



OPEN ACCESS

EDITED BY

David W. Ussery,
University of Arkansas for Medical Sciences,
United States

REVIEWED BY

Adela Finstrlová,
Masaryk University, Czechia
Adolphe Zeze,
Félix Houphouët-Boigny National Polytechnic
Institute,
Côte d'Ivoire

*CORRESPONDENCE

Itumeleng Matle
✉ matlei@arc.agric.za

RECEIVED 07 March 2023

ACCEPTED 26 June 2023

PUBLISHED 20 July 2023

CITATION

Carroll LM, Pierneef R, Mafuna T,
Magwedere K and Matle I (2023) Genus-wide
genomic characterization of *Macrocooccus*:
insights into evolution, population structure,
and functional potential.
Front. Microbiol. 14:1181376.
doi: 10.3389/fmicb.2023.1181376

COPYRIGHT

© 2023 Carroll, Pierneef, Mafuna, Magwedere
and Matle. This is an open-access article
distributed under the terms of the [Creative
Commons Attribution License \(CC BY\)](#). The
use, distribution or reproduction in other
forums is permitted, provided the original
author(s) and the copyright owner(s) are
credited and that the original publication in this
journal is cited, in accordance with accepted
academic practice. No use, distribution or
reproduction is permitted which does not
comply with these terms.

Genus-wide genomic characterization of *Macrocooccus*: insights into evolution, population structure, and functional potential

Laura M. Carroll^{1,2,3,4}, Rian Pierneef⁵, Thendo Mafuna⁶,
Kudakwashe Magwedere⁷ and Itumeleng Matle^{8*}

¹Department of Clinical Microbiology, SciLifeLab, Umeå University, Umeå, Sweden, ²Laboratory for Molecular Infection Medicine Sweden (MIMS), Umeå University, Umeå, Sweden, ³Umeå Centre for Microbial Research, Umeå University, Umeå, Sweden, ⁴Integrated Science Lab, Umeå University, Umeå, Sweden, ⁵Biotechnology Platform, Agricultural Research Council, Onderstepoort Veterinary Research, Onderstepoort, South Africa, ⁶Department of Biochemistry, University of Johannesburg, Auckland Park, South Africa, ⁷Directorate of Veterinary Public Health, Department of Agriculture, Land Reform and Rural Development, Pretoria, South Africa, ⁸Bacteriology Division, Agricultural Research Council, Onderstepoort Veterinary Research, Onderstepoort, South Africa

Introduction: *Macrocooccus* species have been isolated from a range of mammals and mammal-derived food products. While they are largely considered to be animal commensals, *Macrocooccus* spp. can be opportunistic pathogens in both veterinary and human clinical settings. This study aimed to provide insight into the evolution, population structure, and functional potential of the *Macrocooccus* genus, with an emphasis on antimicrobial resistance (AMR) and virulence potential.

Methods: All high-quality, publicly available *Macrocooccus* genomes ($n=104$, accessed 27 August 2022), plus six South African genomes sequenced here (two strains from bovine clinical mastitis cases and four strains from beef products), underwent taxonomic assignment (using four different approaches), AMR determinant detection (via AMRFinderPlus), and virulence factor detection (using DIAMOND and the core Virulence Factor Database).

Results: Overall, the 110 *Macrocooccus* genomes were of animal commensal, veterinary clinical, food-associated (including food spoilage), and environmental origins; five genomes (4.5%) originated from human clinical cases. Notably, none of the taxonomic assignment methods produced identical results, highlighting the potential for *Macrocooccus* species misidentifications. The most common predicted antimicrobial classes associated with AMR determinants identified across *Macrocooccus* included macrolides, beta-lactams, and aminoglycosides ($n=81$, 61, and 44 of 110 genomes; 73.6, 55.5, and 40.0%, respectively). Genes showing homology to *Staphylococcus aureus* exoenzyme aureolysin were detected across multiple species (using 90% coverage, $n=40$ and 77 genomes harboring aureolysin-like genes at 60 and 40% amino acid [AA] identity, respectively). *S. aureus* Pantone-Valentine leucocidin toxin-associated *lukF-PV* and *lukS-PV* homologs were identified in eight *M. canis* genomes ($\geq 40\%$ AA identity, $>85\%$ coverage). Using a method that delineates populations using recent gene flow (PopCOGenT), two species (*M. caseolyticus* and *M. armenti*) were composed of multiple within-species populations. Notably, *M. armenti* was partitioned into two populations, which differed in functional potential (e.g., one harbored beta-lactamase family, type II toxin-antitoxin system, and stress response proteins, while the other possessed a Type VII secretion system; PopCOGenT $p<0.05$).

Discussion: Overall, this study leverages all publicly available *Macrocooccus* genomes in addition to newly sequenced genomes from South Africa to identify

genomic elements associated with AMR or virulence potential, which can be queried in future experiments.

KEYWORDS

Macrococcus, *Macrococcus caseolyticus*, *Macrococcus armenti*, antimicrobial resistance, virulence, cattle, whole-genome sequencing, taxonomy

1. Introduction

Members of the *Macrococcus* genus are Gram-positive, catalase-positive, oxidase-positive, and coagulase-negative cocci (Mazhar et al., 2018; Ramos et al., 2021). The *Macrococcus* genus is a member of the Staphylococcaceae family and was first proposed as a novel genus in 1998, when its four original species (*M. caseolyticus*, *M. equipercicus*, *M. bovicus*, and *M. carouselicus*) were differentiated from members of the closely related *Staphylococcus* genus using numerous genetic and phenotypic characteristics (e.g., 16S rDNA sequencing, DNA–DNA hybridization, pulsed field gel electrophoresis, oxidase activity, cell wall composition, plasmid profiles; Kloos et al., 1998; Mazhar et al., 2018). Since the four original *Macrococcus* spp. were described in 1998, eight additional *Macrococcus* spp. have been identified ($n=12$ total validly published *Macrococcus* spp. per the List of Prokaryotic names with Standing in Nomenclature [LPSN], <https://lpsn.dsmz.de/genus/macrocooccus>; accessed 10 December 2022; Parte et al., 2020): *M. brunensis* (Mannerova et al., 2003), *M. hajekii* (Mannerova et al., 2003), *M. lamae* (Mannerova et al., 2003), *M. canis* (Gobeli Brawand et al., 2017), *M. bohemicus* (Maslanova et al., 2018), *M. epidermidis* (Maslanova et al., 2018), *M. goetzii* (Maslanova et al., 2018), and *M. armenti* (Keller et al., 2022).

Macrococcus spp. have historically been viewed as animal commensals (Mazhar et al., 2018) and have been isolated from a range of mammals (e.g., the skin of cows, pigs, horses, llamas, dogs) and the products derived from them (e.g., dairy products and meat; Kloos et al., 1998; Mannerova et al., 2003; Cotting et al., 2017; Mazhar et al., 2018; Ramos et al., 2021; Keller et al., 2022). However, the role of *Macrococcus* spp. as opportunistic pathogens has been discussed increasingly in recent years (MacFadyen et al., 2018; Ramos et al., 2021). In veterinary clinical settings, *Macrococcus* spp. have been isolated from infections (e.g., mastitis, otitis, and dermatitis cases, abscesses) in numerous animals, including cattle, sheep, and dogs (Gomez-Sanz et al., 2015; Cotting et al., 2017; Schwendener et al., 2017; Ramos et al., 2021). Notably, in 2018, *Macrococcus* spp. were reportedly isolated from human clinical samples for the first time, when *M. goetzii*, *M. epidermidis*, *M. bohemicus*, and *M. caseolyticus* subsp. *hominis* were isolated from infections at several body sites (i.e., wound sites, gynecological cases, and mycoses cases; Maslanova et al., 2018). Since then, *M. canis* has additionally been isolated from a human clinical case (i.e., a skin infection; Jost et al., 2021).

