

UC Davis

UC Davis Previously Published Works

Title

Meta-analysis of the *Ralstonia solanacearum* species complex (RSSC) based on comparative evolutionary genomics and reverse ecology

Permalink

<https://escholarship.org/uc/item/6mj3c2mk>

Journal

Microbial Genomics, 8(3)

ISSN

2057-5858

Authors

Sharma, Parul
Johnson, Marcela A
Mazloom, Reza
et al.

Publication Date

2022-03-17

DOI

10.1099/mgen.0.000791

Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

Meta-analysis of the *Ralstonia solanacearum* species complex (RSSC) based on comparative evolutionary genomics and reverse ecology

Parul Sharma^{1,2}, Marcela A. Johnson^{1,2}, Reza Mazloom³, Caitilyn Allen⁴, Lenwood S. Heath³, Tiffany M. Lowe-Power^{5,*} and Boris A. Vinatzer^{1,*}

Abstract

Ralstonia solanacearum species complex (RSSC) strains are bacteria that colonize plant xylem tissue and cause vascular wilt diseases. However, individual strains vary in host range, optimal disease temperatures and physiological traits. To increase our understanding of the evolution, diversity and biology of the RSSC, we performed a meta-analysis of 100 representative RSSC genomes. These 100 RSSC genomes contain 4940 genes on average, and a pangenome analysis found that there are 3262 genes in the core genome (~60% of the mean RSSC genome) with 13128 genes in the extensive flexible genome. A core genome phylogenetic tree and a whole-genome similarity matrix aligned with the previously named species (*R. solanacearum*, *R. pseudosolanacearum*, *R. syzygii*) and phylotypes (I–IV). These analyses also highlighted a third unrecognized sub-clade of phylotype II. Additionally, we identified differences between phylotypes with respect to gene content and recombination rate, and we delineated population clusters based on the extent of horizontal gene transfer. Multiple analyses indicate that phylotype II is the most diverse phylotype, and it may thus represent the ancestral group of the RSSC. We also used our genome-based framework to test whether the RSSC sequence variant (sequevar) taxonomy is a robust method to define within-species relationships of strains. The sequevar taxonomy is based on alignments of a single conserved gene (*egl*). Although sequevars in phylotype II describe monophyletic groups, the sequevar system breaks down in the highly recombinogenic phylotype I, which highlights the need for an improved, cost-effective method for genotyping strains in phylotype I. Finally, we enabled quick and precise genome-based identification of newly sequenced RSSC strains by assigning Life Identification Numbers (LINs) to the 100 strains and by circumscribing the RSSC and its sub-groups in the LINbase Web service.

DATA SUMMARY

The authors confirm that all raw data and code and protocols have been provided within the manuscript. All publicly available sequencing data used for analysis have been supplemented with accession numbers to access the data. The assembled genome of strain 19-3PR_UW348 was submitted to NCBI under Bioproject PRJNA775652 Biosample SAMN22612291. This Whole Genome Shotgun project has been deposited at GenBank under accession JAJMMU000000000. The version described in this paper is version JAJMMU010000000. Supplementary Material can be found at: <https://doi.org/10.6084/m9.figshare.19068134> [1].

INTRODUCTION

Named species generally correspond to groups of bacteria with pairwise genome similarity over a 95% average nucleotide identity (ANI) threshold and that also share a core set of phenotypes [2]. Bacterial plant pathogens rarely conform to this

Received 24 November 2021; Accepted 02 February 2022; Published 17 March 2022

Author affiliations: ¹School of Plant and Environmental Sciences, Virginia Tech, Blacksburg, VA, USA; ²Graduate Program in Genetics, Bioinformatics and Computational Biology, Virginia Tech, Blacksburg, VA, USA; ³Department of Computer Science, Virginia Tech, Blacksburg, VA, USA; ⁴Department of Plant Pathology, University of Wisconsin-Madison, Madison, WI, USA; ⁵Department of Plant Pathology, University of California Davis, Davis, CA, USA.

***Correspondence:** Boris A. Vinatzer, vinatzer@vt.edu; Tiffany M. Lowe-Power, tlowepower@ucdavis.edu

Keywords: bacterial wilt; core genome; sequevar; species concepts; taxonomy; recombination.

Abbreviations: ANI, average nucleotide identity; CDI, contact-dependent inhibition; GTDB, genome taxonomy database; HGT, horizontal gene transfer; LIN, life identification number; RSSC, *Ralstonia solanacearum* species complex; sequevar, sequence variant group.

Data statement: All supporting data, code and protocols have been provided within the article or through supplementary data files. Four supplementary tables are available with the online version of this article.

000791 © 2022 The Authors



This is an open-access article distributed under the terms of the Creative Commons Attribution License. This article was made open access via a Publish and Read agreement between the Microbiology Society and the corresponding author's institution.

Impact Statement

The *Ralstonia solanacearum* species complex (RSSC) includes dozens of economically important pathogens of many cultivated and wild plants. The extensive genetic and phenotypic diversity that exists within the RSSC has made it challenging to subdivide this group into meaningful subgroups with relevance to plant disease control and plant biosecurity. This study provides a solid genome-based framework for improved classification and identification of the RSSC by analysing 100 representative RSSC genome sequences with a suite of comparative evolutionary genomic tools. The results also lay the foundation for additional in-depth studies to gain further insights into the evolution and biology of this heterogeneous complex of destructive plant pathogens.

description. In contrast, many plant pathogenic bacteria belong to species complexes whose members share phenotypes but have pairwise ANI values below 95%. Further, one of the most important phenotypes for plant pathologists, host range, varies widely among members of the same plant pathogen species.

The bacterial wilt pathogens in the *Ralstonia solanacearum* species complex (RSSC) are a notable example and the objects of this study. RSSC pathogens share a specialized habitat, the water-transporting xylem vessels and stem apoplasts of angiosperm plants, as well as a common pathology, lethal wilt symptoms [3]. Nonetheless, pairwise ANI of RSSC strains can be as low as 90.7%, and host ranges can vary dramatically between closely related strains that have pairwise ANI over 95% [4]. At the same time, many phylogenetically distant strains with pairwise ANI below 95% share host ranges [4, 5].

Genomic analyses place RSSC strains into four statistically supported phylogenetic clades that each share ANI values above 95% and correspond to geographical regions where the clades diversified [6]. These clades are known as phylotypes I, II, III and IV and have geographical origins in Asia, the Americas, Africa and the Indonesian archipelago/Japan, respectively. Phylotype II can be further subdivided into IIA and IIB corresponding to two sub-clades [7]. Taxonomists formally divided the species complex into three species: *R. solanacearum*, corresponding to phylotype II; *R. pseudosolanacearum*, corresponding to phylotypes I and III; and *R. syzygii*, corresponding to phylotype IV [8] (Fig. 1).

Describing the RSSC phylotypes as three named species conforms to taxonomic practice since RSSC clades are separated by genomic metrics and several physiological traits correlate with the clades [6]. On the other hand, one could argue that there are not consistent differences in relevant pathogen behaviour and ecology between clades to justify their division into separate species. Moreover, re-classification using new names leads to inconsistent naming of strains in the literature and in databases. The resulting confusion can interfere with one of the main goals of taxonomy: clear communication about organisms.

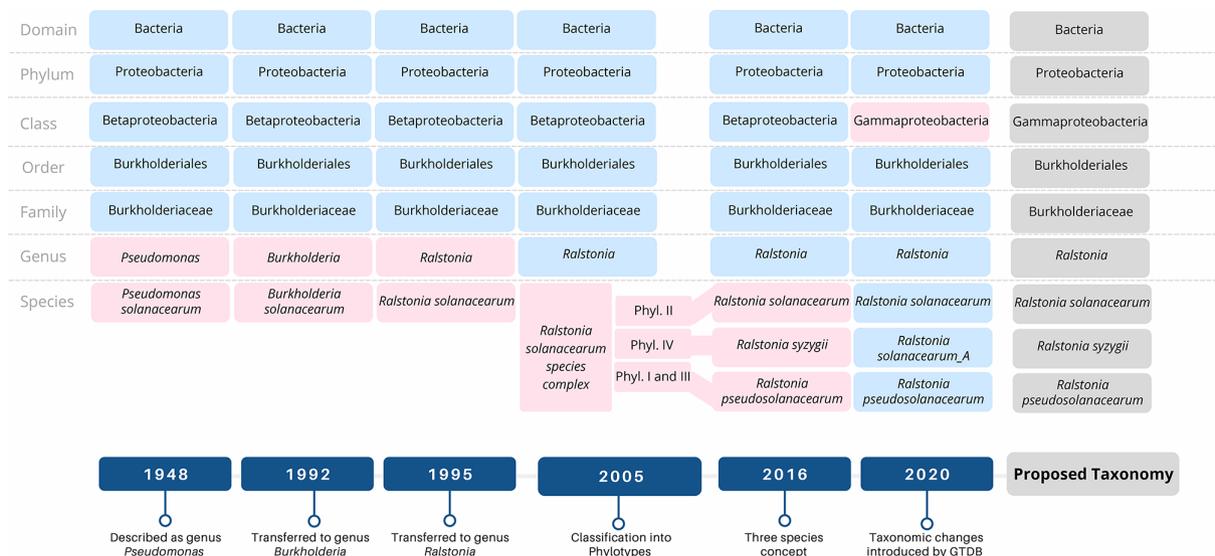


Fig. 1. Major taxonomic revisions for the *Ralstonia solanacearum* species complex (RSSC). The bottom half depicts the timeline when these major changes were introduced, and the top half illustrates the predominant taxonomy used for each era. For each revision, pink boxes highlight changes to the classification, and blue boxes show levels that were unchanged. The taxonomic classification proposed throughout this paper is highlighted in grey.

