across the genus (Fig. 4c), whilst GC1 shows a slight skew across the genus, varying by  $\sim 1\%$ . Codon position GC3 showed a clear increase in the frequency of A/T-ending codons in low GC species compared with high GC species, as expected<sup>44</sup> (Fig. 4d). Hence, the difference in average GC across the genus is largely explained by variation in codon position GC3. For example, GC3 in *S. grimesii* is 20% lower than in *S. marcescens* L9 (Fig. 4c).

Taken together, the variations in metabolic capability and GC content between both species and niche-adapted lineages are indicative of long-term niche adaptation within evolutionary timescales.

### Pan-genome analysis highlights lineage-specific gene gain and loss as well as intra-genus gene flow

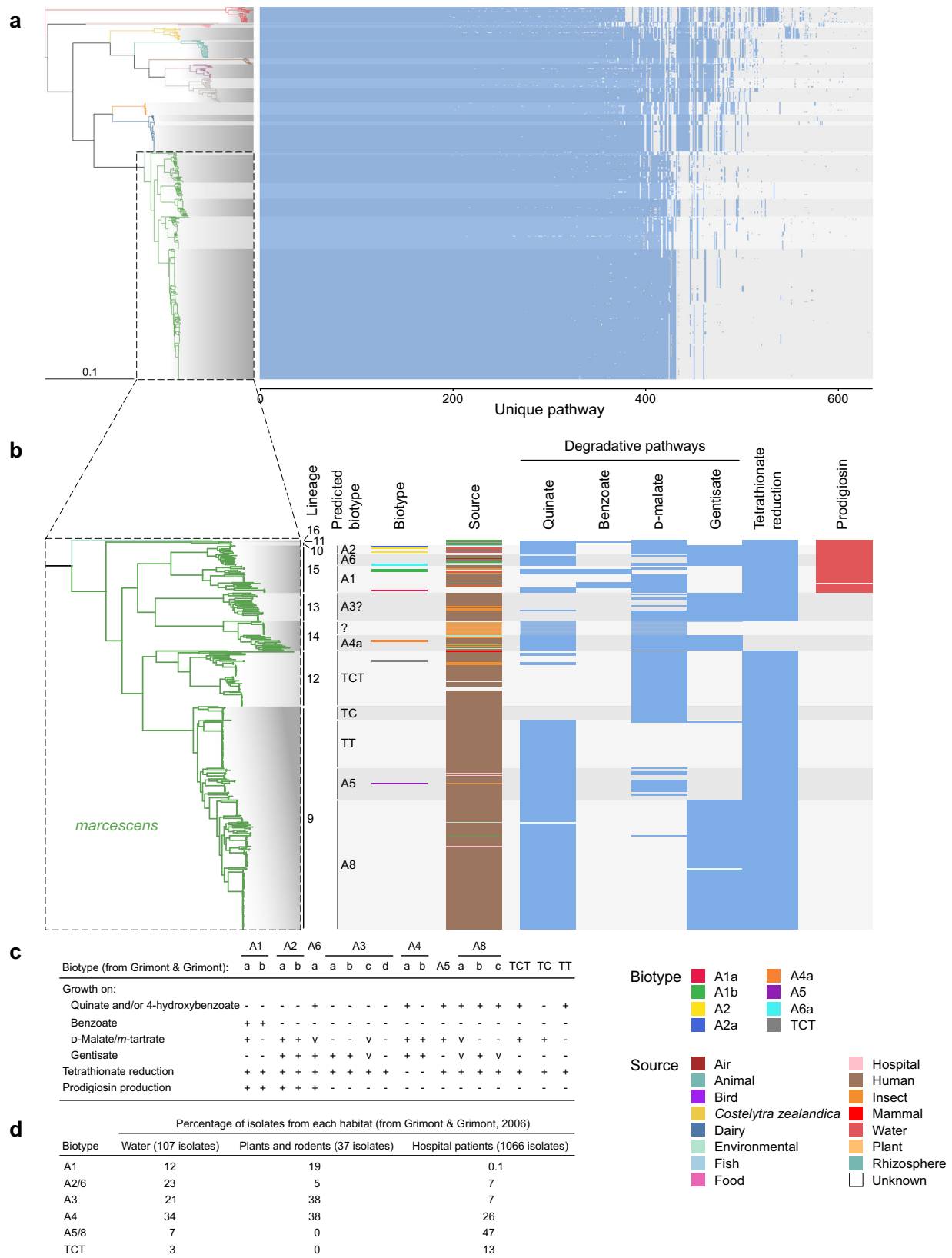
The results so far suggest that the pan-genome of *Serratia* lineages is phylogenetically constrained, yet members of *Enterobacteriaceae* are known to have a highly plastic gene content through horizontal gene transfer (HGT). To investigate this further, we sought to understand genus-wide species plasticity. 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<sup>41</sup> 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).

Variations in the pan-genome between different *Serratia* species were evident. For example, whilst *S. entomophila* and *S. ficaria* display similar core gene branch lengths, indicative of similar evolutionary timescales, *S. entomophila* has a closed pan-genome whilst *S. ficaria* has an open pan-genome (Fig. 2c). The difference in the size of the accessory genome between these two species is 726 genes, with *S. entomophila* and *S. ficaria* and having accessory genome sizes of 2764 and 3490 genes, respectively, despite *S. ficaria* being represented by fewer genomes in the analysis (Supplementary Table 1). This suggests that different *Serratia* species may have varying propensities for 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. Here we use the term ‘gene flow’ to refer to the apparent movement of genes across the genus based on their phylogenetic distribution, which may reflect both gene losses or gains, by mechanisms such as, but not limited to, deletion, horizontal gene transfer or recombination. 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).

In an attempt to understand the mechanisms by which genes are gained and lost, we focused initially on *S. marcescens*. We investigated the genetic context of the metabolic gene loci associated with different biotypes. In doing so, we identified a hypervariable locus analogous to the plasticity zone seen in *Yersinia*<sup>45</sup>. Variation in this locus explained some of the biochemical differences seen within *S. marcescens*. This plasticity zone is located between two sets of tRNAs: one encoding tRNA-Pro<sub>ggg</sub>, the other encoding tRNA-Ser<sub>tga</sub> and tRNA-Thr<sub>tgt</sub>. It



encodes the genes required for gentisate degradation (*nag* gene cassette) and/or tetrathionate reduction (*ttr* gene cassette), present in the same order and orientation across the species, located alongside three sets of genes that are variably present across the *S. marcescens* phylogeny (Fig. 5). These three sets comprise: (1) four genes including one encoding a cyclic AMP (cAMP) phosphodiesterase; (2) an

acyltransferase; and (3) a two-gene toxin cassette. A gene predicted to encode a DNA damage-inducible protein I (*dinI*) is always present, downstream of the *ttr/nag/cAMP* genes and upstream of the acyltransferase gene. In a small number of instances, frameshifts have truncated or split coding genes in this region. Additionally, prophage sequences can also be found flanking these variable sets of genes in

**Fig. 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. Shown alongside the maximum-likelihood tree in Fig. 1. **b** Presence/absence of selected complete metabolic pathways across *Serratia marcescens*. Presence of degradative pathways and that of tetrathionate reduction

are coloured in blue, whilst presence of the *pig* gene cluster containing genes required for the biosynthesis of prodigiosin is coloured in red to reflect the red colour of this pigmented molecule. Pathways were selected according to a subset of the biochemical tests originally used to group *Serratia* isolates into Biotypes<sup>21</sup> (v, variable). **d** Habitat source for different *S. marcescens* biotypes<sup>21</sup>. Tables in **c** and **d** are adapted from Grimont and Grimont (ref. 21).

some genomes (Fig. 5). Interestingly, in L13 and L9, when the *nag* genes are present, an additional gene, encoding a protein with predicted 3-chlorobenzoate degradation activity, is present 3' of the other genes in the cassette (Fig. 5).

Further evidence of gene flow can be seen in *S. marcescens* (Supplementary Figs. 7, 8). Certain genes core to *S. marcescens* L10, L11 and L14 were also found in members of other lineages, including *S. marcescens* L15 and L9, and *S. proteamaculans* L6 (Supplementary Fig. 8). On closer inspection, the genes shared with *S. marcescens* L15 and *S. proteamaculans* L6 comprise a Type VI Secretion System (T6SS). Whilst polyphyletic across *S. marcescens*, this T6SS is syntenic when found in *S. marcescens* but is encoded in a different region of the chromosome when present in *S. proteamaculans*. In both cases, this T6SS is encoded adjacent to a tRNA, and also an integrase in *S. proteamaculans*, potentially suggestive of horizontal transfer across the genus from *marcescens*. There are also 42 genes core to the clinically-associated *S. marcescens* L9, for which 37 are also found polyphyletically across the rest of *S. marcescens* (Supplementary Fig. 8). Many of these genes are predicted to be components of fimbrial usher systems (Supplementary Fig. 8).