In addition to their pathogenic potential, some *Macrococcus* spp. carry antimicrobial resistance (AMR) genes (Schwendener et al., 2017; MacFadyen et al., 2018; Mazhar et al., 2018; Jost et al., 2021; Ramos et al., 2021). Methicillin resistance in *Macrococcus* spp. is of particular concern, as several mobilizable methicillin resistance determinants (e.g., penicillin-binding protein homologs *mecB*, *mecD*) have been identified in *Macrococcus* spp. (MacFadyen et al., 2018; Mazhar et al., 2018; Ramos et al., 2021). In this context, methicillin-resistant *Macrococcus* strains become particularly concerning: not only can they

potentially serve as opportunistic human and veterinary pathogens, but they can potentially transfer mobilizable AMR genes to other organisms, including taxa with a higher virulence potential (e.g., pathogenic *Staphylococcus aureus*; MacFadyen et al., 2018; Mazhar et al., 2018; Ramos et al., 2021; Schwendener and Perreten, 2022).

Several studies have employed genomic approaches to gain insight into the evolution and population structure of *Macrococcus*; however, these studies relied on a limited number of genomes (Maslanova et al., 2018; Schwendener and Perreten, 2022) and/or focused on specific taxa within the genus (e.g., *M. caseolyticus*; MacFadyen et al., 2018; Zhang et al., 2022). Furthermore, very few studies—genomic or otherwise—describing *Macrococcus* spp. strains isolated in Africa are available (Tshipamba et al., 2018; Ouoba et al., 2019; Ali et al., 2022). Here, we used whole-genome sequencing (WGS) to characterize six *Macrococcus* spp. strains isolated from bovine-associated sources in South Africa. To gain insight into *Macrococcus* at a genomic scale, we compare our six genomes to all publicly available *Macrococcus* genomes ($n=110$ total genomes). Overall, our study provides insight into the evolution, population structure, and functional potential of all species—both validly published and putative novel—within the *Macrococcus* genus in its entirety.

2. Materials and methods

2.1. Strain isolation

Macrococcus strains sequenced in this study were isolated from bovine clinical mastitis samples ($n=2$) and beef products ($n=4$) and submitted to the Onderstepoort Veterinary Research (OVR) General Bacteriology Laboratory for routine diagnostic services (Supplementary Table S1). From each sample, 10 g (ratio 1:10) were homogenized in buffered peptone water, and then aliquots of 0.1 mL were inoculated onto Baird-Parker agar and Brilliance MRSA 2 agar (both Oxoid, ThermoFisher, Johannesburg) and incubated for 24 h at 37°C. Presumptive macrococci colonies were streaked onto blood agar supplemented with 5% sheep blood (Oxoid, ThermoFisher, Johannesburg), incubated for 24 h at 37°C, and identified by phenotypic characteristics as described by Poyart et al. (2001). Briefly, Gram staining, catalase test, hemolysis, coagulase test, and API 32 ID STAPH (bioMérieux) were used to identify the isolates as macrococci.

2.2. Genomic DNA extraction, whole-genome sequencing, data pre-processing, and quality control

Genomic DNA was prepared from overnight cultures using the QIAGEN® DNeasy blood and tissue kit (Germany) according to the manufacturer's instructions (see section "Strain isolation" above;

Supplementary Table S1). WGS of isolates was performed at the Biotechnology Platform, Agricultural Research Council, Onderstepoort, South Africa. DNA libraries were prepared using TruSeq DNA library preparation kits (Illumina, San Diego, CA, USA), followed by sequencing on a HiSeq 2500 instrument (Illumina, San Diego, CA, USA).

Raw Illumina paired-end reads derived from each of the strains isolated here ($n=6$) were supplied as input to Trimmomatic v0.38 (Bolger et al., 2014). Trimmomatic was used to remove Illumina adapters (ILLUMINACLIP:TruSeq3-PE-2.fa:2:30:10:2:keepBothReads), leading and trailing low quality or N bases (i.e., Phred quality <3; LEADING:3 TRAILING:3), and reads <36 bp in length (MINLEN:36). FastQC v0.11.9¹ was used to evaluate the quality of the resulting trimmed paired-end reads (Supplementary Table S2).

The resulting trimmed paired-end reads associated with each strain were assembled into contigs via Shovill v1.1.0,² using the following parameters (all other parameters were set to their default values): (i) SKESA v2.4.0 (Souvorov et al., 2018) as the assembler (“--assembler skesa”); (ii) a minimum contig length of 200 (“--minlen 200”); (iii) a minimum contig coverage value of 10 (“--mincov 10”). QUAST v5.0.2 (Gurevich et al., 2013) was used to evaluate the quality of each resulting assembled genome (using a minimum contig length parameter of 1 bp), and the “lineage_wf” workflow in CheckM v1.1.3 (Parks et al., 2015) was used to evaluate genome completeness and contamination. MultiQC v1.12 (Ewels et al., 2016) was used to evaluate the quality of all six *Macrococcus* genomes in aggregate (Supplementary Tables S1, S2).

2.3. Acquisition and quality control of publicly available *Macrococcus* spp. genomes

All publicly available GenBank genomes submitted to the National Center for Biotechnology Information (NCBI) Assembly database as members of *Macrococcus* were downloaded ($n=102$ genomes; accessed 27 August 2022; Kitts et al., 2016; Schoch et al., 2020). Additionally, all genomes assigned to the *Macrococcus* genus within the Genome Taxonomy Database (GTDB) v207 (Parks et al., 2022), which were not included in the initial set of 102 genomes, were downloaded ($n=8$ of 88 total GTDB genomes). Together, this search of NCBI and GTDB yielded a preliminary set of 110 publicly available, putative *Macrococcus* genomes.

All 116 putative *Macrococcus* genomes (i.e., 110 publicly available genomes, plus the six genomes sequenced here) were characterized using QUAST and CheckM as described above (see section “Genomic DNA extraction, whole-genome sequencing, data pre-processing, and quality control” above). Six publicly available *Macrococcus* genomes showcased CheckM completeness <95% and/or QUAST N50 <20 Kbp; these genomes were excluded from further analysis ($n=104$ publicly available genomes used in subsequent analyses; Supplementary Table S3). One genome (NCBI GenBank

Assembly accession GCA_002119805.1) had >5% CheckM contamination (i.e., 5.11%; Supplementary Table S3). However, because this genome represented the type strain of *M. canis* and was a complete genome, it was used in subsequent steps. Overall, after removing low-quality genomes, the search of NCBI and GTDB, in combination with the six genomes sequenced here, yielded a final set of 110 *Macrococcus* genomes used in subsequent steps (Supplementary Tables S1, S3).

2.4. Taxonomic assignment

The following genomospecies delineation methods were applied to the set of 110 *Macrococcus* genomes (i.e., all 104 high-quality, publicly available *Macrococcus* genomes, plus the six genomes sequenced here; Supplementary Tables S1, S3): (i) the Genome Taxonomy Database Toolkit (GTDB-Tk), a popular genomospecies delineation tool, which relies primarily on a 95 average nucleotide identity (ANI) genomospecies threshold; (ii) bactaxR, which uses pairwise ANI values calculated between a set of genomes to delineate genomospecies *de novo* at any user-specified genomospecies threshold; (iii) the specI taxonomy, a marker gene-based taxonomic assignment approach; (iv) PopCOGenT (Populations as Clusters Of Gene Transfer; Arevalo et al., 2019), a method that relies on a metric of recent gene flow to identify species units. Each of these methods is explained in detail below.

For the GTDB-Tk workflow, all 110 *Macrococcus* genomes were assigned to species using the GTDB-Tk v2.1.0 “classify_wf” workflow (default settings) and version R207_v2 of GTDB (Chaumeil et al., 2019; Parks et al., 2022). GTDB-Tk confirmed that all 110 genomes identified here belonged to the *Macrococcus* genus (i.e., either “g__*Macrococcus*” or “g__*Macrococcus_B*,” per GTDB’s nomenclature; these corresponded to the only two GTDB genus designations, which contained the term “*Macrococcus*”; Supplementary Table S4).

For the bactaxR workflow, pairwise ANI values were calculated between all 110 *Macrococcus* genomes using the command-line implementation of OrthoANI v1.40 (Lee et al., 2016) with default settings (Supplementary Table S5). The resulting pairwise ANI values were supplied as input to the bactaxR v0.2.1 package (Carroll et al., 2020) in R v4.1.2 (R Core Team, 2021); bactaxR was used to construct a dendrogram and graph of all genomes based on pairwise ANI (dis)similarities, using the ANI.dendrogram and ANI.graph functions, respectively. bactaxR’s ANI.dendrogram function was further used to construct *de novo* genomospecies clusters using a 95 ANI genomospecies threshold (selected because this genomospecies threshold has been widely adopted by the microbiological community; Supplementary Table S6; Jain et al., 2018). OrthoANI was additionally used to calculate ANI values between all 110 *Macrococcus* genomes identified here (query genomes) relative to all *Macrococcus* spp. type strain genomes available in NCBI (reference genomes, $n=16$ type strain genomes, accessed 4 October 2022; Supplementary Tables S3, S5).