There is no simple resolution to this conflict. There are almost as many opinions about what a bacterial species is, and if bacterial species even exist, as there are taxonomists [9]. However, in today's taxonomic practice, a pragmatic species 'definition' is used. Bacterial species are commonly defined as groups of bacteria that have over 95% ANI to the name-bearing type strain of one species, have less than 95% ANI to type strains of all other named species, and share a set of measurable phenotypes that distinguish them from members of other named species [2, 10]. Fortunately, genome sequence analysis now allows us to go far beyond ANI to infer many characteristics of groups of bacteria and to circumscribe bacterial species using a variety of species concepts, including the evolutionary, the ecological and the pseudo-biological species concepts.

The evolutionary species concept considers species as independently evolving units [11]. Therefore, the investigation of evolutionary relationships or phylogenetics is the main approach for describing species based on this concept. The economic and technological accessibility of genome sequencing has allowed scientists to replace older approaches, such as DNA–DNA hybridization and 16S rRNA gene sequencing, with phylogenetic reconstructions based on whole genomes. Yet, even using all genes shared by a group of organisms may not precisely reflect their complete evolutionary relationships because of horizontal genetic exchange between sub-lineages [12]. However, it is hard to argue that there is anything that comes closer to representing evolutionary relationships than building a phylogenetic tree based on all gene sequences shared by the genomes under investigation; in other words, building a core genome phylogeny.

In a herculean effort, the genome taxonomy database (GTDB) team has built a phylogeny using protein sequences that roughly correspond to the core genome of all sequenced bacteria [13, 14]. This effort has helped correct incongruencies in the taxonomic lineages of validly published species descriptions, which are often based only on single gene 16S rRNA gene sequences. The names and lineages of these species descriptions can be found in the official List of Prokaryotic Names with Standing in Nomenclature (LPSN), which are reflected in large part by NCBI taxIDs [15]. Each time the GTDB finds a genome that does not belong to a named species because it has a lower than 95% ANI to the type strain of a species, it creates a new species cluster with a placeholder name, e.g. *Escherichia coli*_A. With respect to the RSSC, the GTDB changed a higher rank of the RSSC taxonomy: based on evolutionary distances inferred from genome sequences, the GTDB demoted the *Betaproteobacteria* to a subgroup nested within the class of the *Gammaproteobacteria* [14] (Fig. 1). This shift in RSSC taxonomy was adopted by the microbial community profiling database SILVA with release 138 [16]. Importantly, the GTDB does not resolve evolutionary relationships beyond the 95% ANI threshold (i.e. within species) since its goal is to improve 'traditional' taxonomy based on the established ranks from kingdom to species and not to resolve evolutionary relationships within species.

The sequevar system was developed as a phylogeny-based taxonomy for within-species classification of the RSSC. This system coarsely estimates phylogenetic relationships of strains based on a multiple sequence alignment of a single DNA marker (a 750 bp region of the *egl* endoglucanase gene). Strains with similar sequences are assigned to sequence variant groups (sequevars) [17]. This can be considered a taxonomy focused on the 'Evolutionary within-species concept', with the expectation that some of the predicted relationships are inaccurate due to horizontal gene transfer (HGT). As the plant pathology community transitions from population genetics to population genomics, the ability of the sequevar system to estimate within-species phylogeny can be validated, which is one goal of this paper.

The ecological species concept defines a species as a group of bacteria that adapted to the same ecological niche [18]. Genomic comparisons can also provide insight into ecological species since bacterial adaptation necessarily involves a combination of gene gain/loss and allelic differentiation of gene sequences. For example, a pangenome analysis identifies gene families that are present or absent in different sets of genomes. These genome sets may represent groups that have adapted to different ecological niches and may thus represent different ecological species. Recently, the novel reverse ecology approach has gained traction [19]. This approach aims to identify populations that are in the process of adapting to an ecological niche based on frequent exchange of advantageous mutations during selective sweeps [20]. Putting this concept into practice, Arevalo and colleagues developed a tool, PopCOGenT, that assigns bacteria to distinct populations by identifying recent recombination events within sets of genomes and cessations of recombination between other sets of genomes [21]. Since the reverse ecology approach defines populations based on gene exchange, it also relates to the pseudo-biological species concept [22], which connects bacteria to the biological species concept, which was developed for sexually reproducing eukaryotes. In the pseudo-biological species concept, gene exchange by homologous recombination during sexual crosses is replaced with gene exchange by HGT [23]. For example, the *Pseudomonas syringae* species complex has been proposed to represent a single species because HGT of virulence genes has been found to occur across the entire complex [24].

Because plant pathogenic bacteria with pairwise ANI values above 95% can have starkly distinct host ranges, plant pathologists have developed *ad hoc* within-species classification systems. In most pathogen groups, the 'pathovar' concept is used to describe sub-species groups that cause the same disease on the same range of host plant species [25]. The 'race' system is often used to describe strains within a pathovar that cause disease on different crop genotypes within the same species (e.g. in *Pseudomonas syringae* pv. *phaseolicola* [26]). The RSSC was never divided into pathovars, but for many years the term race was used in an attempt to divide strains by host range at the plant species level. This was never practically useful and eventually the RSSC race system broke down for two reasons. First, RFLP and sequence data revealed the 'races' did not correspond to phylogenetic divisions

[4, 27]. Second, most RSSC strains have very broad host ranges; it is not unusual for one strain to be able to cause disease on monocot and dicot hosts (e.g. banana and tomato [28] or potato and ginger [29]). As a result, most strains end up in a single unhelpful 'Race 1' bin that includes members of all four phylotypes described above. In parallel, the RSSC was also subclassified into biovars based on *in vitro* physiological tests [30]. Once again, these biovars did not correspond to phylogenetic subgroups.

To alleviate the problem with the many different opinions about what should be considered a species, the confusion due to recurrent reclassification and the various within-species classification schemes that are hard to use for non-specialists, we have developed a stable and neutral genome-based framework to circumscribe any of the above groups and to easily translate from one classification system to another. This system is based on genome similarity-based codes, called Life Identification Numbers (LINs) [31]. LINs consist of a series of positions with each position representing a different ANI threshold. ANI thresholds increase moving from the left to the right of a LIN. Therefore, bacteria with very low pairwise ANI do not share any LIN position (below 70% ANI). Bacteria with intermediate ANI (e.g. 95%) have identical LINs to an intermediate position (e.g. position F). Nearly identical bacteria (e.g. 99.99% ANI) have LINs that are identical up to, but not including, the rightmost LIN positions (e.g. position R or S). Therefore, LINs can precisely circumscribe any bacterial group with pairwise ANI values from 70%, corresponding approximately to families and genera, to around 99.99%, corresponding approximately to clonal lineages. In the LINbase Web server, LINs have been implemented for numerous microbial genomes, including the representative genomes of GTDB [32].

The goal of this paper is to investigate RSSC classification through the lens of the different species concepts and within-species concepts by applying comparative evolutionary genomics and a reverse ecological approach to a set of representative, publicly available RSSC genomes. To translate this meta-analysis into applied utility, we then circumscribed the identified groups in the LINbase Web server, so that users can easily identify any new isolate based on its sequenced genome as a member of a named species, phylotype, population or any other group within the RSSC.

METHODS

Selection of representative genomes

All publicly available genomes belonging to the three species (*Ralstonia solanacearum*, *Ralstonia pseudosolanacearum* and *Ralstonia syzygii*) were downloaded from the Assembly database of NCBI on 5 September 2020. Assembled genomes of strain *R. syzygii* R24 and Blood Disease Bacterium R229 were downloaded from the Microscope Microbial Genome Annotation and Analysis Platform – MaGe [33]. All genome assemblies were assessed for quality using the CheckM (version 1.0.13) tool [34]. Genomes with completeness over 98%, contamination below 6%, number of contigs below 670, and N50 scores above 20000 and ambiguous bases below 7% of the genome size were retained. This genome set was further reduced by removing almost identical genomes to obtain a more even representation of the currently known genomic diversity of the RSSC. This was done using the LINflow tool (version 1.1.0.3) [35], retaining only one genome for each group of genomes that had reciprocal ANI values of over 99.975%. Preference was given to genomes of higher sequence quality and for which more published biological data were available. To increase the representation of genomic diversity, the genome of strain 19-3PR_UW348 was sequenced using the Pacbio Sequel II sequencing platform and assembled using Canu (version 2.0) [36] and included in the analysis as well (NCBI accession number JAJMMU000000000). As outgroups, we chose the genome of *Ralstonia mannitolilytica* reference strain SN82F48 and Rsol85. 58_RSOL is one of the strains most closely related to the RSSC without being a member of the RSSC. (We note that it is wrongly classified as *R. solanacearum* in NCBI.)