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

To understand the potential contribution of plasmids to the plasticity seen in this genus, we searched for plasmid contigs in our genus-wide dataset. This uncovered 409 putative plasmids in 228 genomes and 9 species, 301 (73%) of them present in *S. marcescens* (Fig. 6; Supplementary Data 2). The collection of identified plasmids displays a wide range of sizes (-1–310 kb) and GC content (-30–66%), indicating diversity. However, the distribution of these traits varied amongst *Serratia* species (Supplementary Figs. 13–15). For example, plasmids identified in *S. marcescens* and *liquefaciens* show a markedly broader range of size and GC content compared with those detected in *S. entomophila* and *quinivorans*. Seventy out of a total of 113 predicted plasmid replicons were found within *S. marcescens* L9 and L12, which are the 'clinical' lineages in which 97% and 81% percent of the isolates, respectively, are known to be human- or clinically-associated. In terms of mobility, 296 (72%) of the plasmids were predicted to be conjugative or mobilisable (Fig. 6; Supplementary Data 2), highlighting their potential role in HGT. Consistent with this notion, the predicted host range for this collection of plasmids ranges from single genus to multi-phyla, with the most heterogeneous host range profile observed for plasmids found *S. marcescens* (Fig. 6).

A network visualisation of the all-versus-all Mash distances<sup>46</sup> calculated for the *Serratia* plasmids was used to explore their diversity. The resulting network comprises 113 clusters, of which 53 (47%) correspond to singletons, illustrating the diversity of *Serratia* plasmids (Fig. 6). Differences in plasmid abundance between clusters were evident from the network, as four top clusters included 36% of the plasmids identified in *Serratia* genomes. Overall, the plasmids clustering was concordant with their size and GC content but also with the host species (Fig. 6c, Supplementary Fig. 15), suggesting limited between-species plasmid transfer within *Serratia*. Nevertheless, some multi-species clusters were identified, perhaps hinting at recent plasmid acquisition events. A cluster formed by plasmids of four non-*marcescens* species was the largest in the network. This cluster mainly consists of large MOBP conjugative plasmids related to the amber disease

associated plasmid (pADAP), which is required for virulence of *S. entomophila* and *S. proteamaculans* in the larvae of the grass grub *Costelytra zealandica*<sup>22</sup>. Interestingly whilst pathogenic potential in *Costelytra zealandica* is a defining trait of *S. entomophila*, the presence of a pADAP-related plasmid was not universal or a defining trait for either *S. entomophila* or *S. proteamaculans*, being found in members of *S. entomophila*, *S. quinivorans* and *S. proteamaculans*, and a single *S. liquefaciens* genome.

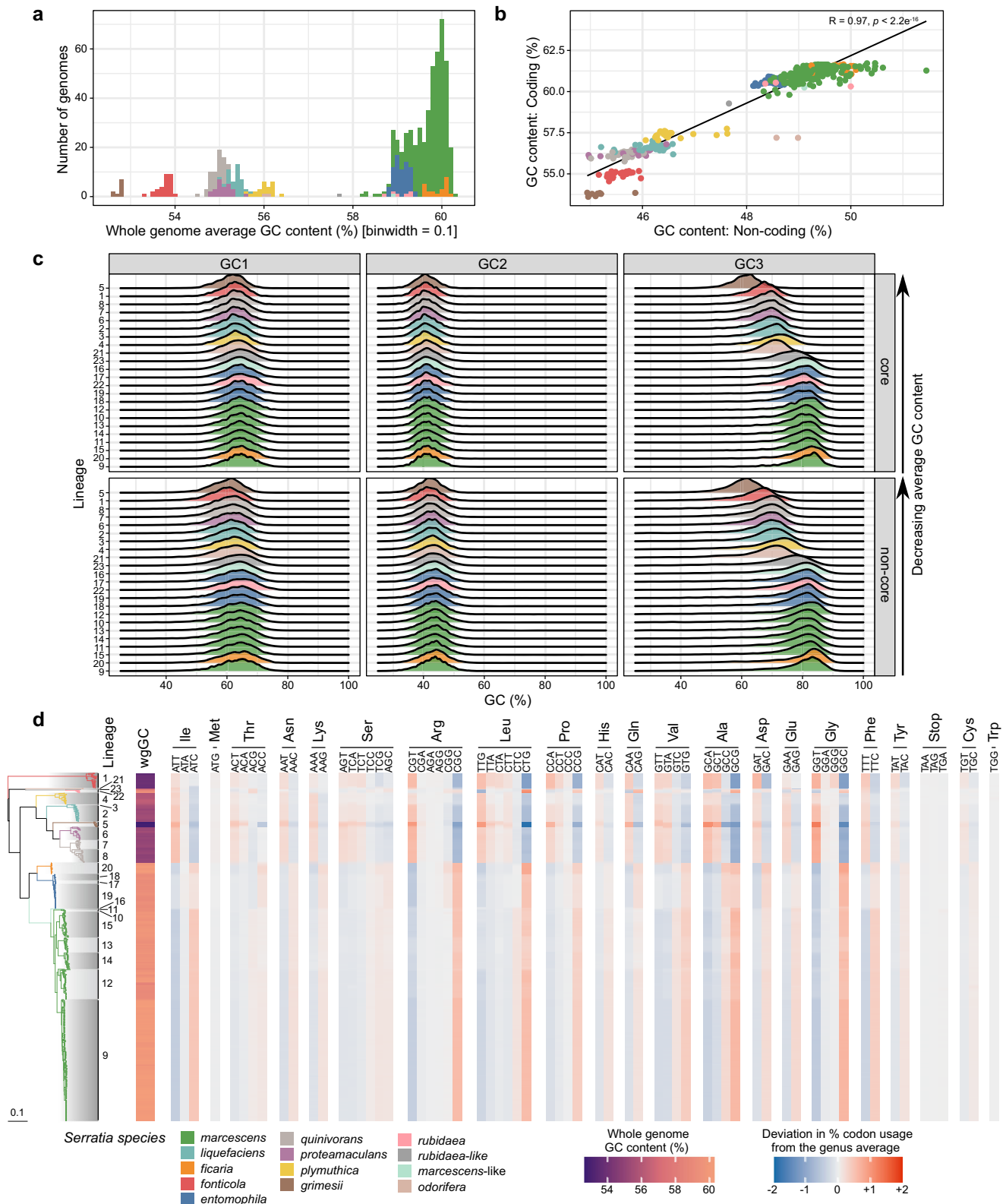
Notably, the predicted host range of the plasmids brings an additional perspective on their potential dynamics within the genus. Most plasmids identified in *S. marcescens* appear to be restricted to this species within *Serratia*. Yet many of them have a predicted host range that goes beyond the taxonomic rank of family, implying transfer outside the genus, including two clusters of small ColRNAI plasmids predicted to cross multiple phyla (Fig. 6). In contrast, the largest plasmid cluster (pADAP-like), featuring multiple non-*marcescens* species, seems to be restricted to the *Serratia* genus. Altogether, this picture may suggest that the ecological niche of *S. marcescens* has favoured plasmid exchange with diverse hosts outside the genus but has also promoted plasmid containment within the species in *Serratia*. The diversity of plasmids identified in *S. marcescens* and their predicted host range thus implies a major role for this species in the gene flow outside the genus and to a lesser but relevant extent within it.

### A genomics perspective on a historical phenotype

A famous characteristic often popularly associated with *Serratia* spp. is the production of the red pigment prodigiosin<sup>17</sup>. However, in fact, prodigiosin production has only been observed in *S. marcescens* biogroups A1, A2 and A6, some *rubidaea* and some *plymuthica* isolates<sup>21</sup>. The *pig* gene cluster comprises fourteen genes (*pigA-pigN*) required for the production of prodigiosin<sup>47</sup>. Searching across the genus for *pig* gene cluster loci and flanking regions showed that, consistent with the earlier biotyping observations, the *pig* cluster is only encoded in certain *S. marcescens*, *S. rubidaea* and *S. plymuthica* genomes (Figs. 3, 7) which are associated with biotypes or biogroups known to be pigmented. In each case, the *pig* gene cluster presents exactly the same contiguous order of genes (*pigA-pigN*). However, notably, it is found in separate genomic loci in each of the three different species, as evidenced by the fact that the flanking regions upstream and downstream of the *pig* cluster are not syntenic between the different species (Fig. 7b). Representative *pig* gene clusters from each species share ~77–80% identity at the nucleotide level (Fig. 7b), which is similar to the shared nucleotide percentage identity between these species at fully syntenic regions in the chromosome. Given the different genomic contexts in which the *pig* cluster is present in these three different species, it would appear that the *pig* gene cluster has been acquired horizontally on at least three separate occasions.