For the specI workflow, each *Macrococcus* genome was assigned to a marker gene-based species cluster (specI cluster) using classify-genomes (<https://github.com/AlessioMilanese/classify-genomes>; accessed 3 June 2020; Milanese et al., 2019) and version 3 of the specI taxonomy (Mende et al., 2013). specI clusters reported by classify-genomes were treated as species assignments (Supplementary Table S7).

1 <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

2 <https://github.com/tseemann/shovill>

For the PopCOGenT workflow, the “PopCOGenT” module within PopCOGenT (latest version downloaded 31 August 2022; Arevalo et al., 2019) was used to identify gene flow units among all 110 *Macrococcus* genomes. The resulting “main clusters” reported by PopCOGenT (i.e., gene flow units, which attempt to mimic the classical species definition used for animals and plants) were treated as species assignments (Supplementary Table S8; Arevalo et al., 2019). Two PopCOGenT main clusters (i.e., Main Clusters 0 and 2; Supplementary Table S8) contained >1 subcluster (i.e., within-species populations identified via PopCOGenT, referred to hereafter as “subclusters”); each of these main clusters was additionally queried individually using the “flexible genome sweeps” module in PopCOGenT to identify subcluster-specific orthologues, using an “alpha” (significance) value of 0.05 (Supplementary Tables S9, S10; Arevalo et al., 2019).

2.5. *In silico* multi-locus sequence typing

Each of the 110 *Macrococcus* genomes (Supplementary Tables S1, S3; see section “Acquisition and quality control of publicly available *Macrococcus* spp. genomes” above) was supplied as input to mlst v2.22.0³ for *in silico* multi-locus sequence typing (MLST). Default settings were used so that mlst could auto-select a MLST scheme from PubMLST (Jolley and Maiden, 2010; Jolley et al., 2018). Of the 110 genomes, 62 and 23 genomes were queried using the *M. caseolyticus* (“mcaseolyticus”) and *M. canis* (“mcanis”) PubMLST schemes, respectively; for 25 genomes, no scheme could be applied (Supplementary Table S11).

2.6. Genome annotation

Prokka v1.14.6 (Seemann, 2014) was used to annotate each *Macrococcus* genome ($n=110$, Supplementary Tables S1, S3; see section “Acquisition and quality control of publicly available *Macrococcus* spp. genomes” above), using the “Bacteria” database and default settings. The “.gff” and “.faa” files produced by Prokka, along with the assembled contigs associated with each strain, were supplied as input to AMRFinderPlus v3.10.40 (Feldgarden et al., 2019), which was used to identify antimicrobial resistance (AMR) determinants in each genome, using the “plus” option (“--plus,” i.e., to enable a search of the extended AMRFinderPlus database, which includes genes involved in virulence, biocide, heat, metal, and acid resistance) and the Prokka annotation format (“--annotation_format prokka”; Supplementary Table S12).

Amino acid (AA) sequences of virulence factors in the Virulence Factor Database (VFDB) core database (Liu et al., 2019) were downloaded ($n=4,188$ AA sequences in the VFDB core database; accessed 4 September 2022). CD-HIT v4.8.1 (Li and Godzik, 2006; Fu et al., 2012) was used to cluster all VFDB core database AA sequences using the “cd-hit” command, a sequence identity threshold of 0.4 (“-c 0.4”), and a word length of 2 (“-n 2,” the word size recommended for a 0.4 sequence identity threshold; <https://github.com/weizhongli/>

[cdhit/blob/master/doc/cdhit-user-guide.wiki](https://github.com/weizhongli/cdhit/blob/master/doc/cdhit-user-guide.wiki)). The “makedb” command in DIAMOND v2.0.15 (Buchfink et al., 2015) was used to construct a DIAMOND database of the VFDB core database in its entirety, and the “diamond blastp” command was used to query AA sequences derived from each *Macrococcus* genome (i.e., “.faa” files produced by Prokka) against the entire VFDB core database, using the following parameters (default values were used for all other parameters): ultra-sensitive mode (“--ultra-sensitive”), one reported maximum target sequence (“--max-target-seqs 1,” corresponding to the best match produced by DIAMOND: <https://github.com/bbuchfink/diamond/issues/29>), a minimum percent AA identity threshold of 60% (“--id 60”), and a minimum subject coverage threshold of 50% (“--subject-cover 50”). Each search was repeated using all combinations of (i) minimum percent AA identity thresholds of 0, 40, and 60%, and (ii) minimum subject coverage thresholds of 50 and 90% (Supplementary Tables S13–S18). Because many VFDB virulence factors are composed of multiple genes, and because some genes in VFDB may be highly similar/redundant, virulence factor presence and absence was considered at the whole virulence factor level, where a gene within a given virulence factor was considered to be “present” if any gene within its CD-HIT cluster could be detected in a given genome using DIAMOND. For example, the *S. aureus* exotoxin Panton-Valentine leukocidin (PVL) is a two-component toxin (Loffler et al., 2010; Shallcross et al., 2013). In the VFDB core database, PVL (VFDB ID VF0018) is composed of two genes: *lukF-PV* and *lukS-PV* (VFDB IDs VFG001276 and VFG001277, respectively). If any gene within the CD-HIT cluster of *lukF-PV* was detected in a *Macrococcus* genome, *lukF-PV* was considered “present”; likewise, if any gene within the CD-HIT cluster of *lukS-PV* was detected, *lukS-PV* was considered “present.” If both genes were “present,” PVL as a whole was considered to be 100% present. If one gene was “present,” PVL was considered to be 50% present. If neither gene was “present,” PVL was absent (0% present).

Biosynthetic gene clusters (BGCs) were detected in all 110 *Macrococcus* genomes using the command-line implementations of: (i) antiSMASH v6.1.0, using the “bacteria” taxon option (“--taxon bacteria”) and gene finding via Prodigal’s metagenomic mode option (“--genefinding-tool prodigal-m”; Blin et al., 2021); (ii) GECCO v0.9.2, using the “gecco run” command and the cluster probability threshold lowered to 0.3 (“-m 0.3”; all other settings were set to their defaults; Carroll et al., 2021). GenBank files (“.gbk”) for all BGCs identified by antiSMASH and GECCO were supplied as input to BiG-SCAPE v1.1.2 (Navarro-Munoz et al., 2020), which was used to cluster the 309 BGCs identified here, as well as experimentally validated BGCs in the MIBiG v2.1 database (“--mibig”) into Gene Cluster Families (GCFs) using default parameter values (Supplementary Table S19; Kautsar et al., 2020).

2.7. Genus-level phylogeny construction

Panaroo v1.2.7 (Tonkin-Hill et al., 2020) was used to identify orthologous gene clusters and construct a core genome alignment (“-a”) among the 110 *Macrococcus* genomes (see section “Acquisition and quality control of publicly available *Macrococcus* spp. genomes” above), plus *S. aureus* str. DSM 20231 as an outgroup genome (NCBI RefSeq Assembly accession GCF_001027105.1; $n=111$ total genomes). The following input/parameters were used (all other parameters were

³ <https://github.com/tseemann/mlst>

set to their default values): (i) each genome's ".gff" file produced by Prokka as input (see section "Genome annotation" above); (ii) MAFFT as the aligner ("--aligner mafft"; [Kato and Standley, 2013](#)); (iii) strict mode ("--clean-mode strict"); (iv) a core genome threshold of 95% ("--core_threshold 0.95"); (v) a protein family sequence identity threshold of 50% ("--f 0.5"). The core gene alignment produced by Panaroo ("core_gene_alignment.aln") was supplied as input to IQ-TREE v1.5.4 ([Nguyen et al., 2015](#)), which was used to construct a maximum likelihood (ML) phylogeny, using the General Time-Reversible (GTR) nucleotide substitution model ("--m GTR"; [Tavaré, 1986](#)) and 1,000 replicates of the ultrafast bootstrap approximation ("--bb 1,000"; [Minh et al., 2013](#)). The resulting ML phylogeny was visualized using the iTOL v6 webserver (<https://itol.embl.de/>; [Letunic and Bork, 2021](#)).