Pangenome analysis and reconstruction of the core-genome phylogenetic tree

The selected RSSC genomes were subjected to a pangenome analysis using PIRATE (version 1.0.4) [37]. To prepare the genome sequences for input to PIRATE, genomes were annotated using the PROKKA gene annotation tool (version 1.14.6) [38] with default settings. The annotated files were then used to obtain a core gene alignment whereby all genes present in at least 98% of the genomes were considered as core genes. The following parameters were used: -a to obtain a multiFASTA core gene alignment file as output and -k for faster homology searching with the --diamond option specified. The final core gene alignment file was used as input for IQtree (version 2.0.3) [39] using automated model selection to obtain a maximum-likelihood phylogenetic tree. The final phylogenetic tree was visualized using the ggtree [40] package in R. For the pangenome analysis, the PIRATE output file with all gene families was used to obtain the differences in gene content between different phylotypes. For phylotypes I and II, a gene was considered as a core phylotype gene if it was present in more than 95% of the genomes in a phylotype. Because of the much smaller number of genomes in phylotypes III and IV, presence in all but one genome was used as a rule. A score of 1 was assigned for gene presence and a score of 0 for gene absence. This assessment was performed for each gene in the pangenome for all four phylotypes (I, II, III and IV), resulting in a presence-absence matrix with genes as rows and phylotypes as columns (Table S2, available in the online version of this article). The matrix was then visualized through an upset analysis using the UpSetR [41] package in R.

ANI analysis

Pairwise ANI was measured for all representative genomes using *pyani* (version 0.2.10) [42] with default settings. The resulting matrix was used to construct a heatmap of ANI values using the function *heatmap.2* under the *gplots* package [43] in R.

Recombination analysis

First, a recombination analysis of the RSSC was performed within the core genome. The core gene alignment and the phylogenetic tree obtained in the pangenome analysis were used as input to *ClonalframeML* (version 1.12) [44] with default parameters. The inferred recombination regions were used in two different analyses: (1) to find the genes in these regions using *SAMtools* (version 1.12) [45] with the command *intersect*; and (2) to build a recombination-free phylogenetic tree by masking the recombination regions using *cfml-maskrc* [46] and using the new recombination-free alignment as input to *raxml-ng* (version 1.0.3) [47] with the following parameters: `--all --model GTR+G --bs-trees 1000`. The tree was visualized using the *ggtree* [40] package in R.

Next, a recombination analysis was performed separately for each phylotype including the entire genome. For each phylotype, three different reference genomes (four for phylotype II; Table S4) were picked based on the CheckM results. The genome assemblies for each phylotype in FASTA format were obtained from NCBI and used as input to *snippy* (version 4.6.0) [48] with the command *snippy-multi* to generate a whole genome SNP alignment mapped to each of the different reference genomes separately. The whole genome SNP alignment was used as input to *gubbins* (version 3.0.0) [49] to obtain the regions under recombination for each phylotype. The *SAMtools intersect* [45] function was used to find the genes in these regions.

Reverse ecology analysis

To obtain population predictions, inferred from the pairwise measurement of HGT, all of the representative genomes were used as input to *PopCOGenT* (downloaded from <https://github.com/philarevalo/PopCOGenT> in March 2021 [21]).

Sequevar analysis

Automated sequevar assignments were generated using a custom bash script that takes a query genome sequence and compares it to a database of *egl* gene sequences (compiled by E. Wicker, CIRAD, France [50]) using the command line version of the Basic Local Alignment Search Tool: *BLAST* (version 2.9.0+) [51]. Sequevar assignment was made based on the best hit with 99–100% alignment, and results were cross-checked with data from the literature when available.

LIN assignment and LINgroup circumscriptions

All representative genomes and their metadata were uploaded into *LINbase* [32] for automated LIN assignment. *LINgroups* corresponding to groups identified here were circumscribed including a name, a description and a link to the present paper.

RESULTS AND DISCUSSION

A core-genome phylogeny to determine evolutionary relationships

To classify the RSSC based on the evolutionary, ecological and pseudo-biological species concepts, we needed to identify high-quality genome sequences that best represent the described genetic diversity. We started with 167 publicly available genome sequences (Table S1), from which we removed 11 low-quality genomes that were fragmented into many contigs, had low genome completeness scores, had high contamination scores or had a high number of ambiguous bases. From the remaining 156 genomes, we selected 100 genomes (Fig. 2) best representing the known diversity of the species complex and limiting redundancy due to several nearly identical genomes present in the original set.

To uncover the phylogenetic relationships among the representative strains, we performed a pangenome analysis. This analysis revealed that 3262 orthologous genes constitute the RSSC core genome (Table S2). Previous attempts to determine the RSSC core genome using fewer genomes yielded smaller estimates of 2370 [52] and 1940 [4], respectively. This goes against the well-established trend that core genome size decreases as the number of included genomes increases [53]. This unexpected result is probably because we used the pangenome software *PIRATE* [37], which was designed to recognize members of orthologous groups even in the presence of low sequence identity as is the case with a species complex such as the RSSC. Regardless, bioinformatic studies are inherently approximations and definitive identification of orthologous genes will depend on gene by gene functional analyses such as trans-complementation experiments.

A phylogenetic tree based on these core genes (Fig. 2) clustered strains into clades corresponding to the four known phylotypes, with 59 strains belonging to phylotype I; 28 strains belonging to phylotype II (among which nine and 16 strains were in phylotypes IIA and IIB, respectively, and three strains were in a cluster that was distinct from either IIA or IIB); five strains belonging to phylotype III; and eight strains belonging to phylotype IV. During this analysis, we identified one genome sequence that may be the result of a chimeric assembly between a phylotype I strain and a phylotype II strain: CRMRS218. This genome was published

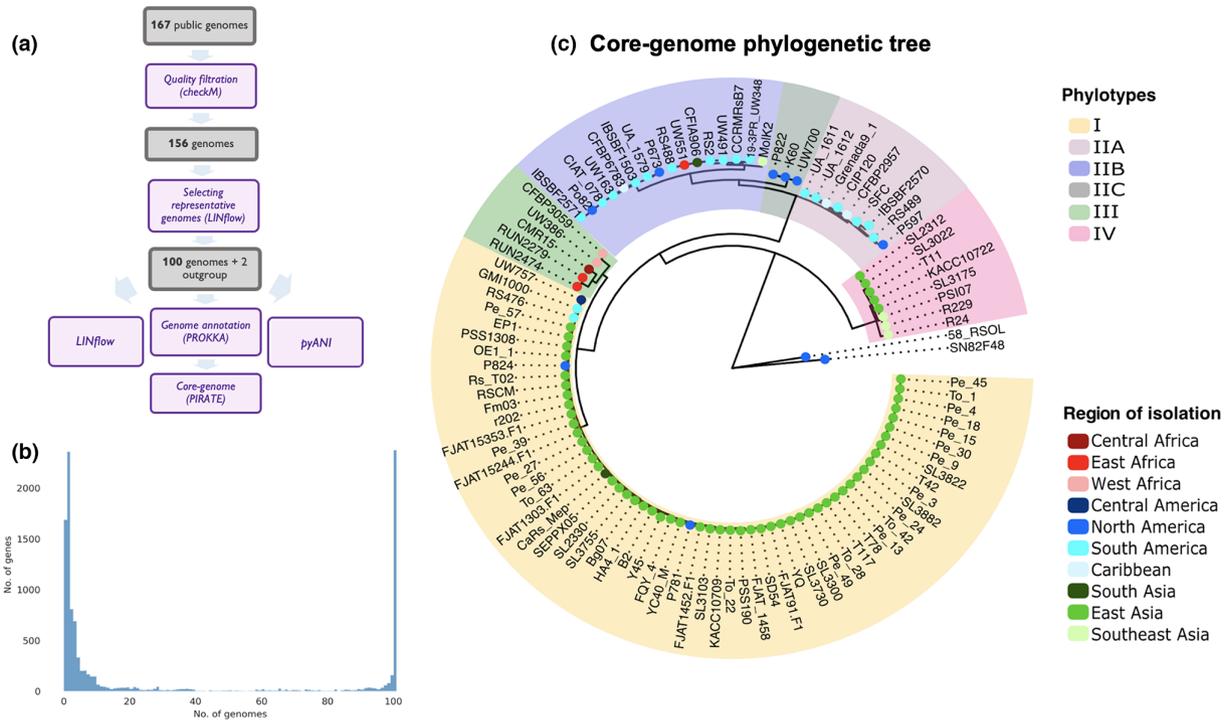


Fig. 2. Core genome analysis for the representative genomes of the RSSC. (a) Selection of the representative genomes. Purple boxes indicate the software used, and the grey boxes show the number of genomes left at each step. (b) The number of genomes that carry each gene in the pangenome. (c) Phylogenetic tree obtained with the core-genome analysis. All clades with high bootstrap values are included in the tree. Phylotypes of the strains are highlighted in different colours representing phylotypes I, IIA, IIB, III and IV. Based on the analysis, strains P822, K60 and UW700 are classified as phylotype IIC. Coloured dots at the node of each strain represent the region of isolation.

as a phylotype I strain [54], but in the core genome tree it formed a singleton branch basal to all phylotype II strains. Because of this ambiguity, the strain was excluded from further analysis.