### Discussion

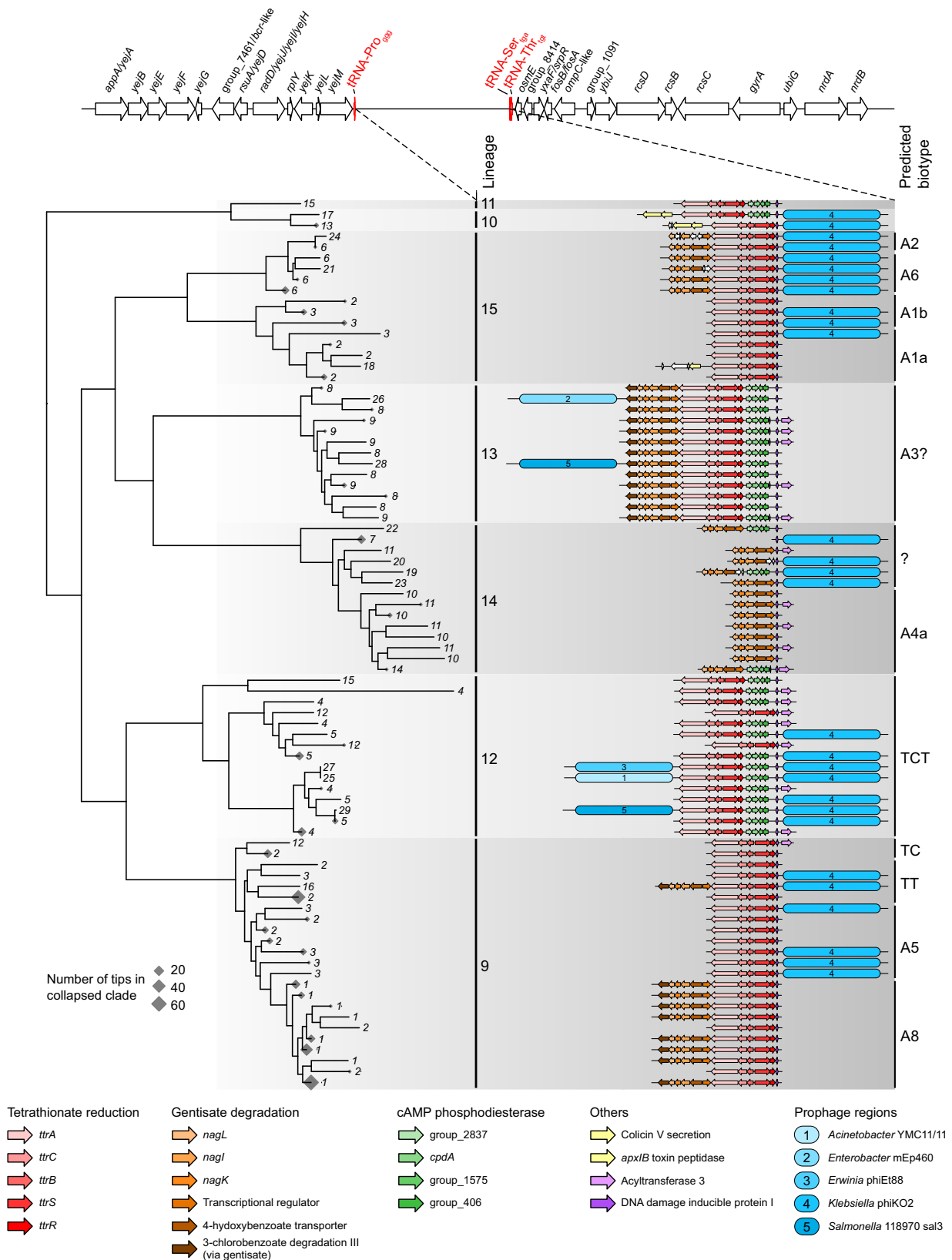
With advancements in technology, the methods used to delineate and decipher prokaryotic species boundaries have changed over time, as researchers attempt to resolve the shortcomings of earlier approaches and build upon the understanding of biology at any given point in time. This study has, in part, investigated the relationship by which species level boundaries have been determined within a genus, namely phenotypic characterisation and whole genome sequencing. It also



**Fig. 4 | *Serratia* is split by GC content.** **a** Histogram of GC content (average over whole genome) across *Serratia*. **b** Plot of GC content of coding regions against that of intergenic (non-coding) regions for each genome; a Pearson correlation test (two-sided) was performed, giving correlation coefficient  $R = 0.97$  and  $p$  value  $< 2.2e^{-16}$  ( $n = 664$ ). **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 (wgGC) content is also shown in the left-most column. Source data are provided as a Source Data file.





highlights how, in order to make appropriate conclusions from these approaches, the currently available data requires to be constantly filtered, checked and reviewed.

Following its original identification in the early 19th century, the nomenclature and number of species within *Serratia* underwent several iterations as additional strains with similar, yet distinct

phenotypes were identified and added to an expanding membership of the genus<sup>17,21</sup>. Then in the 1970s and 80s, comprehensive biochemical and phenotypic characterisation, along with the use of DNA-DNA hybridisation, allowed the genus to be defined as a collection of ten clearly defined species. Since the advent of the genomic era, despite the potential of genomic approaches to

**Fig. 5 | A tRNA-associated hypervariable region ('plasticity zone') encodes gene cassettes for metabolic pathways used for biotyping within *S. marcescens*.** The gene arrangement between the conserved tRNA-Pro<sub>ggg</sub> and tRNA-Ser<sub>gga</sub> in *S. marcescens* is plotted against a maximum-likelihood sub-phylogeny from the tree 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<sub>ggg</sub> and tRNA-Ser<sub>gga</sub> have been pruned from the tree. Each tip number represents a unique combination of genes in the locus. 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.

resolve fine-scaled differences between taxa, no similar-scale work within *Serratia* has been attempted, nor do we have a robust phylogenetic framework against which we are able to recognise novel *Serratia* spp or emerging lineages. Such a framework is also required to resolve confusion over existing species. For example, strain DSM 21420, a nematode-associated strain proposed to belong to *S. nematodiphila*<sup>48</sup>, sits within the broadly non-clinical *S. marcescens* L15, suggesting that it does not in fact represent a separate species. Conversely, the identification of singleton ANI phylogroups and FastBaps lineages (*S. marcescens*-like L16 and *S. rubidaea*-like L23) highlights that there is likely further species diversity to be discovered. This may be partly due to geography and lack of sampling: the strain that occupies L16 (MSU97) was sourced from a plant in the Carrao River in Venezuela<sup>49</sup>, a region which is not highly sampled. These singletons may also represent individual strains that have arisen from single or multiple recombination events, resulting in hybrid or chimeric strains with genomes containing sections of multiple *Serratia* lineages or even other taxa.