The genus-level phylogeny produced using Panaroo was compared to genus-level trees constructed using other methods, specifically: (i) PEPPAN ([Zhou et al., 2020](#)), a pipeline that can construct pan-genomes from genetically diverse bacterial genomes (e.g., spanning the diversity of an entire genus), and (ii) GTDB-Tk, which, in addition to taxonomic assignment, produces a multiple sequence alignment (MSA) of 120 bacterial marker genes detected in all input genomes ([Chaumeil et al., 2019](#)). For (i) PEPPAN, ".gff" files produced by Prokka were used as input ($n = 111$ total genomes, including the *S. aureus* outgroup; see section "Genome annotation" above). Default settings were used, except for the "--match_identity" option (the minimal identity of an alignment to be considered during pan-genome construction), which was set to "0.4," and the "--orthology" option (the algorithm for separating paralogous genes from orthologous genes), which was set to "ml" (i.e., the maximum-likelihood algorithm, reportedly the most accurate; [Zhou et al., 2020](#)). The PEPPAN_parser command was used to produce a Core Genome Allelic Variation (CGAV) tree (using a core genome threshold of 95%; "--a 95"), a gene presence/absence tree ("--tree"), and pan- and core-genome rarefaction curves ("--curve"; [Simonsen et al., 2008](#); [Tettelin et al., 2008](#); [Camacho et al., 2009](#); [Price et al., 2010](#); [Steinberger and Soding, 2017](#)). All aforementioned PEPPAN/PEPPAN_parser steps were repeated three separate times: (a) once as described above, but without the outgroup genome, and (b) using a lower minimal identity threshold (i.e., 20%, "--match_identity 0.2"), with and without the outgroup genome. The resulting trees were annotated using iTOL, and the resulting rarefaction curves were plotted in R using ggplot2 v3.4.0 ([Supplementary Figures S1–S6](#); [Wickham, 2016](#)). For the (ii) GTDB-Tk phylogeny, GTDB-Tk was run as described above, with the addition of the outgroup genome (see section "Taxonomic assignment" above). The resulting AA MSA produced by GTDB-Tk was supplied to IQ-TREE, which was used to construct a ML phylogeny as described above, but with the "LG+F+R4" AA substitution model (i.e., the optimal AA substitution model selected using IQ-TREE's implementation of ModelFinder, based on Bayesian Information Criteria [BIC] values; [Yang, 1995](#); [Le and Gascuel, 2008](#); [Soubrier et al., 2012](#); [Kalyaanamoorthy et al., 2017](#)). iTOL was used to plot the resulting phylogeny ([Supplementary Figure S7](#)).

2.8. Functional enrichment analyses

As mentioned above, two PopCOGenT main clusters (i.e., Main Clusters 0 and 2) contained >1 subcluster (see section "Taxonomic

assignment" above; [Supplementary Table S8](#)). To gain insight into the functional potential of subcluster-specific genes, which had been acquired post-speciation and differentially swept through subclusters identified via PopCOGenT (i.e., flexible genes identified via PopCOGenT, see section "Taxonomic assignment" above; [Supplementary Tables S9, S10](#)), functional enrichment analyses were conducted.

Briefly, for each relevant PopCOGenT main cluster (i.e., Main Cluster 0 and Main Cluster 2; see section "Taxonomic assignment" above), open reading frames (ORFs) produced by PopCOGenT for all members of the given main cluster were supplied as input to the eggNOG-mapper v2.1.9 web server (<http://eggno-mapper.embl.de/>; accessed 26 November 2022; [Huerta-Cepas et al., 2019](#); [Cantalapiedra et al., 2021](#)). eggNOG-mapper was used to functionally annotate each ORF, using default settings for all parameters except the input data type option (which was set to "CDS," as DNA sequences were used as input) and the "Gene Ontology evidence" option, which was set to "Transfer all annotations (including inferred from electronic annotation)."

For each PopCOGenT subcluster within the given main cluster, enrichment analyses were conducted to identify Gene Ontology (GO) terms ([Ashburner et al., 2000](#); [The Gene Ontology Consortium, 2018](#)) assigned via eggNOG-mapper, which were overrepresented among the PopCOGenT flexible genes identified within that particular subcluster: flexible genes identified within the given subcluster were treated as positive instances (PopCOGenT $p < 0.05$; [Supplementary Tables S9, S10](#)), and all other genes within the main cluster were treated as negative instances. Only genes with ≥ 1 assigned GO term were maintained. GO terms enriched within the positive instances (i.e., the subcluster-specific flexible genes identified via PopCOGenT; [Supplementary Tables S9, S10](#)) were identified via the "runTest" function in the topGO v2.46.0 R package ([Alexa et al., 2006](#)), using a Fisher's exact test (FET) with the "weight01" algorithm. Tests were conducted using each of the Biological Process (BP), Molecular Function (MF), and Cellular Component (CC) ontologies, using a minimum topGO node size of 3 for each ontology (i.e., "nodeSize=3," where topGO prunes the GO hierarchy from the terms with <3 annotated genes). GO terms were considered to be significantly enriched in the flexible genome of a PopCOGenT subcluster if the resulting FET p -value was < 0.05 ; no additional multiple testing correction was applied, as the "weight01" algorithm accounts for GO graph topology and produces p -values, which can be viewed as inherently corrected or not affected by multiple testing ([Alexa et al., 2006](#)). This approach was repeated for each subcluster within PopCOGenT Main Clusters 0 and 2 ([Supplementary Tables S20–S24](#)).

2.9. Species-level phylogeny construction

Species-level phylogenies were additionally constructed for the following: (i) GTDB's *M. caseolyticus* genomospecies, as it was composed of multiple PopCOGenT subclusters and contained five of the six South African genomes sequenced in this study ($n = 58$ genomes; [Supplementary Tables S4, S8](#)); (ii) bactaxR Cluster 13, corresponding to an unknown GTDB genomospecies, which contained the *M. armenti* type strain, because it, too, was composed of multiple PopCOGenT subclusters ($n = 8$ genomes; [Supplementary Tables S6, S8](#)); (iii) bactaxR Cluster 2, as it contained one of the South African genomes sequenced in this study ($n = 4$ genomes; [Supplementary Table S6](#)).

For *M. caseolyticus* and bactaxR Cluster 13 (i.e., *M. armentii*), Panaroo was used to construct a core gene alignment as described above (see section “Genus-level phylogeny construction” above), using a protein family sequence identity threshold of 70% (“-f 0.7”), all genomes assigned to the respective species cluster as input, and the following outgroup genomes: (i) a *Macrococcus* spp. genome from bactaxR Cluster 2 for *M. caseolyticus* (NCBI GenBank Assembly accession GCA_019357535.1), and (ii) a *M. canis* genome for bactaxR Cluster 13 (NCBI GenBank Assembly accession GCA_014524485.1; [Supplementary Figure S7](#)). Each resulting core gene alignment was supplied as input to IQ-TREE, and ML phylogenies were constructed and annotated as described above (see section “Genus-level phylogeny construction” above).

For *M. caseolyticus*, which was composed of >2 PopCOGenT subclusters, RhierBAPs v1.1.4 ([Tonkin-Hill et al., 2018](#)) was additionally employed to cluster the 58 *M. caseolyticus* genomes using two clustering levels. Briefly, Panaroo was used to construct a core gene alignment as described above but with the outgroup genome omitted ($n=58$ total *M. caseolyticus* genomes). Core SNPs were identified within the resulting core gene alignment using snp-sites v2.5.1 ([Page et al., 2016](#); using the “-c” option), and the resulting core SNP alignment was supplied as input to RhierBAPs.