Based on the geographical origin of strains, the phylogenetic tree is consistent with the hypothesis that the phylotypes diversified in different global regions [5, 55]. In fact, most phylotype I strains were isolated in continental Asia, phylotype II strains in the Americas, phylotype III strains in Africa, and phylotype IV strains in Indonesia and the Pacific Islands (Fig. 2). It is important to point out that the strains used here are not equally distributed between and within continents and thus the phylotypes are not equally distributed. For example, strains belonging to phylotype III isolated in Africa are underrepresented (5% of total strains) compared to other phylotypes. East Asian strains represent 90% of the analysed phylotype I strains, with most sequenced strains isolated in either South Korea or China. Although phylotype I is common in South Asia, only 1.7% of the sequenced phylotype I strains were isolated in South Asia. This uneven representation probably reflects a bias in publicly available genome sequences from different geographical regions and does not reflect the actual geographical distribution and diversity of RSSC strains.

The phylotype II circumscription was consistent with the classification of strains based on the LPSN and GTDB classification systems of belonging to the named species *R. solanacearum*. Similarly, all phylotype I and III strains were consistent with the LPSN and GTDB classification of belonging to the recently named species *R. pseudosolanacearum* [8]. Phylotype IV strains correspond to *R. syzygii* as per the LPSN taxonomy and '*R. solanacearum_A*' as per the GTDB. It is important to note that many strains that are members of *R. syzygii* and *R. pseudosolanacearum* are listed as *R. solanacearum* in NCBI, because the genomes were submitted before the reclassification and adoption of the new species names by the scientific community.

Pangenome analyses provide a basis to investigate adaptation to ecological niches

One of the currently unanswered questions about the RSSC is to what degree the four phylotypes diverged from each other because of adaptation to different niches or because of allopatry. As a small step towards answering this question, we determined the congruences and differences in gene content between and within phylotypes.

Overall, the RSSC contained a total of 13128 gene families, which represent the RSSC pangenome. This is a considerably smaller pangenome size estimation compared to a previous study that found the RSSC pangenome to include 16757 genes based on 19 genomes [4]. This can again be explained by the pangenome software we used, PIRATE [37], which was designed to recognize

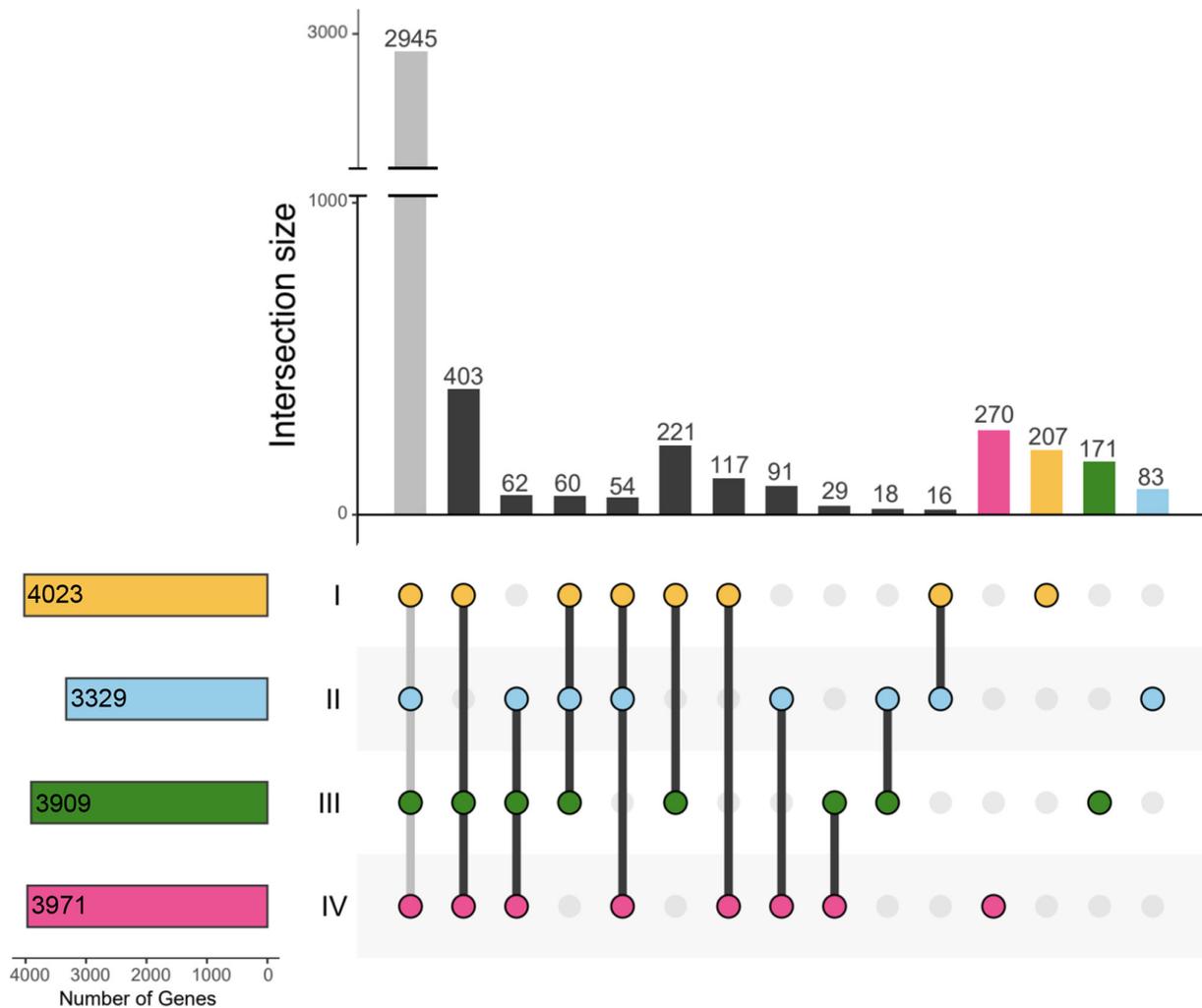


Fig. 3. Pangenome analysis represented using an Upset plot to highlight how many genes are shared between phylotypes I, II, III and IV. Each bar on the vertical bar chart represents the number of genes shared by the combination of phylotypes shown below the chart. The horizontal bar chart indicates the size of the phylotype-core genomes.

members of orthologue groups across diverse genomes and, therefore, led to a smaller number of gene families that each contain a larger number of member genes.

The respective pangenome sizes of the individual phylotypes are: 4023 (I), 3329 (II), 3909 (III) and 3971 (IV). An Upset plot was used to visualize the number of genes that are either shared by all strains of one phylotype and absent from all other phylotypes, i.e. the phylotype-specific core genes, or that are shared between subsets of phylotypes (Fig. 3). Due to the above-mentioned differences in the extent to which the genomic diversity within each phylotype was sampled, it is difficult to draw firm conclusions. Based on the available data, the core genome of phylotype II (3329 genes) was considerably smaller than those of phylotypes I, III and IV (3909–4023 genes).

At the species level, *R. solanacearum* (phylotype II) has a core genome size (3329 genes) very similar to the core genome size of *R. pseudosolanacearum* (phylotypes I and III combined) (3408). A surprising finding is the large core genome size of the species *R. syzygii*, which includes strains that cause the most phenotypically diverse diseases (Sumatra disease of cloves, banana blood disease and classical bacterial wilts) [56]. However, the large size of the *R. syzygii*/phylotype IV core genome (3971) could be an artefact due to the small number of phylotype IV genomes available and the use of a less stringent core-genome cutoff (presence in all but one of the eight genomes, i.e. 87.5%).

When comparing gene content between phylotypes, phylotypes I and III share the most core genes with each other that are not core genes of the other phylotypes (221 genes). This is consistent with the shared membership of phylotypes I and III in the species *R. pseudosolanacearum*. Phylotypes I, III and IV constitute the group of phylotypes that have the most genes in common