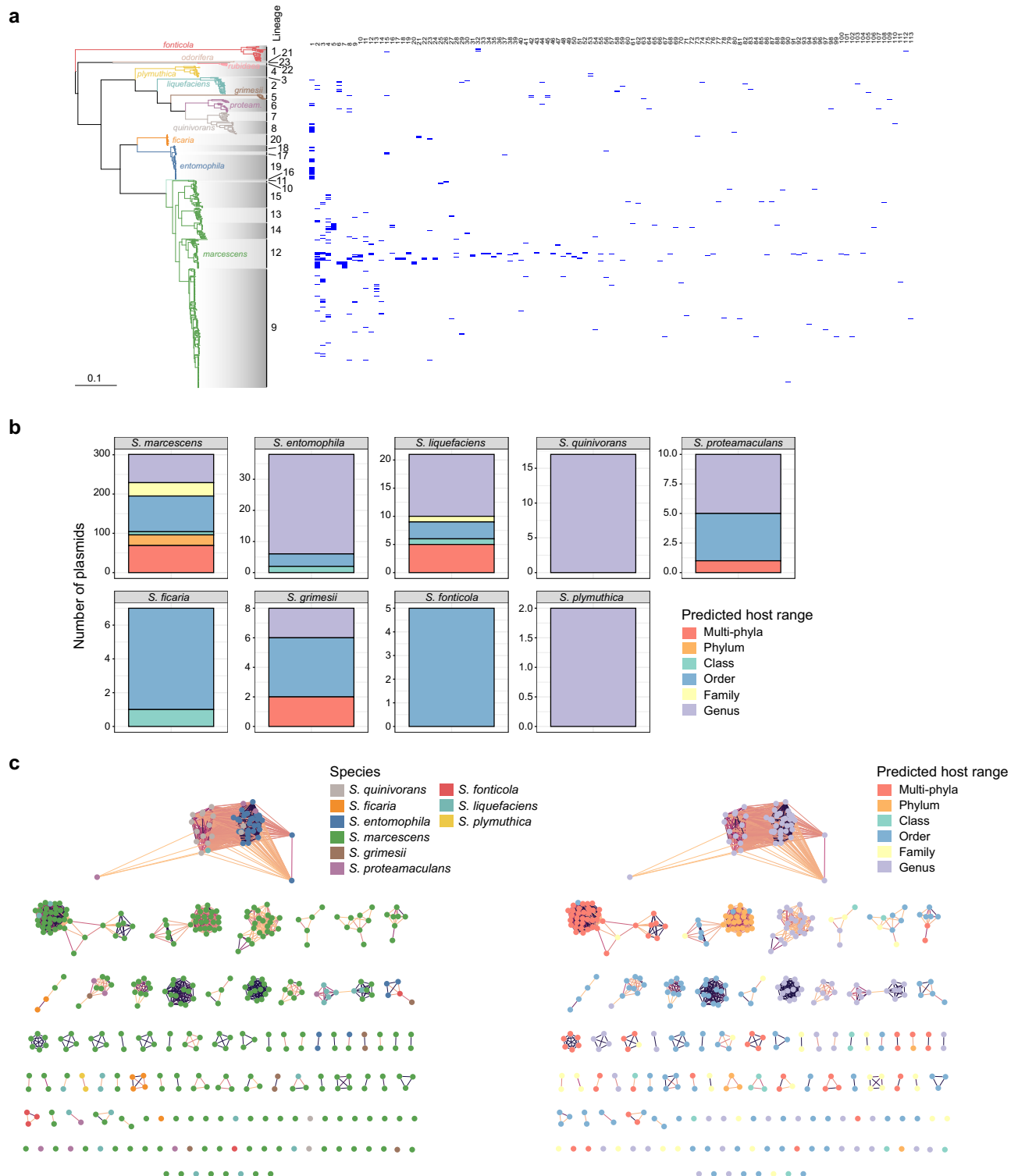
It is interesting to consider how the computational approaches used here to classify and describe the genus parallel the original biotyping. In the earlier studies, in vitro DNA-DNA hybridisation was used to assess genomic relatedness between novel *Serratia* strains<sup>5,11</sup>, an approach for which ANI is in many ways an in silico proxy, whilst the connection between in silico prediction of metabolic potential and the lab-based tests detecting the corresponding metabolic pathway in the original biochemical-based biotyping is obvious. Furthermore, in some cases, these biotypes highlight further clusters within lineages that match branching within the phylogeny. For example, biotypes C1c, EB and RB, and biotypes A1b, A1a, A6 and A2 are all monophyletic within *S. proteamaculans* L6 and *S. marcescens* L15, respectively (Figs. 3b, 5). This highlights just how accurate the original biochemical-based typing was for defining species. However a small number of incongruencies are apparent when comparing predicted and reported biotypes within the *liquefaciens* complex (Supplementary Fig. 10). For example, two strains labelled as biotype RQ (*S. quinivorans*) sit within *S. liquefaciens* biotype C1ab, whilst two genomes situated in *S. proteamaculans* (biotype RQ), were labelled as biotype C1c (*S. grimesii*), perhaps reflecting labelling issues during the long history of the collection.

This accuracy is particularly striking given that we have observed that presence or absence of metabolic pathways (corresponding to the historic biotyping tests used) can be due to repeated gene gain or loss in the same locus over short evolutionary distances. For example, the genes required for the degradation of gentisate and the reduction of tetrathionate are gained and lost within and between lineages in *S. marcescens*, in the same locus and also in the same conserved order (Fig. 5). This would explain why the original phenogrouped biotypes based on biochemical typing had 'variable' results for certain metabolic tests, such as gentisate degradation being observed to be variable in the clinical biotypes A8a and A8c<sup>40</sup>. This locus-specific pathway gain and loss in historic isolates is also seen in more contemporary strains (Fig. 5). The maintenance of this plasticity zone suggests that there are transient and frequently re-occurring environmental selective pressures where the benefit and cost of these pathways are great enough to provide selection both for and against them. In other words, the data suggest that both the loss and re-acquisition of these elements are of benefit to *S. marcescens* at various times. Alternatively, it is also

possible that the pattern observed may reflect multiple events of gene loss alongside these prophage integrations.

It is also noteworthy that the environment from which strains were isolated across our assembled dataset tends to match the environments and niches with which each biotype was historically associated<sup>21</sup>. Of particular interest is the observation that the predominantly hospital-associated biotypes of *S. marcescens* that were defined in the 1980s (A5, A8, TCT) sit within L9 and L12 defined in the current study. These lineages are mainly comprised of recently-sequenced genomes from hospital settings, including a large collection of clinically-derived *S. marcescens* isolates from the UK that represent the recent emergence of hospital-adapted clones exhibiting recent acquisition of MDR phenotypes<sup>18</sup>. The fact that these lineages of clinically-associated *S. marcescens* were identified back in the 1970s and 80s shows that the original biochemical characterisation of *Serratia* captured the emergence of *S. marcescens* lineages that have subsequently been reported to be associated with human disease many times in recent years<sup>18,19,50–55</sup>. The apparent specialisation of *S. marcescens* L9 to be a clinically-adapted pathogen is further highlighted by plasmid replicon identification and the types of lineage-specific core genes observed. The identification of numerous plasmid replicons in these lineages (L9, L12 and L14), as opposed to the rest of the genus, is perhaps unsurprising, given that most known plasmids are associated with multi-drug resistance and hospital environments. Fimbrial genes are well-known pathogenicity factors and multiple different fimbrial genes are found to be core to L9 but accessory to multiple other *S. marcescens* lineages. This potential gene flow from L9 across the rest of *S. marcescens* may be one reason why isolates from more environmental *S. marcescens* lineages are still isolated from nosocomial settings. In these other lineages, *S. marcescens* is still an opportunistic pathogen, with nosocomial isolates being genetically similar to strains that have colonised or infected plants, insects or other environments. Indeed, bee-associated *S. marcescens* cause infections in bees in a similar manner to how *S. marcescens* can cause bloodstream infections in preterm neonates<sup>29</sup>. Taking the historic biotyping data along with the population structure defined here, the combined data suggest that *S. marcescens* is highly plastic in its nature yet can also become specialised in a particular niche.

Speciation and niche specialisation events or processes are seen across the phylogeny, as highlighted by the long branch lengths between divisions, separations in GC content, variation in metabolic potential, and enrichment for certain isolation source sites in different lineages. These divisions likely represent ancient speciation events that have occurred as *Serratia* has spread to be ubiquitous worldwide. As mentioned above, changes in GC content can be a response to long-term niche adaptation, however, there is no commonly held theory or understanding of the possible reasons that underpin this. One possible factor that may have influenced the variation in GC content observed across *Serratia* is a difference in ideal growth temperature: higher GC *Serratia* species tend to be able to grow better at higher temperatures than lower GC *Serratia* species<sup>21</sup>. Another possibility is that the observed GC-dependent change in codon usage, which does not alter protein sequence or function, is indicative of a shift to an optimal set of codons for each particular *Serratia* species, although the evolutionary pressure that would drive such a shift is not clear. Importantly, however, this division in GC content does not seem to be a barrier for gene flow in *Serratia*, since genes core to the high GC species *S. entomophila*



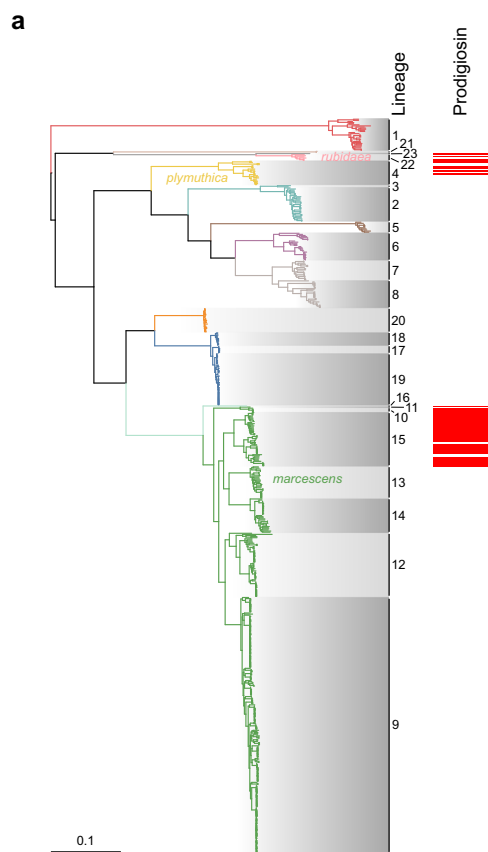
**Fig. 6 | Predicted plasmids across *Serratia*.** **a** Distribution of the 131 plasmid clusters identified against the maximum-likelihood phylogeny of *Serratia* shown in Fig. 1. **b** Number and predicted host range of plasmids identified in *Serratia* genomes. **c** Diversity of *Serratia* plasmids 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). Source data are provided as a Source Data file.