For bactaxR Cluster 2, all genomes were fairly closely related (>99.2 ANI via OrthoANI); thus, Snippy v4.6.0⁴ was used to identify core SNPs among all four genomes within this species cluster, using the closed chromosome of one of the bactaxR Cluster 2 genomes as a reference (NCBI Nucleotide accession NZ_CP079969.1; [Li and Durbin, 2009](#); [Li et al., 2009](#); [Quinlan and Hall, 2010](#); [Li, 2011](#); [Cingolani et al., 2012](#); [Garrison and Marth, 2012](#); [Li, 2013](#); [Tan et al., 2015](#); [Page et al., 2016](#); [Li, 2019](#); [Seemann, 2019](#)). For the bactaxR Cluster 2 genome sequenced in this study, trimmed paired-end reads were used as input; for the three publicly available genomes, assembled genomes were used as input. Snippy was run using default settings, and the resulting cleaned alignment was supplied as input to Gubbins v3.1.3 ([Croucher et al., 2015](#)) to remove recombination using default settings. The resulting recombination-free alignment produced by Gubbins was queried using snp-sites as described above, and the resulting core SNP alignment was supplied as input to IQ-TREE. IQ-TREE was used to construct a ML phylogeny using an ascertainment bias correction based on the number of constant sites in the Snippy alignment (“-fconst 645,581,381,484,377,636,655,813”), one thousand replicates of the ultrafast bootstrap approximation (“-bb 1,000”), and the optimal nucleotide substitution model selected using ModelFinder (“-m MFP”; the K3Pu model, based on its BIC value; [Kimura, 1981](#); [Kalyanamoothy et al., 2017](#)). The resulting phylogeny was displayed and annotated using FigTree v1.4.4.⁵ The aforementioned steps were repeated, with the genome sequenced in this study omitted, as the remaining three genomes were highly similar on a genomic scale (>99.99 ANI via OrthoANI for the three publicly available bactaxR Cluster 2 genomes; note that Gubbins was not used here, as there were only three genomes available). Pairwise core SNP distances between genomes were calculated in R using the dist.gene function (with “method” set to “pairwise”) in ape v5.6.2 ([Paradis et al., 2004](#); [Paradis and Schliep, 2019](#)). Snippy was additionally

used to identify SNPs between other closely related genomes identified in the study (i.e., >99.9 ANI via OrthoANI), using default settings.

3. Results

3.1. Multiple GTDB species are represented among bovine-associated South African *Macrococcus* strains

Of the *Macrococcus* strains isolated in South Africa that underwent WGS (i.e., two veterinary isolates from bovine clinical mastitis cases, plus four food isolates from beef products), five were assigned to the *M. caseolyticus* genomospecies using the Genome Taxonomy Database Toolkit (GTDB-Tk; [Table 1](#); [Supplementary Table S4](#)). These five genomes each shared 98.0–98.6 average nucleotide identity (ANI) with the closed type strain genome of *M. caseolyticus* (calculated via OrthoANI relative to the *M. caseolyticus* type strain genome with NCBI RefSeq Assembly accession GCF_016028795.1; [Supplementary Table S5](#)), which is well above the 95 ANI threshold typically used for prokaryotic species delineation ([Jain et al., 2018](#)). When compared to each other, the five *M. caseolyticus* genomes sequenced here shared 97.9–99.4 ANI via OrthoANI. One genome (S135) was assigned to PubMLST *M. caseolyticus* sequence type 2 (ST2), while another (S139) was assigned to ST16; the remaining three *M. caseolyticus* genomes belonged to unknown STs ([Supplementary Table S11](#)).

Notably, however, one food isolate (S115) could not be assigned to any known species within GTDB ([Table 1](#); [Supplementary Table S4](#)). Strain S115 was isolated in 2015 from beef biltong, a South African spiced intermediate moisture, ready-to-eat (RTE) meat product, which was being sold in a retail outlet in South Africa’s Limpopo province ([Table 1](#); [Supplementary Table S1](#)). When compared to the five *M. caseolyticus* genomes sequenced here, S115 shared <95 ANI with each (via OrthoANI; [Supplementary Table S5](#)). When compared to the type strain genomes of all *Macrococcus* species, S115 was most closely related to *M. caseolyticus* subsp. *hominis* str. CCM 7927 (NCBI RefSeq Assembly accession GCF_002742395.2), sharing 95.3 ANI via OrthoANI ([Supplementary Table S5](#)). Comparatively, S115 shared 94.6 ANI with the closed *M. caseolyticus* type strain genome (via OrthoANI, NCBI RefSeq Assembly accession GCF_016028795.1; [Supplementary Table S5](#)).

Overall, these results indicate that, among the bovine-associated South African *M. caseolyticus* genomes sequenced here, (i) considerable within-species diversity exists (e.g., multiple STs are represented, novel STs are present, ANI values between strains sequenced in this study are not particularly high); (ii) one or two *Macrococcus* genomospecies are represented, depending on the species delineation method used (i.e., GTDB or ANI-based comparisons to type strain genomes; [Table 1](#)).

3.2. Human clinical, veterinary clinical, and food spoilage-associated strains are represented among *Macrococcus* spp. genomes

To compare the bovine-associated South African *Macrococcus* genomes sequenced here to *Macrococcus* genomes collected from other

⁴ <https://github.com/tseemann/snippy>

⁵ <http://tree.bio.ed.ac.uk/software/figtree/>

TABLE 1 South African *Macrococcus* spp. genomes sequenced in this study ($n = 6$).

Strain	Year of isolation	Province	Animal	Sample type	Isolation source ^a	Establishment category	GTDB species ^b
S99	1991	Gauteng	Cattle	Veterinary clinical sample	Milk from mastitis case	Farm	<i>M. caseolyticus</i>
S125	1992	Gauteng	Cattle	Veterinary clinical sample	Milk from mastitis case	Farm	<i>M. caseolyticus</i>
S120	2015	Gauteng	Cattle	Meat sample	RTE beef biltong	Retail outlet	<i>M. caseolyticus</i>
S139	2015	Gauteng	Cattle	Meat sample	Minced beef	Butchery	<i>M. caseolyticus</i>
S135	2015	Free State	Cattle	Meat sample	Processed beef patties	Retail outlet	<i>M. caseolyticus</i>
S115	2015	Limpopo	Cattle	Meat sample	RTE beef biltong	Retail outlet	<i>M. spp. nov.</i>

^aRTE, ready-to-eat.

^bSpecies assigned using the Genome Taxonomy Database (GTDB) Toolkit (GTDB-Tk) v2.1.0 and version R207_v2 of GTDB; all six genomes were assigned to GTDB's "Macrococcus_B" genus.

sources in other world regions, the six genomes sequenced here were aggregated with all high-quality, publicly available *Macrococcus* genomes ($n = 110$ total genomes; Figure 1; Supplementary Table S3). Overall, the complete set of 110 *Macrococcus* genomes represented strains collected from at least 10 countries, with most genomes originating from Europe (88 of 110 genomes, 80.0%; Figure 1; Supplementary Tables S1, S3).

A vast majority of the genomes (97 of 110 genomes, 88.2%) originated from animal- and animal product-associated sources, with over half of all strains originating from bovine-associated sources (60 of 110 total genomes, 54.5%; Figure 1; Supplementary Tables S1, S3). Numerous animal-associated strains, including two strains sequenced here, were reportedly clinical in origin (e.g., isolated from bovine mastitis cases, canine ear infection cases, and wound infections in donkeys; Table 1; Supplementary Tables S1, S3). Several animal-associated strains, including four sequenced here, were isolated from food products with the potential for human consumption (i.e., beef and pork meat, cow milk, cheese); one strain was isolated from a food product with a known defect (i.e., "ropy" milk; Table 1; Supplementary Tables S1, S3).

Interestingly, six of the 110 *Macrococcus* genomes (5.5%) were derived from human-associated strains (Figure 1; Supplementary Table S3). At least five of these strains were isolated in conjunction with human clinical cases, including: (i) a hemolytic, methicillin-resistant *M. canis* strain isolated from a 52-year-old immunocompromised patient with cutaneous maculopapular and impetigo lesions (Switzerland, 2019); (ii) a *M. bohemicus* strain from an 80–85 year-old patient with a traumatic knee wound (Czech Republic, 2003); (iii) a *M. goetzii* strain from a foot nail mycosis case in a 30–35 year-old patient (Czech Republic, 2000); (iv) a *M. caseolyticus* subsp. *hominis* strain from an acute vaginitis case in a 40–45 year old patient (Czech Republic, 2003); (v) a *M. epidermidis* strain associated with mycose in a 66–70 year old patient (Czech Republic, 2001; Figure 1; Supplementary Table S3; Maslanova et al., 2018; Jost et al., 2021).