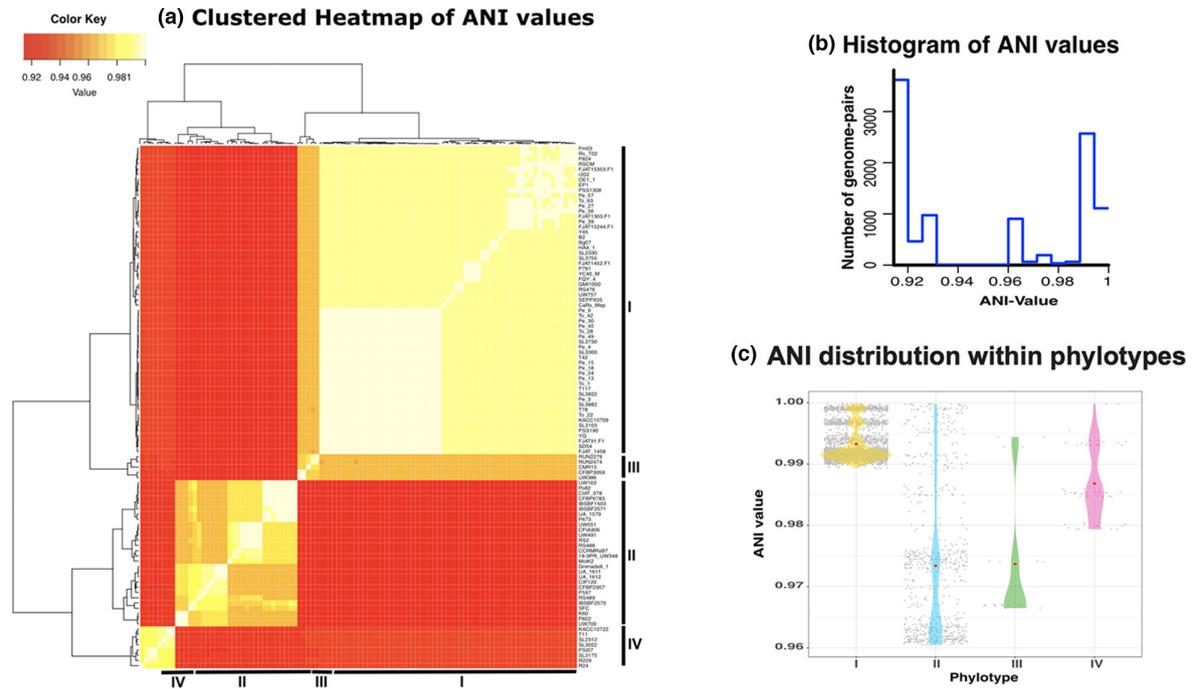


Fig. 4. Average nucleotide identity (ANI) analysis for representative RSSC genomes. (a) Heatmap of pairwise ANI values for all genomes. (b) Histogram of pairwise ANI values among all paired genome combinations. (c) Pairwise ANI distribution within each phylotype. Grey dots represent pairwise ANI between genomes belonging to the same phylotype, and red dots show the mean ANI for each phylotype.

that are absent from the core genome of the remaining phylotype, i.e. phylotype II in this case (403 genes). This is consistent with phylotype II having the smallest core genome and being the most diverse phylotype with regard to gene content.

ANI analysis confirms species boundaries and genome similarity-based clusters

After determining phylogenetic relationships and comparing gene content between strains to provide the basis for investigating the RSSC from an evolutionary and ecological perspective, we calculated pairwise ANI between all 100 genomes (Fig. 4, Table S3). Since ANI is based on the average genetic distance of all DNA sequences shared between pairs of strains, it provides an orthogonal measure of genomic relationships beyond a core genome tree, which is limited to the genes shared by all 100 strains. In agreement with the core genome analysis, pairwise ANI clustered the genomes into the four phylotypes. Importantly, although phylotypes I and III formed distinct clusters, all strains in these two phylotypes had pairwise ANI values above 95%, which is consistent with these phylotypes being part of the same species.

Phylotype I strains had higher average pairwise ANI (99.35%) than other phylotypes (97.73% for phylotype II, 97.30% for phylotype III and 98.67% for phylotype IV). Phylotype I appears to be the most genetically homogenous phylotype, but, as pointed out above, the genomic similarity could be an artefact stemming from the limited geographical distribution of most phylotype I genomes. If the high ANI among phylotype I strains is maintained as more South Asian strains are sequenced, this may indicate that phylotype I emerged more recently in evolutionary time, possibly from within the wider genetic diversity of phylotype III.

Strains within phylotype II are characterized by relatively low ANI. Pairwise ANI indicates that there are three main subgroups. Strains in the sequevar 7 clade (type strain K60, UW700, P822) had high pairwise ANI with each other (mean ANI 99.73%) and lower ANI with IIA and IIB strains (mean ANI 97.53 and 96.18%, respectively), which is consistent with sequevar 7 strains clustering as a sub-clade separate from phylotypes IIA and IIB.

Recombination analyses provide a basis to identify biological and ecological species

Most RSSC strains are naturally transformable [57], and prior population genetics and genomics studies at the global, regional and field scales have indicated that RSSC genomes are highly recombinogenic [55, 58–60]. To investigate whether the core genome phylogenetic tree was biased by recombination within the RSSC, we used ClonalFrameML to identify core genes that lack evidence of recombination. ClonalFrameML found recombination regions in 1559 core genes (Table S4). The recombination regions detected by ClonalFrameML were masked and a recombination-free tree is shown in Fig. 5 (tree on the right). While this tree maintained the main clades from the core genome tree shown in Figs. 2 and 5 (tree on the left), the Southeast US clade

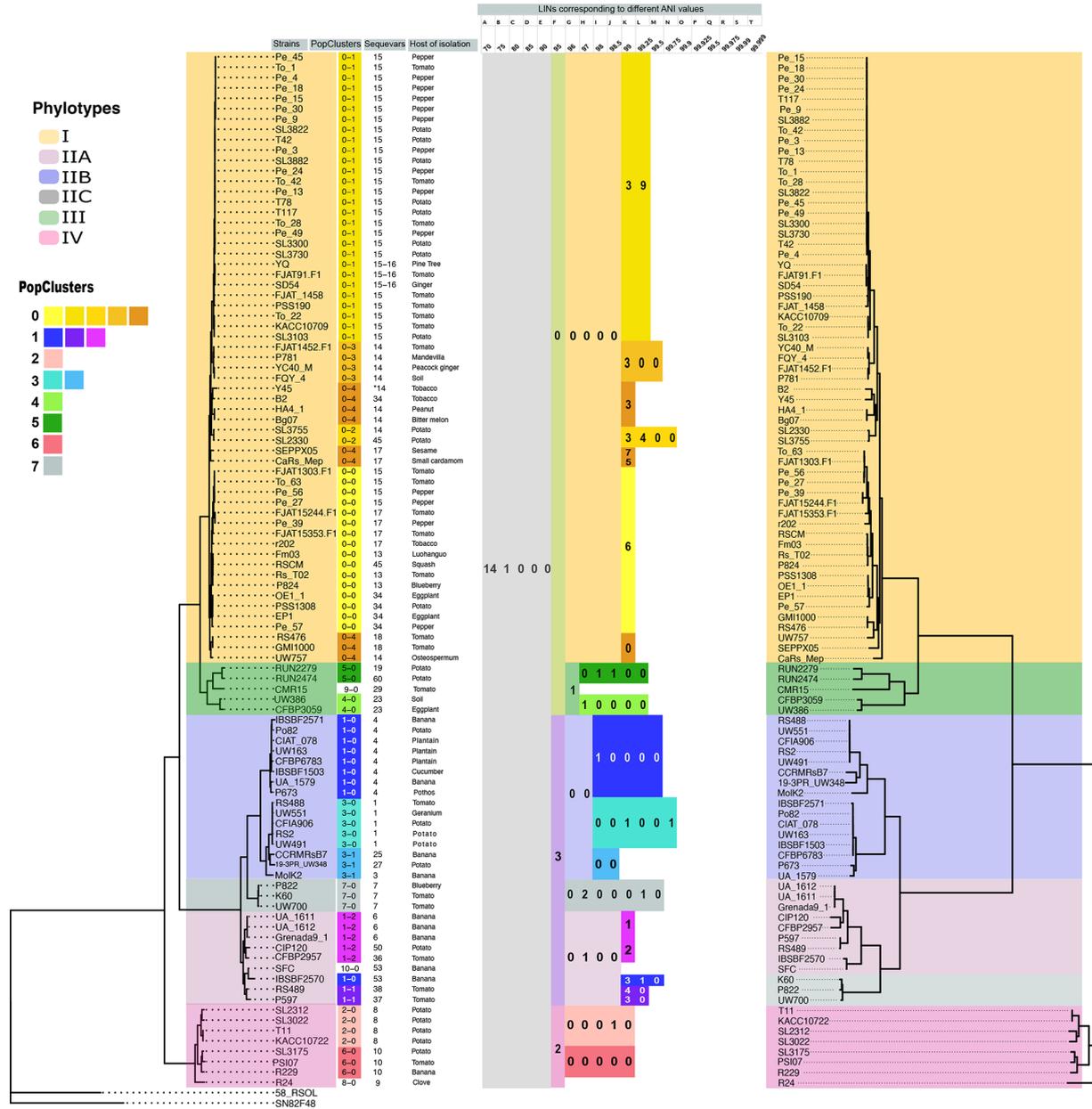


Fig. 5. Comparison of core-genome tree, recombination-free tree, population clusters, sequevar types and delineation of RSC groups using LINs. The tree on the left is a vertical version of the core-genome phylogenetic tree from Fig. 2. To the right of each strain name are assignments to population clusters, sequevars and then the respective hosts of isolation. LINs corresponding to each group (the RSC, named species, phylotypes, sub-phylotypes and population clusters) are listed using colours matching each group. Newly sequenced genomes can be identified as members of these groups at