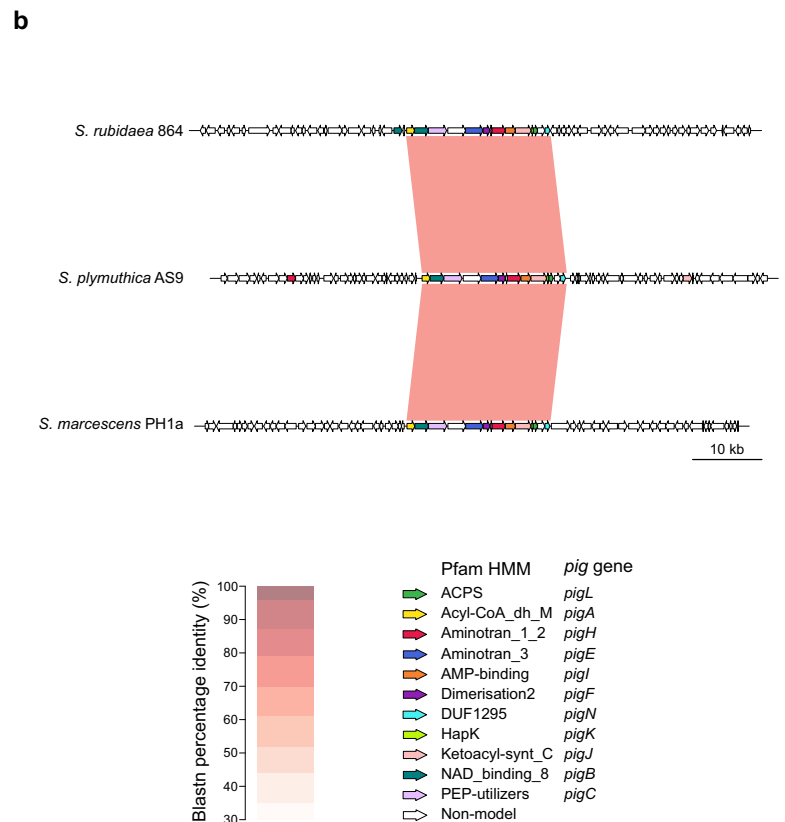
can also be found in polyphyletic and variable patterns across the genus, including in low GC *Serratia* species (Supplementary Fig. 11). However, it is formally possible that these genes could have been horizontally acquired from non-*Serratia* sources.

This study also provides definitive genomic evidence to explain the variation in a classical *Serratia* phenotype, namely the production

of the red pigment prodigiosin (Fig. 7). The high level of synteny within the *pig* gene cluster together with the absence of homology in the flanking regions indicates that the ability to produce prodigiosin has been acquired on at least three separate occasions within *Serratia*, namely in subsets of *S. marcescens*, *S. plymuthica* and *S. rubidaea*. Given the relatively low shared nucleotide identity, it is unclear when



**Fig. 7 | The prodigiosin gene cluster is present variably across *Serratia* and in different genomic loci. a** Prodigiosin (*pig*) gene clusters identified using Hamberger are plotted against the maximum-likelihood phylogeny of *Serratia* shown in



**Fig. 1. b** Pairwise blastn comparison of *pig* loci (core *pig* cluster +/-20 kb) from representative members of the three species containing *pig* genes, extracted using Hamberger.

and how these genes were incorporated into the chromosome, and whether each event reflects gene flow within the genus or separate acquisition from an external source. This genomic evidence of separate acquisition of the *pig* clusters matches the historical metadata noting that *S. marcescens*, *S. plymuthica* and *S. rubidaea* all variably produce a red pigment<sup>21</sup>. Prodigiosin has been reported to display many functions, including anti-protozoal, anti-fungal, anti-bacterial, immunosuppressive, and anti-cancer activity<sup>47</sup>. The biological advantage for these individual *Serratia* species, or subsets thereof, to be able to produce prodigiosin is unclear, however it could reflect a degree of convergent evolution within *Serratia*, or perhaps the varied potential functions of prodigiosin may provide different fitness benefits to different species. Further evidence for convergent evolution in the genus is provided by the observation that members of both *S. proteamaculans* and *S. entomophila* carry pADAP, which is required for the pathogenesis of grass grub larvae<sup>22,56</sup>.

In conclusion, we have demonstrated the power of combining phenotypic metadata with a comprehensive and balanced genomics-based phylogeny to define an important and diverse bacterial genus, its plasticity and its niche adaptation. The dataset and phylogeny that we present here will provide a vital platform for future work, including in the tracking of further emergence of pathogenic *Serratia* or changes in the portfolio of anti-microbial resistance genes or pathogenicity factors.

## Methods

### Bacterial strains

Bacterial isolates sequenced in this study are listed in Supplementary Data 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 and 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 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 ~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 Data 1. 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). Adapter 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 were repaired and dA-tailed using NEBNext End Repair/dA-tailing module, following by ligation of barcodes. DNA concentration and clean up steps were performed using AMPureXP beads (New England Biolabs). 12 samples (from 12 isolates) were multiplexed on a single MinION run. Basecalling and demultiplexing were performed by Albacore v2. In all cases, kits were used according to the manufacturers' instructions.

### Sequence data quality control

Read sets obtained from all samples were compared to the MiniKraken database by Kraken v1.1.1<sup>57</sup>, and then corrected using Bracken v1.0.0<sup>58</sup> which assigns reads to a specific reference sequence, species or genus. If reads were not able to be assigned to a taxonomic class, they were classed as 'unclassified'. Any read sets that belonged to genera other than *Serratia* were discarded from any further analysis, along with any assemblies obtained from those read sets.

Any read sets with more than an estimated five percent of heterozygous SNPs across the whole genome were removed from further analysis, in addition to any assemblies obtained from those read sets. Heterozygous SNPs were calculated using a software pipeline from the pathogen informatics team at the Wellcome Sanger Institute. Specifically, read sets from each *Serratia* sample were aligned to the genome of *S. marcescens* Db11. Reads were aligned using bwa v0.7.17<sup>59</sup>, and parsed using samtools v0.1.19<sup>60</sup> and bcftools v0.1.19<sup>60</sup>. Reads were considered as heterozygous if there were at least two variants at the same base, both supported by a number of reads that was fewer than 90 percent of the total reads mapped to that site. Read coverage to each strand was considered independently. The minimum total coverage required was 4×, and the minimum total coverage for each strand was 2×. Calculated heterozygous SNP coverage was then predicted by scaling the number of observed heterozygous SNPs against the proportion of the reference that was covered by read mapping.

Eight genome sequences from the Pasteur collection dataset and one from the UK hospitals set were removed due to the above criteria. In addition, a number of the isolates resuscitated from the Pasteur collection were duplicate samples of the same strain. After inspection of preliminary phylogenetic trees from core-gene alignments (see below), a further 56 genomes were removed from the Pasteur collection dataset due to being duplicates of the same-named strain.

### Publicly available genome sequences

Previously published, publicly-available assembled genome sequences were downloaded from the NCBI GenBank database (<https://ftp.ncbi.nlm.nih.gov/genomes/genbank/>) as of 19/03/2019. Genomes were downloaded if the species was attributed to any of the following: *Serratia sp.*, *odorifera*, *rubidaea*, *plymuthica*, *liquefaciens*, *grimesii*, *oryzae*, *proteamaculans*, *quinivorans*, *nematodiphila*, *ficaria*, *entomophila* or *marcescens*. Assemblies smaller than 4.5 Mbp or larger than 6.5 Mbp were removed from the analysis, along with any assemblies comprised of more than 250 contigs. Quast v4.6.0<sup>61</sup> was used to extract statistics

for genomes and genomic assemblies, specifically whole genome GC content, number of contigs and assembly size. Initial phylogenetic trees with additional non-*Serratia* reference sequences (*Yersinia enterocolitica*, *Rahnella aquatilis* and *Dickeya solani*) were computed, and genomes determined by visual inspection as being non-*Serratia* or close to non-*Serratia* members of *Enterobacteriaceae* were removed from any subsequent analysis. Ten genomes were excluded on this basis, including several so-called *Serratia sp.* and *Serratia oryzae*.

### Genome assembly and annotation

The assembly method used for genome assembly and annotation for each genome is detailed in Supplementary Data 1. For samples sequenced using short-read only data, genomes were assembled in two different ways depending on their origin. Isolates in the Institut Pasteur collection were assembled through assembly pipelines at the Wellcome Sanger Institute. For each sample, sequence reads were used to create multiple assemblies using VelvetOptimiser v2.2.5 (<https://github.com/tseemann/VelvetOptimiser>) and Velvet v1.2<sup>62</sup>. An assembly improvement step was applied to the assembly with the best N50 and contigs were scaffolded using SSPACE v2.0<sup>63</sup> and sequence gaps filled using GapFiller v1.11<sup>64</sup>. For isolates from UK hospitals that were only sequenced by short-read technology, these short reads were assembled using SPAdes v3.6.1<sup>65</sup>, using default settings.

For hybrid short- and long-read assemblies of selected isolates from UK hospitals, genomes were assembled using Unicycler v0.4.7<sup>66</sup>. Long-read-only assemblies from MinION or PacBio long reads were generated first, using Canu v1.6<sup>67</sup>, with the expected genome size set as 5.4 Mbps, the minimum read length and overlap length set to 100 bp, and 'corOutCoverage' set to 1000. Long-read assemblies were then used as input to Unicycler, using the `-existing_long_read_assembly` flag. Sets of paired-end Illumina reads were then used as input to Unicycler alongside this long-read assembly and also the long reads. The `-mode` flag was set to 'normal'. In the event that Unicycler was not able to produce circularised assemblies, Circlator v1.5<sup>68</sup> was used to circularise assemblies.