Overall, *Macrococcus* WGS efforts have overwhelming queried animal or animal product-associated strains, although several human clinical strains have undergone WGS (Figure 1).

3.3. *Macrococcus* genomospecies clusters may overlap at a conventional 95 ANI threshold

Overall, using GTDB-Tk, *Macrococcus* encompassed 15 genomospecies: 12 defined genomospecies, plus three undefined/

putative novel genomospecies defined using a conventional 95 ANI threshold (Figure 1; Supplementary Table S4). One of these putative novel GTDB-Tk genomospecies encompassed strain S115 sequenced here (denoted as bactaxR Cluster 2 in Figure 1), plus publicly available genomes submitted to NCBI as *M. caseolyticus* (Supplementary Table S3). All members of this genomospecies shared <95 ANI with the *M. caseolyticus* type strain genome but >95 ANI with the *M. caseolyticus* subsp. *hominis* type strain genome (via OrthoANI, NCBI RefSeq Assembly accessions GCF_016028795.1 and GCF_002742395.2, respectively; Figure 2). The second putative novel GTDB-Tk genomospecies (denoted as bactaxR Cluster 13 in Figure 1) contained the type strain of *M. armenti* (NCBI GenBank Assembly accession GCA_020097135.1); considering *M. armenti* was published as a novel species in 2022, it is likely this genomospecies will be described as such in future versions of GTDB (Keller et al., 2022). The third putative novel GTDB-Tk genomospecies contained a single genome (denoted as bactaxR Cluster 14 in Figure 1), which had been submitted to NCBI as *M. caseolyticus* (strain Ani-LG-066, NCBI GenBank Assembly accession GCA_021366795.1); however, this genome shared <75 ANI with the *M. caseolyticus* and *M. caseolyticus* subsp. *hominis* type strain genomes and shared <81.0 ANI with all other *Macrococcus* spp. genomes (via OrthoANI; Figures 1, 2). Strain Ani-LG-066, which was isolated in 2016 from milk samples taken from a lactating dairy cow in Liege, Belgium, thus likely represents a truly novel *Macrococcus* genomospecies (Figure 1; Supplementary Tables S3, S5).

Importantly, for three of the four genomospecies delineation methods used here (i.e., GTDB-Tk, bactaxR, and PopCOGenT), genomes assigned to separate genomospecies could share >95 ANI with each other (Figure 2; Supplementary Figures S8, S9), indicating that some *Macrococcus* genomospecies defined at a conventional 95 ANI threshold overlap. specI did not yield overlapping genomospecies at 95 ANI (Supplementary Figure S9); however, nearly a third of *Macrococcus* genomes ($n = 32$ of 110 genomes, 29.1%) could not be assigned to a species via specI (Figure 1; Supplementary Figure S9; Supplementary Table S7).

Taken together, these results indicate that (i) three of the four genomospecies delineation methods queried here (i.e., GTDB-Tk, bactaxR with OrthoANI and a 95 ANI threshold, and PopCOGenT) produced similar, albeit not identical, results when applied to *Macrococcus* (Figure 1); (ii) the same three genomospecies delineation methods produced "overlapping genomospecies," in which some genomes could share >95 ANI with members of another genomospecies (Figure 2; Supplementary Figure S9).