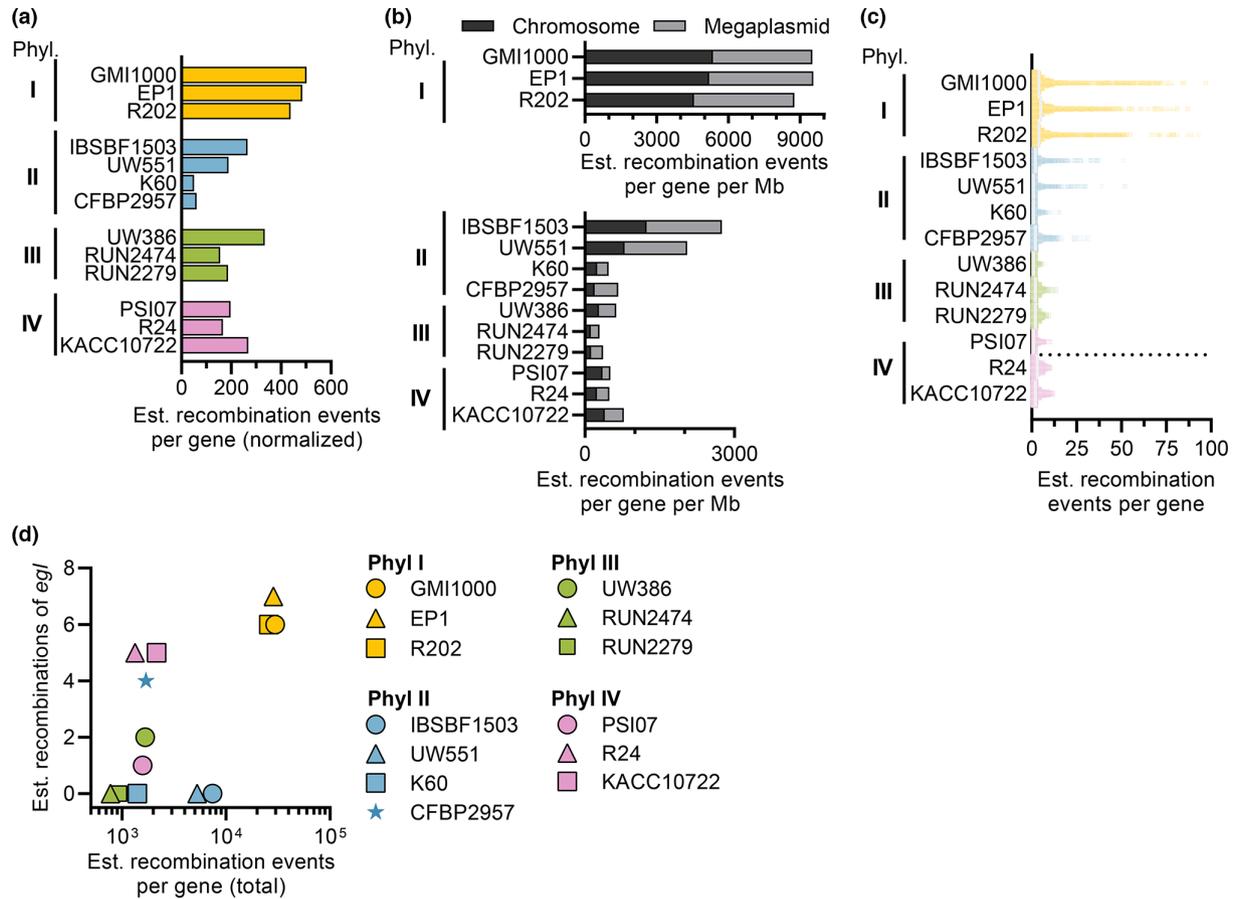


Fig. 6. Comparison of estimated recombination for representative RSSC genomes from each phylotype. Genes with putative recombination events were identified using Gubbins [49]. (a) The number of recombination events for each genome, normalized by the number of genomes of each phylotype in the genome set. (b) The number of recombination events on the chromosome vs. megaplasmid, normalized by the length of the replicon. (c) Estimated number of recombination events detected for each gene (dots). (d) Comparison of the number of recombination events for the sequevar marker gene (*egl*) vs. the total recombination events for each genome.

of under-sampled populations. However, there were two PopCOGenT clusters that were polyphyletic: PopCluster 1–0 contained eight IIB-4 strains and a IIA-53 strain IBSBF2570, while PopCluster 0–4 contained nine phylotype I strains from three distinct branches on the core genome tree.

Genes that are frequently transmitted horizontally between strains may play a role in adaptation (the ecological species concept). Therefore, in addition to PopCOGenT, we ran the independent recombination tool Gubbins [49] to detect recombination in the RSSC using 13 reference genomes (three genomes for phylotypes I, III and IV and four genomes for phylotype II). The results are summarized in Fig. 6. Table S4 contains the estimated number of recombinations for each gene in the 13 reference genomes. As expected, mobile genetic elements (transposases, integrases and phage-associated proteins) were highly recombinogenic genes. Many of the highly recombining genes are type III secreted effectors, which RSSC strains use to manipulate plant host physiology and immunity. The high plasticity of type III effector repertoires is well known in RSSC strains [61]. Additionally, glycoside hydrolases, polygalacturonases and endoglucanases displayed evidence of frequent recombination. Endoglucanases are involved in adaptation of *Xanthomonas* species to vascular vs. apoplastic niches [62], but variation in plant cell-wall-degrading enzyme repertoires has not been investigated for the RSSC. Several classes of genes potentially involved in inter-microbial interactions were recombinogenic: type VI secretion system genes such as Vgr, PAAR and putative effector/immunity pairs [63], and haemagglutinin-like proteins that are hypothesized to be contact-dependent inhibition (CDI) systems in the RSSC [60]. Another group of recombinogenic genes encode non-ribosomal peptide synthetases and polyketide synthases that can synthesize secondary metabolites involved in intermicrobial competition, among other functions [64]. Investigating the functional diversity of the recombining genes may shed light on how interactions with plant hosts, microbial competitors and novel abiotic environments shape the evolution of RSSC lineages.

Speculation on the relative evolutionary ages of phylotypes

Overall, our comparative genomics analyses suggest that either phylotype II (*R. solanacearum*) or phylotype IV is the most ancestral phylotype within the RSSC. Phylotype II genomes have the lowest average pairwise ANI value and phylotype II has the smallest core genome. Their lower recombination rate is also in line with higher sequence diversity since higher sequence diversity decreases the success of homologous recombination. All these results suggest that phylotype II is more diverse compared to the other phylotypes and, thus, could have emerged first. These findings are also consistent with an earlier study in which 29 RSSC genomes and 73 MALDI proteomes were compared [6]. Surprisingly, however, phylotype IV is on the most basal branch in the core genome tree (Fig. 2), as it was in a previous multi-locus sequence analysis tree [55]. This suggests that phylotype IV is the most ancestral phylotype. This inconsistency could be due to uneven sampling among phylotypes. The genomic diversity in phylotype IV may be under-sampled, and if additional genomes of diverse phylotype IV strains were to be sequenced, it might become more diverse than phylotype II. On the other hand, the basal position of phylotype IV might have been influenced by the choice of outgroup strains. If phylotype IV strains acquired genes from environmental *Ralstonia* closely related to the chosen outgroup strains, recombination could make phylotype IV seem more closely related to the outgroup strains than they are by vertical inheritance. Therefore, we cannot firmly conclude which phylotype is most ancestral based on available data. On the other hand, there is one clear interpretation about relative ages of the phylotypes. Phylotype I is the least diverse phylotype that also branches off the latest as a lineage from phylotype III, making it probably the phylotype that most recently emerged and expanded.

Comparing sequevars (*egl* trees) with the core genome phylogeny and populations

The global plant pathology community has widely adopted the sequevar taxonomic system to classify plant pathogenic *Ralstonia* strains at the within-species level. Over 5000 strains from more than 88 regions have been assigned to over 70 sequevar groups [65]. Because the sequevar system is based on a single genetic marker (750bp of the *egl* gene), and RSSC genomes often recombine, we predicted that the *egl* gene may have recombined between strains. Indeed, *egl* recombination events were detected in 3/3 phylotype I, 1/4 phylotype II, 1/3 phylotype III and 3/3 phylotype IV reference genomes used in the Gubbins analysis (Fig. 6d, Table S4). We and other plant pathologists have deposited over 4500 '*egl* gene, partial cds' sequences from RSSC isolates to the NCBI nucleotide database, but our results suggest that recombination of *egl* within the RSSC may limit the sequevar taxonomy's ability to accurately estimate phylogenetic relationships.

With evidence that *egl* may be horizontally transmitted between RSSC strains, we investigated whether the sequevar system and trees reconstructed with *egl* sequences reflect phylogenetic relationships of strains. We extracted the partial *egl* nucleotide sequences from each of the 100 RSSC genomes and aligned them with reference sequences to assign sequevars to each genome (Table S1, Fig. 5). The sequevar assignments were monophyletic in the tested genomes for phylotypes II (28 genomes assigned to 12 sequevars), III (five genomes assigned to four sequevars) and IV (eight genomes assigned to three sequevars). Sequevar I-18 and sequevar I-13 mapped to single branches of the tree, so these sequevars may be monophyletic. However, most phylotype I sequevars were highly polyphyletic. Five of the phylotype I sequevars (I-14, I-15, I-17, I-34 and I-45) were assigned to distinct branches within phylotype I.