Assembled genomes were then annotated using Prokka v1.13.3<sup>69</sup>.

### Pan-genome analysis

Pan-genomes were calculated from 664 *Serratia* sequences using Panaroo v1.2.3<sup>70</sup>, with Prokka-annotated genomes as input. For initial protein clustering within Panaroo, a protein sequence identity similarity threshold was set at 95 percent (0.95) and the length difference cutoff was set at 95 percent (0.95). The subsequent clustering of these groups into protein families was performed using a threshold of 70 percent identity (0.7). The `-clean-mode` flag was set to 'moderate'. A core-gene alignment was created using the `-a` flag, specifying mafft as the aligner using the `-aligner` flag, with core genes specified by being present in at least 95 percent of genomes (631/664). Pan-genome gene accumulation curves were generated using the *specaccum* function from the R package Vegan v2.5.7<sup>71</sup>, with 100 random permutations.

Population structure-aware classification of genes across the genus was performed upon the gene presence/absence matrix created by Panaroo through the use of the twilight analysis package<sup>41</sup>. Groups were defined by the lineages set by the third level of FastBaps clustering (see below), and singleton lineages were included in the analysis (`-min_size 1`). The core and rare thresholds were set at 0.95 and 0.15, respectively.

Preliminary core-gene alignments using the pan-genome software Roary v3.12.0<sup>72</sup>, including all downloaded genomes from the NCBI GenBank database, duplicate genomes from the Pasteur collection and non-*Serratia* Enterobacteriaceae members, were computed for initial tree-drawing to remove contaminants and assess whether duplicate strains (from data supplied in strain name information, for example, labels on agar stabs from strains in the Pasteur collection) were found in the same position in the tree. Non-*Serratia* Enterobacteriaceae were

also used to determine the location of the root for all visualisations of the *Serratia* genus phylogenetic tree.

### Clustering, phylogroup determination, core-gene alignment filtering and phylogenetic tree construction

For the *Serratia* phylogeny, a concatenated core-gene alignment from 2252 genes (2,820,212 bp in length) from Panaroo v1.2.3<sup>70</sup> (as described above) was filtered to remove monomorphic sites and only keep columns that were exclusively A, T, G or C using SNP-sites v2.5.1<sup>73</sup>. The resulting alignment was 398,551 bp in length. IQtree v1.6.10<sup>74</sup> was then used for maximum-likelihood tree construction using 1000 ultrafast bootstraps<sup>75</sup> using the TIMe+ASC+R4 model chosen using modelfinder<sup>76</sup>. Both the ultrafast bootstraps and modelfinder were implemented in IQtree. The *Serratia* phylogenetic tree was rooted at the position of a *Yersinia enterocolitica* outgroup root after analysis of preliminary trees based on exclusively polymorphic variant sites (filtered using SNP-sites v2.4.1) from preliminary core-gene alignments (determined using Roary v3.12.0 as described above). Trees were constructed using modelfinder implemented in IQtree v1.6.10, followed by tree construction using IQtree v1.6.10. An additional tree was also drawn using the same alignment but with monomorphic sites retained. Both tree files are available alongside all other raw data using for plotting through Figshare (<https://doi.org/10.6084/m9.figshare.18051824>).

Whole-genome assemblies were compared in a pairwise manner using fastANI v1.3<sup>77</sup>, and phylogroups determined through clustering these comparisons using a cutoff of 95% average nucleotide identity (ANI). Genomic assemblies were then clustered based on this cutoff value, using the script fastANI\_to\_clusters.py which uses the networkx package (<https://networkx.github.io/>), and visualised using Cytoscape v3.7.1<sup>78</sup>. The phylogeny was partitioned into lineages defined through hierarchical bayesian clustering using FastBaps v1.0.4<sup>79</sup>. FastBaps was used to cluster the phylogeny over four levels, with the third level selected for lineage designation. The SNP sites-filtered core-gene alignment was used as input to FastBaps, alongside the rooted phylogenetic tree to provide a guide for the hierarchical partitioning.

### Functional and metabolic pathway analysis

In silico reconstruction of metabolic pathways was performed using Pathway tools v23.5<sup>42</sup>, using a multi-processing wrapper tool mpwt (<https://github.com/AuReMe/mpwt>)<sup>80</sup>. 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](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<sup>81</sup>, 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 Commission (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 mpwt with the ‘-path’ and ‘-taxon\_id’ flags, whilst providing a file containing taxon ids for each genome. The in silico-reconstructed metabolic pathways for all genomes were then collated using compare\_pgdb.R in pathwaytools\_gff2gbk, and downstream analysis conducted in R, as shown in [https://github.com/djw533/Serratia\\_genus\\_paper/figure\\_scripts](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<sup>82</sup>. 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<sup>83</sup>. K-mer-based sketches of the plasmid sequences ( $s = 1000$ ,  $k = 21$ ) were generated with the mash v2.3 sketch algorithm<sup>46</sup>. 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<sup>78</sup> using the ‘connected components cluster’ algorithm from the clusterMaker2 v2 app<sup>84</sup>.

### GC content analysis

Whole-genome GC content was 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 (detailed above). Intragenic nucleotide sequence was extracted for all protein-encoding sequences using gals\_parser\_with\_fasta.py with the ‘-t nuc’ flag. Intergenic GC values were then calculated by using Bedtools<sup>85</sup> complement from Bedtools v2.29.0 to identify the inverse of all coding regions (i.e. all intergenic regions). Bedtools getfasta from Bedtools v2.29.0 was then used to extract the intergenic regions as nucleotide sequence. Average GC values for the total intergenic and intragenic regions were then calculated using get\_gc\_content.py.

### Retrieval of specific gene clusters

Gene clusters containing co-localised *pig* genes (*pigA-M*) were identified using Hamburger ([github.com/djw533/hamburger](https://github.com/djw533/hamburger)), which uses protein HMM profiles for each target gene in the gene cluster. User-set parameters define the minimum number of hmmsearch (from HMMER suite v3.2.1<sup>86</sup>) hits required to report the presence of each system in a genome, in addition to the maximum number of non-hit genes that are permitted between two hit genes in a contiguous set of genes. Gene clusters were reported as prodigiosin clusters for loci encoding at least nine genes containing Pfam domains characteristic of 11 of the 14 *pig* genes with no more than five non-model genes between any hits. Extracted genomic sequences were then compared using blast+ v2.2.31<sup>87</sup> and genoplots v0.8.11<sup>88</sup>. Blastn was used with the flags ‘-task Blastn -perc\_identity 20 -evaluate 10000’. Functions created to use these can be found in micro.gen.extra on <https://github.com/djw533/micro.gen.extra>.

Gene clusters around other genes of interest, such as the plasticity zone in *S. marcescens* located between tRNA-Pro<sub>ggg</sub> and tRNA-Ser<sub>tga</sub>, were extracted using the script pull\_out\_around\_point.py, and, if in the unwanted orientation, flipped using gff\_reverse.py.

### Phage prediction

Phage regions from the 29 representative hypervariable regions between tRNA-Pro<sub>ggg</sub> and tRNA-Ser<sub>tga</sub> in *S. marcescens* were predicted using Phaster<sup>89</sup> on the webserver (<https://phaster.ca/>), using default settings.

### Data visualisation

Phylogenetic trees were visualised using the R package ggtree v2.4.2<sup>90</sup>. Synteny of regions of bacterial genomes extracted by Hamburger were visualised using the R package genoplots v0.8.11<sup>88</sup>. Genetic organisation of genes were plotted using the R package gggenes v0.4.1 (<https://wilcox.org/gggenes/>). Other plots were created using the R package ggplot2 v3.3.5<sup>83</sup>. As mentioned above, networks were viewed using

Cytoscape v3.7.1<sup>78</sup>. Sets were visualised as Upset plots using UpsetR v1.4.0<sup>91</sup>.

### Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

### Data availability

The sequencing data generated in this study have been deposited in the European Nucleotide archive (ENA; <https://www.ebi.ac.uk/ena/>), under accession numbers ERP106480 for the Pasteur collection (<https://www.ebi.ac.uk/ena/browser/view/PRJEB24638>) and ERP135711 for the UK hospitals collection (<https://www.ebi.ac.uk/ena/browser/view/PRJEB51113>). The annotated genome assemblies for both collections have been deposited in the ENA under the project accession numbers above: ERP106480 for the Pasteur collection (<https://www.ebi.ac.uk/ena/browser/view/PRJEB24638>) and ERP135711 for the UK hospitals collection (<https://www.ebi.ac.uk/ena/browser/view/PRJEB51113>). The other whole genome sequences used in this study are available in the NCBI GenBank (<https://ftp.ncbi.nlm.nih.gov/genomes/genbank/>), with the accession numbers for the individual sequences provided in Supplementary Data 1. The gene annotations generated for all 664 genomes analysed in this study, together with all the genome assemblies generated during this study, are available through Figshare (<https://doi.org/10.6084/m9.figshare.18051824>)<sup>92</sup>. All other data used in figure generation, including the output of pan-genome analysis, are available in the same repository (<https://doi.org/10.6084/m9.figshare.18051824>)<sup>92</sup>. The MiniKraken reference database used for species identification is available at [https://ccb.jhu.edu/software/kraken/dl/minikraken\\_20171019\\_8GB.tgz](https://ccb.jhu.edu/software/kraken/dl/minikraken_20171019_8GB.tgz) and the MOB-suite reference database used for plasmid replicon identification is available from Zenodo (<https://doi.org/10.5281/zenodo.3785612>). Source data are provided with this paper.