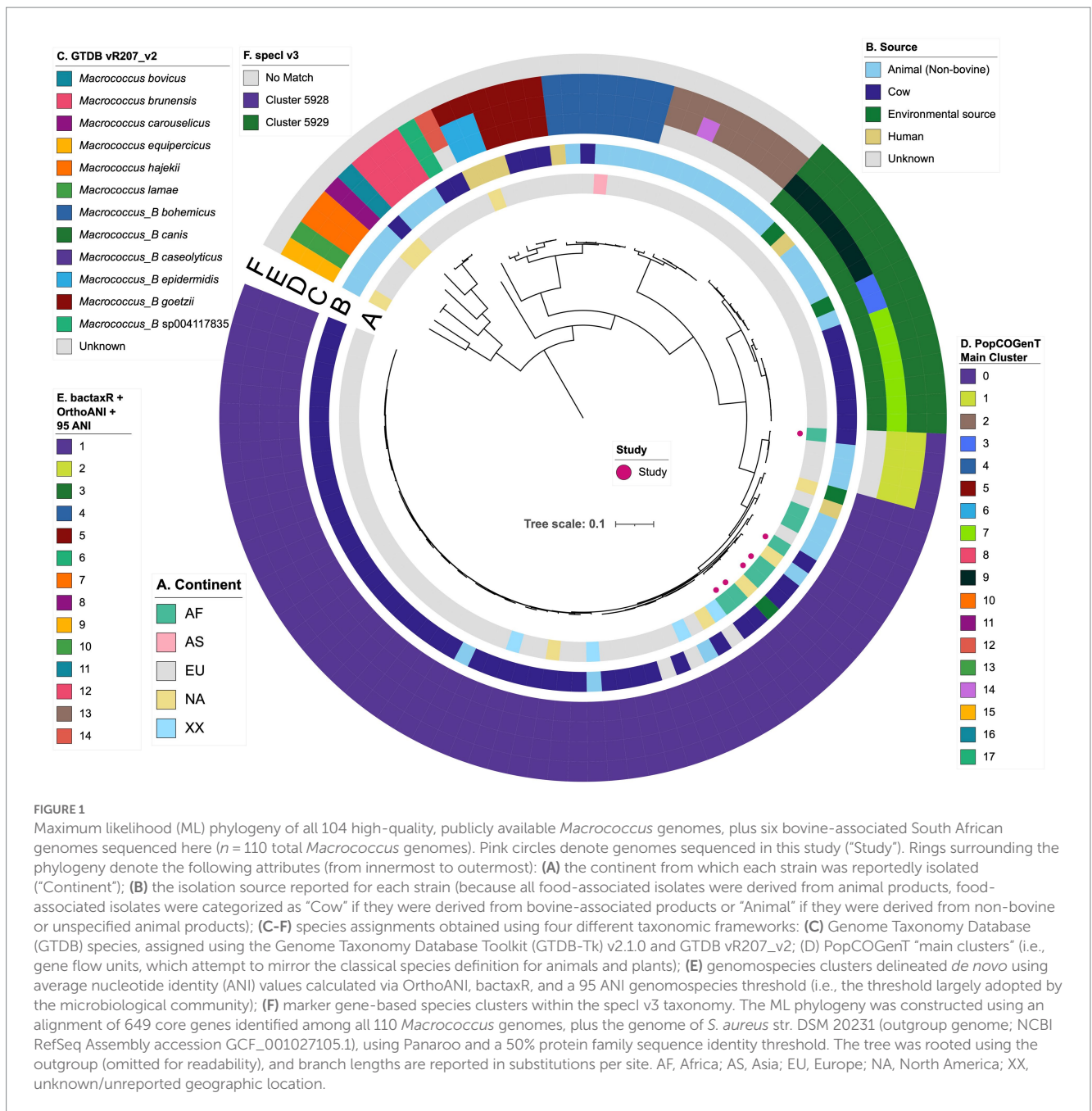


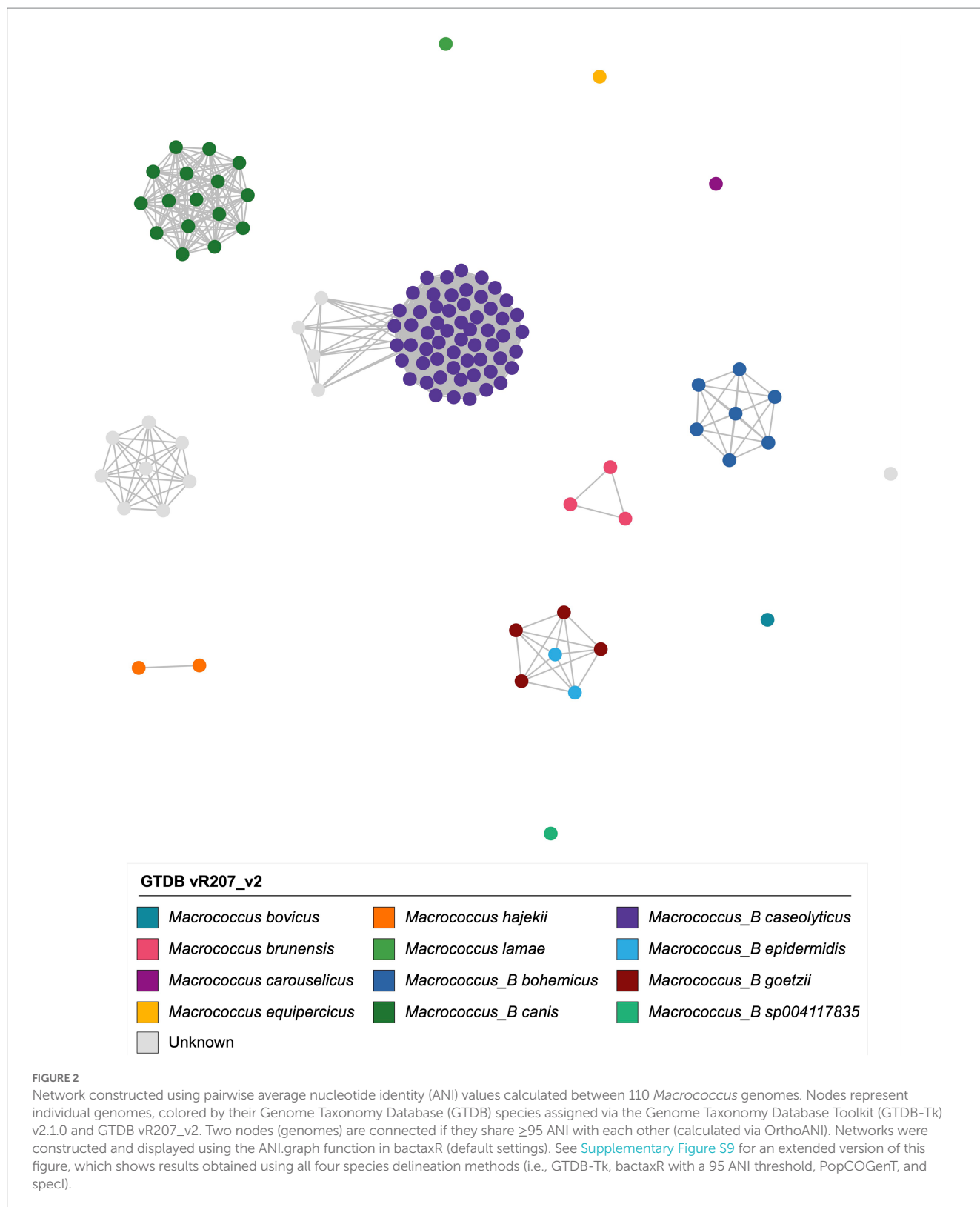
FIGURE 1

Maximum likelihood (ML) phylogeny of all 104 high-quality, publicly available *Macrocooccus* genomes, plus six bovine-associated South African genomes sequenced here ($n = 110$ total *Macrocooccus* genomes). Pink circles denote genomes sequenced in this study ("Study"). Rings surrounding the phylogeny denote the following attributes (from innermost to outermost): (A) the continent from which each strain was reportedly isolated ("Continent"); (B) the isolation source reported for each strain (because all food-associated isolates were derived from animal products, food-associated isolates were categorized as "Cow" if they were derived from bovine-associated products or "Animal" if they were derived from non-bovine or unspecified animal products); (C-F) species assignments obtained using four different taxonomic frameworks: (C) Genome Taxonomy Database (GTDB) species, assigned using the Genome Taxonomy Database Toolkit (GTDB-Tk) v2.1.0 and GTDB vR207_v2; (D) PopCOGenT "main clusters" (i.e., gene flow units, which attempt to mirror the classical species definition for animals and plants); (E) genomospecies clusters delineated *de novo* using average nucleotide identity (ANI) values calculated via OrthoANI, bactaxR, and a 95 ANI genomospecies threshold (i.e., the threshold largely adopted by the microbiological community); (F) marker gene-based species clusters within the specI v3 taxonomy. The ML phylogeny was constructed using an alignment of 649 core genes identified among all 110 *Macrocooccus* genomes, plus the genome of *S. aureus* str. DSM 20231 (outgroup genome; NCBI RefSeq Assembly accession GCF_001027105.1), using Panaroo and a 50% protein family sequence identity threshold. The tree was rooted using the outgroup (omitted for readability), and branch lengths are reported in substitutions per site. AF, Africa; AS, Asia; EU, Europe; NA, North America; XX, unknown/unreported geographic location.

3.4. Multiple *Macrocooccus* spp. contain genomes that are predicted to be multi-drug resistant

Antimicrobial resistance (AMR) and stress response determinants (detected via AMRFinderPlus; Supplementary Table S12) were variably present throughout *Macrocooccus* and were associated with predicted resistance to a variety of antimicrobial classes, heavy metals, and metalloids (Figure 3; Supplementary Figure S10). The most common classes of antimicrobials for which *Macrocooccus* was predicted to harbor resistance determinants included macrolides, beta-lactams, and aminoglycosides ($n = 81$, 61, and 44 of 110 genomes with one or more associated AMR determinants, corresponding to 73.6, 55.5, and 40.0% of *Macrocooccus* genomes, respectively; Figure 3; Supplementary Figure S10;

Supplementary Table S12). The high proportion of genomes harboring an ATP-binding cassette subfamily F protein (ABC-F)-encoding gene (*abc-f*) contributed to the high proportion of genomes with predicted macrolide resistance ($n = 74$ of 110 *Macrocooccus* genomes harbored *abc-f*, 67.3%), although several additional macrolide resistance genes were sporadically present within the genus (Figure 3; Supplementary Figure S10; Supplementary Table S12). The high proportion of genomes showcasing predicted beta-lactam resistance, on the other hand, was largely driven by the presence of *mecD* ($n = 43$ of 110 *Macrocooccus* genomes, 39.1%), although *mecB* and *bla* were also present in >10% of genomes (Figure 3; Supplementary Figure S10; Supplementary Table S12). Aminoglycoside resistance genes were sporadically present among *Macrocooccus* genomes, the most common being *str* ($n = 23$ of 110 genomes, 20.9%; Figure 3; Supplementary Figure S10; Supplementary Table S12).



The most common AMR profiles among *Macrocooccus* genomes harboring one or more AMR determinant were those associated with (i) macrolide and (ii) beta-lactam/macrolide resistance ($n = 18$ and 13 of 110 genomes, corresponding to 16.4 and 11.8% of genomes, respectively; [Figure 3](#); [Supplementary Figure S10](#);

[Supplementary Table S12](#)). However, numerous predicted multidrug-resistance (MDR) profiles were observed, the most common being (i) aminoglycoside/beta-lactam/macrolide and (ii) aminoglycoside/beta-lactam/macrolide/tetracycline resistance ($n = 11$ and 8 of 110 genomes, corresponding to 10.0 and 7.3% of