Overall, our results and previous work [58] indicate that the sequevar system is not informative for describing within-species relationships for phylotype I strains. The polyphyletic phylotype I sequevars are probably due to the inter-related phenomena of phylotype I's low genetic diversity and higher recombination. This suggests that improved methods are needed to classify within-species groups of phylotype I. PCR assays that target insertions/deletions might be a cost-effective method to prioritize strains for whole-genome sequencing [66]. On the other hand, the sequevar system appears to robustly reflect phylogenetic relationships for the diverse phylotype II strains. As more phylotype III and phylotype IV genomes become available, it will be useful to test whether the sequevar system continues to work well in these phylotypes.

Using LINs to circumscribe RSSC groups for easy genome-based identification

In the LIN system, genomes are classified based on genome similarity without deciding on any a priori group boundaries. LINs can thus be used to circumscribe species complexes, species or within-species groups, and place any genome within these groups. If the breadth of a taxon is defined based on an ANI distance from the type strain, this can be done based on the LIN assigned to the type strain. For example, K60 is the type strain of *R. solanacearum*, and the LIN of K60 up to the F position (corresponding to 95% ANI) in the LINbase web server is $14_A1_B0_C0_D0_E3_F$. Therefore, the LIN of the *R. solanacearum* species is $14_A1_B0_C0_D0_E3_F$, and each genome that has the same LIN at these positions can be immediately identified as a member of the species *R. solanacearum*. As shown in Fig. 5, the LIN for *R. pseudosolanacearum* is $14_A1_B0_C0_D0_E0_F$, and the LIN for *R. syzygii* is $14_A1_B0_C0_D0_E2_F$.

If a type strain genome is not available for a group or a group does not have a predetermined ANI breadth (because it is not a species), the group can still be circumscribed based on the LIN positions shared by its members. Since we added the 100 RSSC genomes used in this study to the LINbase web server and assigned LINs to each of them, we were also able to circumscribe the RSSC and its phylotypes, sub-phylotypes, and population clusters so that any newly sequenced genome can be identified not only as a member of a species but also as a member of any of these other groups. In Fig. 6, we report the LINs corresponding to

each of these groups. While the LINs assigned to each individual genome are not shown in the figure, they are stored in Table S1 and in LINbase and can be used to circumscribe even more highly resolved groups corresponding to individual genetic lineages within the RSSC. Whole genome-based LINs could thus be used to replace the single marker gene-based sequevar system, which we have shown to contradict core genome phylogeny for phylotype I.

CONCLUSION

We have shown how a genomic meta-analysis can be used to classify the RSSC according to the evolutionary, biological and ecological species concepts. We circumscribed validly published named species, phylotypes, clades within phylotypes, sequevars (when possible) and populations. We determined how extensively genes are shared within and between phylotypes and which genes most frequently recombine. This work also provided the basis for further in-depth investigations of the RSSC. LINbase makes it straightforward to circumscribe any additional groups based on additional sampling and genome sequencing of the diversity within the RSSC and additional genomic comparisons and phenotypic tests. Any new isolate with a draft genome sequence can then be precisely identified as a member of any of these groups to help inform basic research, disease management and biosecurity regulations.

Funding information

Funding to Boris A. Vinatzer, Caitilyn Allen and Lenwood S. Heath was provided by USDA APHIS (contract AP19PPQS and T00C083). Funding to Boris A. Vinatzer and Lenwood S. Heath was also provided by NSF (DBI-2018522). Funding to Boris A. Vinatzer was also provided in part by the Virginia Agricultural Experiment Station and the Hatch Program of USDA NIFA. Caitilyn Allen was funded by U. Wisconsin–Madison College of Agricultural and Life Sciences. Tiffany M. Lowe-Power was funded by USDA NIFA (grant no. 2022-67013-36272) and UC Davis College of Agricultural and Environmental Sciences and Department of Plant Pathology (laboratory start-up funds).

Acknowledgements

We thank our colleague Emmanuel Wicker (CIRAD, France) for providing reference *egl* sequences and Noah A. Kinscherf, Jessica L. Prom and Alicia N. Truchon for genomic DNA extractions at UW–Madison. Additionally, we thank Stéphane Poussier (University of La Réunion) and Jonathan Jacobs (The Ohio State University) for helpful discussions about sequevar typing of phylotype I strains.

Conflicts of interest

Life Identification Number and LIN are registered trademarks of This Genomic Life, Inc. Lenwood S. Heath and Boris A. Vinatzer report in accordance with Virginia Tech policies and procedures and their ethical obligation as researchers, that they have a financial interest in the company This Genomic Life, Inc., that may be affected by the research reported in this paper. They have disclosed those interests fully to Virginia Tech, and they have in place an approved plan for managing any potential conflicts arising from this relationship.

References

- Sharma P, Flores MA, Mazloom R, Allen C, Health L, et al. Meta analysis of the *Ralstonia solanacearum* species complex (RSSC) based on comparative evolutionary genomics and reverse ecology. *Figshare* 2022.
- Konstantinidis KT, Ramette A, Tiedje JM. The bacterial species definition in the genomic era. *Philos Trans R Soc Lond B Biol Sci* 2006;361:1929–1940.
- Lowe-Power TM, Khokhani D, Allen C. How *Ralstonia solanacearum* exploits and thrives in the flowing plant xylem environment. *Trends Microbiol* 2018;26:929–942.
- Ailloud F, Lowe T, Cellier G, Roche D, Allen C, et al. Comparative genomic analysis of *Ralstonia solanacearum* reveals candidate genes for host specificity. *BMC Genomics* 2015;16:270.
- Remenant B, Coupat-Goutaland B, Guidot A, Cellier G, Wicker E, et al. Genomes of three tomato pathogens within the *Ralstonia solanacearum* species complex reveal significant evolutionary divergence. *BMC Genomics* 2010;11:379.
- Prior P, Ailloud F, Dalsing BL, Remenant B, Sanchez B, et al. Genomic and proteomic evidence supporting the division of the plant pathogen *Ralstonia solanacearum* into three species. *BMC Genomics* 2016;17:90.
- Poussier S, Prior P, Luisetti J, Hayward C, Fegan M. Partial sequencing of the *hrpB* and endoglucanase genes confirms and expands the known diversity within the *Ralstonia solanacearum* species complex. *Syst Appl Microbiol* 2000;23:479–486.
- Safni I, Cleenwerck I, De Vos P, Fegan M, Sly L, et al. Polyphasic taxonomic revision of the *Ralstonia solanacearum* species complex: proposal to emend the descriptions of *Ralstonia solanacearum* and *Ralstonia syzygii* and reclassify current *R. syzygii* strains as *Ralstonia syzygii* subsp. *syzygii* subsp. nov., *R. solanacearum* phylotype IV strains as *Ralstonia syzygii* subsp. *indonesiensis* subsp. nov., banana blood disease bacterium strains as *Ralstonia syzygii* subsp. *celesbesensis* subsp. nov. and *R. solanacearum* phylotype I and III strains as *Ralstonia pseudosolanacearum* sp. nov. *Int J Syst Evol Microbiol* 2014;64:3087–3103.
- Roselló-Mora R, Amann R. The species concept for prokaryotes. *FEMS Microbiol Rev* 2001;25:39–67.
- Stackebrandt E, Frederiksen W, Garrity GM, Grimont PAD, Kämpfer P, et al. Report of the *ad hoc* committee for the re-evaluation of the species definition in bacteriology. *Int J Syst Evol Microbiol* 2002;52:1043–1047.
- Hull DL. The ideal species concept - and why we can't get it. In: Claridge MF, Dawah HA and Wilson MR (eds). *Species: The Units of Biodiversity*. London: Chapman and Hall; . pp. 357–380.
- Stott CM, Bobay L-M. Impact of homologous recombination on core genome phylogenies. *BMC Genomics* 2020;21:829.
- Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, et al. A complete domain-to-species taxonomy for Bacteria and Archaea. *Nat Biotechnol* 2020;38:1079–1086.
- Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, et al. A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat Biotechnol* 2018;36:996–1004.
- Schoch C. NCBI Taxonomy. National Center for Biotechnology Information (US); 2020. <https://www.ncbi.nlm.nih.gov/books/NBK53758/>
- Yilmaz P, Parfrey LW, Yarza P, Gerken J, Pruesse E, et al. The SILVA and “All-species Living Tree Project (LTP)” taxonomic frameworks. *Nucleic Acids Res* 2014;42:D643-8.