### Code availability

All custom scripts for which github repositories are not specified above can be found at [https://github.com/djw533/Serratia\\_genus\\_paper/analysis\\_scripts](https://github.com/djw533/Serratia_genus_paper/analysis_scripts). Other packages used can be found at <https://github.com/djw533/hamburger><sup>93</sup>, <https://github.com/djw533/micro.gen.extra><sup>94</sup>, and [https://github.com/djw533/pathwaytools\\_gff2gbk](https://github.com/djw533/pathwaytools_gff2gbk)<sup>95</sup>. Rscripts used to plot figures can also be found in this repository at [https://github.com/djw533/Serratia\\_genus\\_paper/figure\\_scripts](https://github.com/djw533/Serratia_genus_paper/figure_scripts). Rscripts make use of the tidyverse<sup>96</sup> collection of packages. R version 4.0.3 was used for all analysis and generation of plots.

### References

- Merlino, C. P. Bartolomeo Bizio's letter to the most eminent priest, Angelo Bellani, concerning the phenomenon of the red-colored polenta [translated from the Italian]. *J. Bacteriol.* **9**, 527–543 (1924).
- Grimont, P. A. D. & Dulong de Rosnay, H. L. C. Numerical study of 60 strains of *Serratia*. *J. Gen. Microbiol.* **72**, 259–268 (1972).
- Grimont, P. A. D., Grimont, F. & Dulong de Rosnay, H. L. C. Taxonomy of the genus *Serratia*. *J. Gen. Microbiol.* **98**, 39–66 (1977).
- Grimont, F., Grimont, P. A. D. & Dulong de Rosnay, H. L. C. Characterization of *Serratia marcescens*, *S. liquefaciens*, *S. plymuthica* and *S. marinorubra* by electrophoresis of their proteinases. *J. Gen. Microbiol.* **99**, 301–310 (1977).
- Grimont, P. A. D. et al. Deoxyribonucleic acid relatedness between *Serratia plymuthica* and other *Serratia* species, with a description of *Serratia odorifera* sp. nov. (Type strain: ICPB 3995). *Int. J. System. Bacteriol.* **28**, 453–463 (1978).
- Grimont, P. A. D., Grimont, F. & Starr, M. P. *Serratia ficaria* sp. nov., a bacterial species associated with Smyrna figs and the fig wasp *Blastophaga psenes*. *Curr. Microbiol.* **2**, 277–282 (1979).
- Gavini, F. et al. *Serratia fonticola*, a new species from water. *Int. J. System. Bacteriol.* **29**, 92–101 (1979).
- Holmes, B. Proposal to conserve the specific epithet *liquefaciens* over the specific epithet *proteamaculans* in the name of the organism currently known as *Serratia liquefaciens* (Grimes and Hennerty 1931) Bascomb et al. 1971. Request for an opinion. *Int. J. Systemic Bacteriol.* **30**, 220–222 (1980).
- Grimont, P. A. D., Grimont, F. & Starr, M. P. *Serratia* species isolated from plants. *Curr. Microbiol.* **5**, 317–322 (1981).
- Grimont, P. A. D., Grimont, F. & Irino, K. Biochemical characterization of *Serratia liquefaciens sensu stricto*, *Serratia proteamaculans*, and *Serratia grimesii* sp. nov. *Curr. Microbiol.* **7**, 69–74 (1982).
- Grimont, P. A. D., Irino, K. & Grimont, F. The *Serratia liquefaciens*-*S. proteamaculans*-*S. grimesii* complex: DNA relatedness. *Curr. Microbiol.* **7**, 63–67 (1982).
- 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*. *Int. J. Syst. Bacteriol.* **38**, 1–6 (1988).
- Murdoch, S. L. et al. The opportunistic pathogen *Serratia marcescens* utilizes Type VI secretion to target bacterial competitors. *J. Bacteriol.* **193**, 6057–6069 (2011).
- 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 RhIA-dependent surfactant biosynthesis in *Serratia*. *Environ. Microbiol.* **10**, 1202–1217 (2008).
- Kurz, C. L. et al. Virulence factors of the human opportunistic pathogen *Serratia marcescens* identified by in vivo screening. *EMBO J.* **22**, 1451–1460 (2003).
- Khanna, A., Khanna, M. & Aggarwal, A. *Serratia marcescens*—a rare opportunistic nosocomial pathogen and measures to limit its spread in hospitalized patients. *J. Clin. Diagn. Res.* **7**, 243–246 (2013).
- Mahlen, S. D. *Serratia* infections: from military experiments to current practice. *Clin. Microbiol. Rev.* **24**, 755–791 (2011).
- 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 Res.* **26**, 1101–1109 (2016).
- Karkey, A. et al. Outbreaks of *Serratia marcescens* and *Serratia rubidaea* bacteremia in a central Kathmandu hospital following the 2015 earthquakes. *Trans. R. Soc. Trop. Med. Hyg.* **112**, 467–472 (2018).
- Dubouix, A. et al. Epidemiological investigation of a *Serratia liquefaciens* outbreak in a neurosurgery department. *J. Hosp. Infect.* **60**, 8–13 (2005).
- Grimont, F. & Grimont, P. A. D. The Genus *Serratia*. in *Prokaryotes* (eds. Martin Dworkin, Stanley Falkow, Eugene Rosenberg, Karl-Heinz Schleifer & Erko Stackebrandt), Volume 6, 219–244 (Springer-Verlag, 2006).
- 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 *Photobacterium luminescens*. *J. Bacteriol.* **182**, 5127–5138 (2000).
- 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*. *J. Bacteriol.* **186**, 5116–5128 (2004).
- Núñez-Valdez, M. E. et al. Identification of a putative Mexican strain of *Serratia entomophila* pathogenic against root-damaging larvae of Scarabaeidae (Coleoptera). *Appl. Environ. Microbiol.* **74**, 802–810 (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*. *J. Invertebr. Pathol.* **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. *Appl. Environ. Microbiol.* **84**, e02739-17 (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*. *J. Invertebr. Pathol.* **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. *Int. J. Syst. Evolut. Microbiol.* **52**, 2281–2289 (2002).
31. Lim, Y.-L. L. et al. Complete genome sequence of *Serratia fonticola* DSM 4576T, a potential plant growth promoting bacterium. *J. Biotechnol.* **214**, 43–44 (2015).
32. Abebe-Akele, F. et al. Genome sequence and comparative analysis of a putative entomopathogenic *Serratia* isolated from *Cae-norhabditis briggsae*. *BMC Genom.* **16**, 531 (2015).
33. Petersen, L. M. & Tisa, L. S. Friend or foe? A review of the mechanisms that drive *Serratia* towards diverse lifestyles. *Can. J. Microbiol.* **59**, 627–640 (2013).
34. Cheng, T. H. et al. Genome sequence of *Serratia marcescens* subsp. *sakuensis* strain K27, a marine bacterium isolated from sponge (*Haliclona amboinensis*). *Genome Announc.* **6**, e00022-18 (2018).
35. Matilla, M. A., Udaondo, Z. & Salmond, G. P. C. Genome sequence of the Oocystin A-producing rhizobacterium *Serratia plymuthica* 4Rx5. *Microbiol. Resour. Announcements* **7**, e00997-18 (2018).
36. Chen, S., Blom, J. & Walker, E. D. Genomic, physiologic, and symbiotic characterization of *Serratia marcescens* strains isolated from the mosquito *Anopheles stephensi*. *Front. Microbiol.* **8**, 1483 (2017).
37. Ward, D. V. et al. Metagenomic sequencing with strain-level resolution implicates uropathogenic *E. coli* in necrotizing enterocolitis and mortality in preterm infants. *Cell Rep.* **14**, 2912–2924 (2016).
38. 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 Genet.* **11**, e1005413 (2015).
39. Jain, C., Rodríguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat. Commun.* **9**, 5114 (2018).
40. 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).
41. Horesh, G. et al. Different evolutionary trends form the twilight zone of the bacterial pan-genome. *Microb. Genom.* **7**, 000670 (2021).
42. Karp, P. D. et al. Pathway Tools version 23.0 update: software for pathway/genome informatics and systems biology. *Brief. Bioinform.* **22**, 109 (2021).
43. Foerstner, K. U., von Mering, C., Hooper, S. D. & Bork, P. Environments shape the nucleotide composition of genomes. *EMBO Rep.* **6**, 1208–1213 (2005).
44. Palidwor, G. A., Perkins, T. J. & Xia, X. A general model of codon bias due to GC mutational bias. *PLOS ONE* **5**, e13431 (2010).
45. Reuter, S. et al. Parallel independent evolution of pathogenicity within the genus *Yersinia*. *Proc. Natl Acad. Sci. USA* **111**, 6768–6773 (2014).
46. Ondov, B. D. et al. Mash: Fast genome and metagenome distance estimation using MinHash. *Genome Biol.* **17**, 132 (2016).
47. Harris, A. K. P. et al. The *Serratia* gene cluster encoding biosynthesis of the red antibiotic, prodigiosin, shows species- and strain-dependent genome context variation. *Microbiology* **150**, 3547–3560 (2004).
48. Kwak, Y., Khan, A. R. & Shin, J.-H. Genome sequence of *Serratia nematodiphila* DSM 21420T, a symbiotic bacterium from entomopathogenic nematode. *J. Biotechnol.* **193**, 1–2 (2015).
49. Matilla, M. A., Udaondo, Z., Krell, T. & Salmond, G. P. C. Genome sequence of *Serratia marcescens* MSU97, a plant-associated bacterium that makes multiple antibiotics. *Genome Announc.* **5**, (2017).
50. Cristina, M. L., Sartini, M. & Spagnolo, A. M. *Serratia marcescens* infections in neonatal intensive care units (NICUs). *Int. J. Environ. Res. Public Health* **16**, (2019).
51. Daoudi, A., Benaoui, F., el Idrissi Slitine, N., Sora, N. & Rabou Maoulainine, F. M. An outbreak of *Serratia marcescens* in a Moroccan neonatal intensive care unit. *Adv. Med.* **2018**, 1–4 (2018).
52. Moles, L. et al. *Serratia marcescens* colonization in preterm neonates during their neonatal intensive care unit stay. *Antimicrob. Resistance Infect. Control* **8**, 135 (2019).
53. Martineau, C. et al. *Serratia marcescens* outbreak in a neonatal intensive care unit: New insights from next-generation sequencing applications. *J. Clin. Microbiol.* **56**, 148–154 (2018).
54. Escribano, E. et al. Influence of a *Serratia marcescens* outbreak on the gut microbiota establishment process in low-weight preterm neonates. *PLOS ONE* **14**, e0216581 (2019).
55. Montagnani, C. et al. *Serratia marcescens* outbreak in a neonatal intensive care unit: Crucial role of implementing hand hygiene among external consultants. *BMC Infect. Dis.* **15**, 11 (2015).
56. 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).
57. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.* **15**, 1–12 (2014).
58. Lu, J., Breitwieser, F. P., Thielen, P. & Salzberg, S. L. Bracken: estimating species abundance in metagenomics data. *PeerJ Comput. Sci.* **2017**, e104 (2017).
59. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
60. Li, H. et al. The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
61. Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* **29**, 1072–1075 (2013).
62. Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read assembly using de Bruijn graphs. *Genome Res.* **18**, 821–829 (2008).
63. Boetzer, M., Henkel, C. V., Jansen, H. J., Butler, D. & Pirovano, W. Scaffolding pre-assembled contigs using SSPACE. *Bioinformatics* **27**, 578–579 (2011).
64. Boetzer, M. & Pirovano, W. Toward almost closed genomes with GapFiller. *Genome Biol.* **13**, 1–9 (2012).
65. Bankevich, A. et al. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).
66. Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: resolving bacterial genome assemblies from short and long sequencing reads. *PLOS Comput. Biol.* **13**, e1005595 (2017).



67. Koren, S. et al. Canu: Scalable and accurate long-read assembly by adaptive  $\kappa$ -mer weighting and repeat separation. *Genome Res.* **27**, 722–736 (2017).
68. Hunt, M. et al. Circlator: automated circularization of genome assemblies using long sequencing reads. *Genome Biol.* **16**, 294 (2015).
69. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).
70. Tonkin-Hill, G. et al. Producing polished prokaryotic pangenomes with the Panaroo pipeline. *Genome Biol.* **21**, 1–21 (2020).
71. Dixon, P. VEGAN, a package of R functions for community ecology. *J. Veg. Sci.* **14**, 927–930 (2003).
72. Page, A. J. et al. Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**, 3691–3693 (2015).
73. Page, A. J. et al. SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments. *Microb. Genom.* **2**, e000056 (2016).
74. 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. *Mol. Biol. Evolution* **32**, 268 (2015).
75. Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2: improving the ultrafast bootstrap approximation. *Mol. Biol. Evol.* **35**, 518–522 (2018).
76. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jeremiin, L. S. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods* **14**, 587–589 (2017).
77. 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. *Nat. Commun.* **2018** 9:1 **9**, 1–8 (2018).
78. Shannon, P. et al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
79. Tonkin-Hill, G., Lees, J. A., Bentley, S. D., Frost, S. D. W. & Corander, J. Fast hierarchical Bayesian analysis of population structure. *Nucleic Acids Res.* **47**, 5539 (2019).
80. Belcour, A. et al. Inferring biochemical reactions and metabolite structures to understand metabolic pathway drift. *iScience* **23**, 100849 (2020).
81. 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 Res.* **47**, D309–D314 (2019).
82. 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. *Microb. Genom.* **6**, 1–12 (2020).
83. Wickham, Hadley. *Ggplot2: Elegant Graphics for Data Analysis*. (Springer, 2009).
84. Morris, J. H. et al. ClusterMaker: a multi-algorithm clustering plugin for Cytoscape. *BMC Bioinform.* **12**, 1–14 (2011).
85. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
86. Eddy, S. R. Accelerated profile HMM searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
87. Camacho, C. et al. BLAST+: architecture and applications. *BMC Bioinform.* **10**, 421 (2009).
88. Guy, L., Kultima, J. R., Andersson, S. G. E. & Quackenbush, J. GenPlotR: comparative gene and genome visualization in R. *Bioinformatics* **27**, 2334–2335 (2011).
89. Arndt, D. et al. PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic Acids Res.* **44**, W16 (2016).
90. 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 Ecol. Evol.* **8**, 28–36 (2017).
91. 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).
92. Williams, D. W. et al. The genus *Serratia* revisited by genomics. *Figshare* <https://doi.org/10.6084/m9.figshare.18051824> (2022).
93. David Williams. djw533/hamburger: Initial release. (2022) <https://doi.org/10.5281/zenodo.6981393>
94. David Williams. djw533/micro.gen.extra:. (2022) <https://doi.org/10.5281/zenodo.6981456>.
95. David Williams. djw533/pathwaytools\_gff2gbk:. (2022) <https://doi.org/10.5281/zenodo.6981411>.
96. Wickham, H. et al. Welcome to the Tidyverse. *J. Open Source Softw.* **4**, 1686 (2019).

## Acknowledgements

This work was supported by Wellcome (grant numbers: 104556, Senior Research Fellowship S.J.C.; 220321, Senior Research Fellowship Renewal S.J.C.; 109118, Ph.D. studentship; 206194, N.R.T.), the NIHR (NIHR200639, AMR Capital Award to University of Dundee), and Institut Pasteur and INSERM (P.A.D.G. and F.X.W.). Firstly, we would like to acknowledge the contribution of, and thank, all those colleagues who contributed over many years to the collection of the *Serratia* isolates forming the Institut Pasteur collection of Patrick Grimont. We also thank Alistair Leanord, Teresa Inkster, James Chalmers, Gillian Orange and Nigel Smith for providing recent isolates of *Serratia marcescens* from UK hospitals, and Hazel Auken and George Salmond for sharing isolates reported previously. We thank Sally Kay, Liz McMinn and Florence Juglas for logistical support, the Wellcome Sanger Institute (WSI) sequencing teams for processing these samples, and Christoph Puethe and the WSI Pathogen Informatics team for help with data management. We thank Gal Horesh, Mat Beale and Matt Dorman for expert technical advice and valuable discussions. For the purpose of Open Access, the authors have applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

## Author contributions

D.J.W., N.R.T. and S.J.C. conceived the study; D.J.W. performed the bioinformatics analyses, with contributions from A.C.L. and D.C.L.; P.A.D.G., F.G. and E.A. performed identification and biochemical characterisation of *Serratia* isolates in the Institut Pasteur collection; D.J.W., K.P., E.N. 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. and S.J.C. analysed and interpreted results; D.J.W., N.R.T. and S.J.C. wrote the paper with input from the other authors.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41467-022-32929-2>.

**Correspondence** and requests for materials should be addressed to Nicholas R. Thomson or Sarah J. Coulthurst.

**Peer review information** *Nature Communications* thanks Louis-Marie Bobay, Fernando Gonzalez-Candelas and the other, anonymous, reviewer for their contribution to the peer review of this work.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022