17. Fegan M, Prior P. How complex is the *Ralstonia solanacearum* species complex. *Bacterial wilt disease and the Ralstonia solanacearum species complex* 2005;1:449–461.
18. Andersson L. The driving force: species concepts and ecology. *TAXON* 2019;39:375–382.
19. Vos M. A species concept for bacteria based on adaptive divergence. *Trends Microbiol* 2011;19:1–7.
20. Arevalo P, VanInsberghe D, Polz MF. A reverse ecology framework for bacteria and archaea. In: Polz MF and Rajora OP (eds). *Population Genomics: Microorganisms*. Cham: Springer International Publishing; pp. 77–96.
21. Arevalo P, VanInsberghe D, Elsherbini J, Gore J, Polz MF. A reverse ecology approach based on a biological definition of microbial populations. *Cell* 2019;178:820–834.
22. Staley JT. The bacterial species dilemma and the genomic-phylogenetic species concept. *Philos Trans R Soc Lond B Biol Sci* 2006;361:1899–1909.
23. Bobay L-M, Ochman H. Biological species are universal across Life's domains. *Genome Biol Evol* 2017.
24. Dillon MM, Thakur S, Almeida RND, Wang PW, Weir BS, et al. Recombination of ecologically and evolutionarily significant loci maintains genetic cohesion in the *Pseudomonas syringae* species complex. *Genome Biol* 2019;20:3.
25. Young JM, Takikawa Y, Gardan L, Stead DE. Changing concepts in the taxonomy of plant pathogenic bacteria. *Annu Rev Phytopathol* 1992;30:67–105.
26. Arnold DL, Lovell HC, Jackson RW, Mansfield JW. *Pseudomonas syringae* pv. *phaseolicola*: from “has bean” to supermodel. *Mol Plant Pathol* 2011;12:617–627.
27. Cook D. Genetic diversity of *Pseudomonas solanacearum*: detection of restriction fragment length polymorphisms with DNA probes that specify virulence and the hypersensitive response. *MPMI* 1989;2:113.
28. Albuquerque GMR, Santos LA, Felix KCS, Rollemberg CL, Silva AMF, et al. Moko disease-causing strains of *Ralstonia solanacearum* from Brazil extend known diversity in paraphyletic phylotype II. *Phytopathology* 2014;104:1175–1182.
29. Xu J, Pan ZC, Prior P, Xu JS, Zhang Z, et al. Genetic diversity of *Ralstonia solanacearum* strains from China. *Eur J Plant Pathol* 2009;125:641–653.
30. Hayward AC. Characteristics of *Pseudomonas solanacearum*. *J Appl Bacteriol* 1964;27:265–277.
31. Vinatzer BA, Weisberg AJ, Monteil CL, Elmarakeby HA, Sheppard SK, et al. A proposal for a genome similarity-based taxonomy for plant-pathogenic bacteria that is sufficiently precise to reflect phylogeny, host range, and outbreak affiliation applied to *Pseudomonas syringae* *sensu lato* as a proof of concept. *Phytopathology* 2017;107:18–28.
32. Tian L, Huang C, Mazloom R, Heath LS, Vinatzer BA. LINbase: a web server for genome-based identification of prokaryotes as members of crowdsourced taxa. *Nucleic Acids Res* 2020;48:W529–W537.
33. Vallenet D, Engelen S, Mornico D, Cruveiller S, Fleury L, et al. MicroScope: a platform for microbial genome annotation and comparative genomics. *Database (Oxford)* 2009;2009:bap021.
34. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res* 2015;25:1043–1055.
35. Tian L, Mazloom R, Heath LS, Vinatzer BA. LINflow: a computational pipeline that combines an alignment-free with an alignment-based method to accelerate generation of similarity matrices for prokaryotic genomes. *PeerJ* 2021;9:e10906.
36. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, et al. Canu: scalable and accurate long-read assembly via adaptive *k*-mer weighting and repeat separation. *Genome Res* 2017;27:722–736.
37. Bayliss SC, Thorpe HA, Coyle NM, Sheppard SK, Feil EJ. PIRATE: A fast and scalable pangenomics toolbox for clustering diverged orthologues in bacteria. *Gigascience* 2019;8:giz119.
38. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 2014;30:2068–2069.
39. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, et al. IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era. *Mol Biol Evol* 2020;37:1530–1534.
40. Yu G, Smith DK, Zhu H, Guan Y, Lam T-Y. GGTREE: An R package for visualization and annotation of phylogenetic trees with their covariates and other associated data. *Methods Ecol Evol* 2017;8:28–36.
41. Conway JR, Lex A, Gehlenborg N. UpSetR: an R package for the visualization of intersecting sets and their properties. *Bioinformatics* 2017;33:2938–2940.
42. Pritchard L, Glover RH, Humphris S, Elphinstone JG, Toth IK. Genomics and taxonomy in diagnostics for food security: soft-rotting enterobacterial plant pathogens. *Anal Methods* 2016;8:12–24.
43. Warnes GR, Bolker B, Bonebakker L, Gentleman R, Liaw WHA, et al. (n.d.) Gplots: various R programming tools for plotting data. R package version 2.17.0. computer software.
44. Didelot X, Wilson DJ. ClonalFrameML: efficient inference of recombination in whole bacterial genomes. *PLoS Comput Biol* 2015;11:e1004041.
45. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, et al. Twelve years of SAMtools and BCFtools. *Gigascience* 2021;10:giab008.
46. Kwong J. cfml-maskrc: Masks recombinant regions in an alignment based on ClonalFrameML output; 2021. <https://github.com/kwongj/cfml-maskrc>
47. Kozlov AM, Darriba D, Flouri T, Morel B, Stamatakis A. RAxML-NG: a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* 2019;35:4453–4455.
48. Seemann T. snippy: Rapid haploid variant calling and core genome alignment; 2021. <https://github.com/tseemann/snippy>
49. Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, et al. Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res* 2015;43:e15.
50. Wicker E, N'guessan C, Le Roux-Nio AC, Deberdt P, Sujeeun L, et al. A reference database of *Ralstonia solanacearum* egl-mutS haplotypes for global epidemiological surveillance of bacterial wilts; https://agritrop.cirad.fr/582579/1/Wicker_BD%20egl-mutS_FINAL.pdf
51. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, et al. BLAST+: architecture and applications. *BMC Bioinformatics* 2009;10:421.
52. Bocsanczy AM, Huguet-Tapia JC, Norman DJ. Comparative genomics of *Ralstonia solanacearum* identifies candidate genes associated with cool virulence. *Front Plant Sci* 2017;8:1565.
53. Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, et al. Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the microbial “pan-genome.” *Proc Natl Acad Sci U S A* 2005;102:13950–13955.
54. Albuquerque GMR, Souza EB, Silva AMF, Lopes CA, Boiteux LS, et al. Genome sequence of *Ralstonia pseudosolanacearum* strains with compatible and incompatible interactions with the major tomato resistance source Hawaii 7996. *Genome Announc* 2017;5:e00982-17.
55. Wicker E, Lefeuvre P, de Cambiaire J-C, Lemaire C, Poussier S, et al. Contrasting recombination patterns and demographic histories of the plant pathogen *Ralstonia solanacearum* inferred from MLSA. *ISME J* 2012;6:961–974.
56. Safni I, Subandiyah S, Fegan M. Ecology, epidemiology and disease management of *Ralstonia syzygii* in Indonesia. *Front Microbiol* 2018;9:419.
57. Coupat B, Chaumeille-Dole F, Fall S, Prior P, Simonet P, et al. Natural transformation in the *Ralstonia solanacearum* species

- complex: number and size of DNA that can be transferred. *FEMS Microbiol Ecol* 2008;66:14–24.
58. Guinard J, Latreille A, Guérin F, Poussier S, Wicker E. New Multi-locus Variable-Number Tandem-Repeat Analysis (MLVA) scheme for fine-scale monitoring and microevolution-related study of *Ralstonia pseudosolanacearum* phylotype I populations. *Appl Environ Microbiol* 2017;83:e03095–16.
 59. Guidot A, Coupat B, Fall S, Prior P, Bertolla F. Horizontal gene transfer between *Ralstonia solanacearum* strains detected by comparative genomic hybridization on microarrays. *ISME J* 2009;3:549–562.
 60. Prokhorchik M, Pandey A, Moon H, Kim W, Jeon H, et al. Host adaptation and microbial competition drive *Ralstonia solanacearum* phylotype I evolution in the Republic of Korea. *Microb Genom* 2020;6.
 61. Sabbagh CRR, Carrere S, Lonjon F, Vaillau F, Macho AP, et al. Pangenomic type III effector database of the plant pathogenic *Ralstonia* spp. *PeerJ* 2019;7:e7346.
 62. Gluck-Thaler E, Cerutti A, Perez-Quintero AL, Butchacas J, Roman-Reyna V, et al. Repeated gain and loss of a single gene modulates the evolution of vascular plant pathogen lifestyles. *Sci Adv* 2020;6:eabc4516.
 63. Bernal P, Llamas MA, Filloux A. Type VI secretion systems in plant-associated bacteria. *Environ Microbiol* 2018;20:1–15.
 64. Spraker JE, Sanchez LM, Lowe TM, Dorrestein PC, Keller NP. *Ralstonia solanacearum* lipopeptide induces chlamydospore development in fungi and facilitates bacterial entry into fungal tissues. *ISME J* 2016;10:2317–2330.
 65. Lowe-Power T, Avalos J, Munoz MC, Chipman K. A meta-analysis of the known global distribution and host range of the *Ralstonia* species complex. *bioRxiv* 2020.
 66. Etmnani F, Yousefvand M, Harighi B. Phylogenetic analysis and molecular signatures specific to the *Ralstonia solanacearum* species complex. *Eur J Plant Pathol* 2020;158:261–279.

Five reasons to publish your next article with a Microbiology Society journal

1. The Microbiology Society is a not-for-profit organization.
2. We offer fast and rigorous peer review – average time to first decision is 4–6 weeks.
3. Our journals have a global readership with subscriptions held in research institutions around the world.
4. 80% of our authors rate our submission process as 'excellent' or 'very good'.
5. Your article will be published on an interactive journal platform with advanced metrics.

Find out more and submit your article at microbiologyresearch.org.