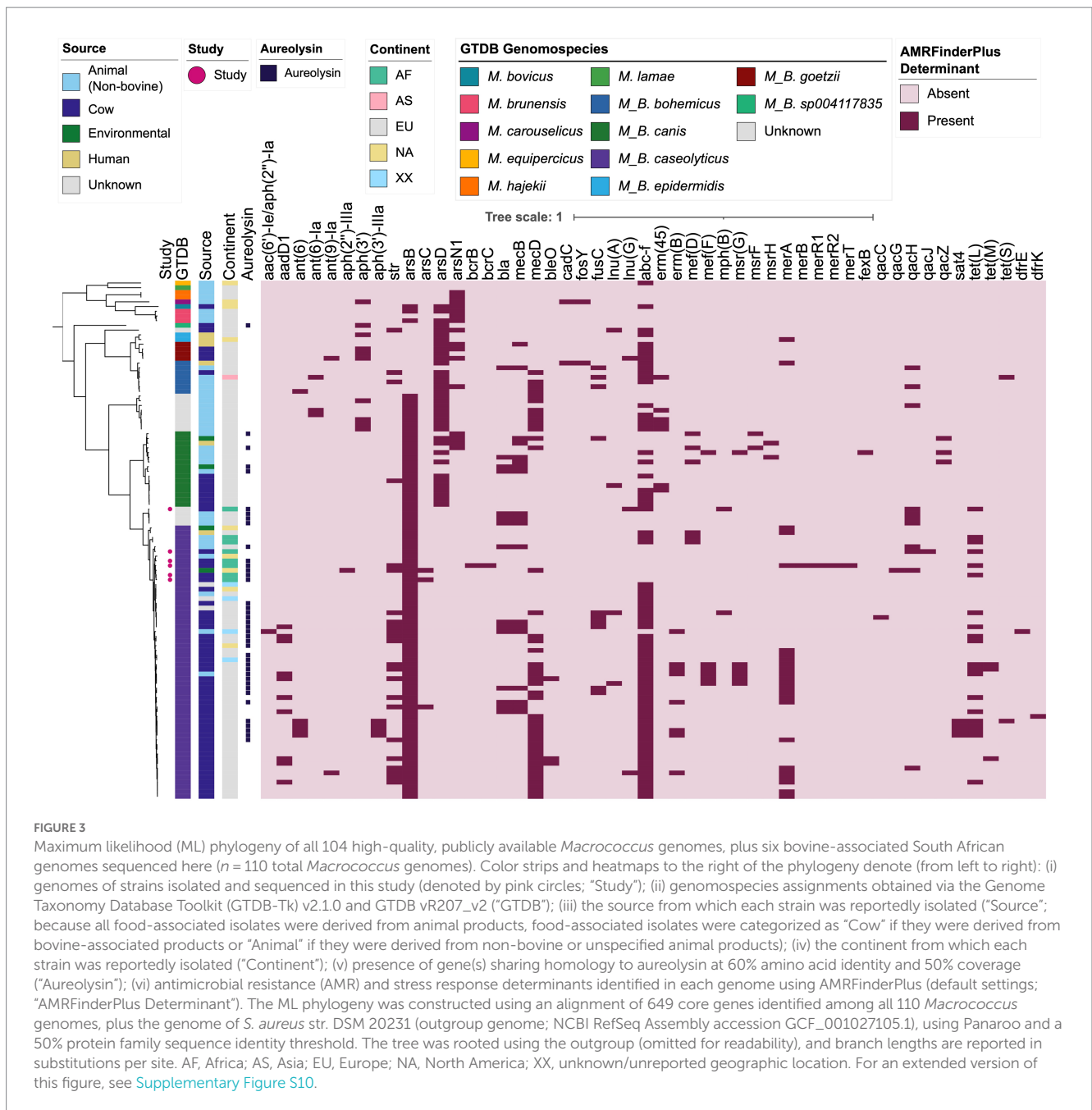


FIGURE 3
 Maximum likelihood (ML) phylogeny of all 104 high-quality, publicly available *Macrocooccus* genomes, plus six bovine-associated South African genomes sequenced here ($n = 110$ total *Macrocooccus* genomes). Color strips and heatmaps to the right of the phylogeny denote (from left to right): (i) genomes of strains isolated and sequenced in this study (denoted by pink circles; "Study"); (ii) genospecies assignments obtained via the Genome Taxonomy Database Toolkit (GTDB-Tk) v2.1.0 and GTDB vR207_v2 ("GTDB"); (iii) the source from which each strain was reportedly isolated ("Source"; because all food-associated isolates were derived from animal products, food-associated isolates were categorized as "Cow" if they were derived from bovine-associated products or "Animal" if they were derived from non-bovine or unspecified animal products); (iv) the continent from which each strain was reportedly isolated ("Continent"); (v) presence of gene(s) sharing homology to aureolysin at 60% amino acid identity and 50% coverage ("Aureolysin"); (vi) antimicrobial resistance (AMR) and stress response determinants identified in each genome using AMRFinderPlus (default settings; "AMRFinderPlus Determinant"). The ML phylogeny was constructed using an alignment of 649 core genes identified among all 110 *Macrocooccus* genomes, plus the genome of *S. aureus* str. DSM 20231 (outgroup genome; NCBI RefSeq Assembly accession GCF_001027105.1), using Panaroo and a 50% protein family sequence identity threshold. The tree was rooted using the outgroup (omitted for readability), and branch lengths are reported in substitutions per site. AF, Africa; AS, Asia; EU, Europe; NA, North America; XX, unknown/unreported geographic location. For an extended version of this figure, see [Supplementary Figure S10](#).

genomes, respectively; [Figure 3](#); [Supplementary Figure S10](#); [Supplementary Table S12](#)). The genome displaying predicted AMR to the most antimicrobial classes was the genome of *M. caseolyticus* strain 5813_BC74, which had reportedly been isolated from bovine bulk tank milk in the United Kingdom in 2016 (NCBI GenBank Assembly accession GCA_002834615.1; [Supplementary Table S3](#)). This genome displayed predicted aminoglycoside/beta-lactam/fusidic acid/lincosamide/macrolide/tetracycline resistance ($n = 6$ antimicrobial classes; [Figure 3](#); [Supplementary Figure S10](#); [Supplementary Table S12](#)).

Predicted AMR phenotypes observed in <10% of all *Macrocooccus* genomes included: (i) fusidic acid resistance (due to the presence of *fusC*; $n = 10$), (ii) lincosamide resistance (per *lnu(A)*, *lnu(G)*; $n = 7$),

(iii) streptothricin resistance (via *sat4*; $n = 4$), (iv) bleomycin resistance (via *bleO*; $n = 3$), (v) trimethoprim (via *dfrE*, *dfrK*) and (vi) fosfomycin resistance (via *fosY*, $n = 2$ genomes each), and (vii) phenicol resistance (via *fexB*, $n = 1$ genome; [Figure 3](#); [Supplementary Figure S10](#); [Supplementary Table S12](#)). Interestingly, one genome (South African strain S99, from a bovine mastitis case in Gauteng in 1991) harbored bacitracin resistance genes (i.e., *bcrB* and *bcrC*; [Figure 3](#); [Supplementary Figure S10](#); [Supplementary Table S12](#)).

Overall, these results indicate that (i) numerous AMR determinants are variably present within and among *Macrocooccus* species; and (ii) *Macrocooccus* genomes may harbor AMR determinants predictive of an MDR phenotype (i.e., resistant to three or more antimicrobial classes; [Figure 3](#); [Supplementary Figure S10](#); [Supplementary Table S12](#)). However,

these results should be interpreted with caution, as AMR potential was not evaluated phenotypically in this study.

3.5. *Staphylococcus aureus* virulence factor homologues can be detected within some *Macrococcus* genomes at low amino acid identity

To gain insight into the virulence potential of *Macrococcus*, the 110 genomes aggregated here were queried for virulence factors present in the VFDB core database (Figure 3; Supplementary Figure S11; Supplementary Tables S13–S18). Notably, genes encoding the *S. aureus* exoenzyme aureolysin could be detected across multiple *Macrococcus* species (using 90% coverage, $n=40$ and 77 genomes harboring aureolysin-encoding genes at 60 and 40% amino acid [AA] identity, respectively; Figure 3; Supplementary Figure S11; Supplementary Tables S16, S18).

Further, genes sharing $\geq 40\%$ AA identity and $>85\%$ coverage with *S. aureus* Pantone-Valentine leucocidin (PVL) toxin-associated *lukF-PV* and *lukS-PV* were identified in eight *M. canis* genomes (per GTDB-Tk; Supplementary Figure S11; Supplementary Table S15). AA identities of the *Macrococcus lukF-PV* relative to *S. aureus lukF-PV* (VFDB ID VFG001276) ranged from 50.5 to 53.9%, with a mean of 52.6%; for *lukS-PV* (VFDB ID VFG001277), AA identities of 49.5–49.8% were observed, with a mean of 49.6%. For all eight *M. canis* genomes in which they were detected, the *lukF-PV* and *lukS-PV* homologs were located next to each other in the genome (Supplementary Figure S11; Supplementary Table S15).

Overall, these results indicate that proteins homologous to virulence factors present in other species (e.g., *S. aureus*) can be detected in some *Macrococcus* genomes. However, the methods employed here are not adequate to properly evaluate the virulence potential of *Macrococcus* strains that possess these homologs; thus, these results should be interpreted with extreme caution.

3.6. *Macrococcus* species differ in pan-genome composition

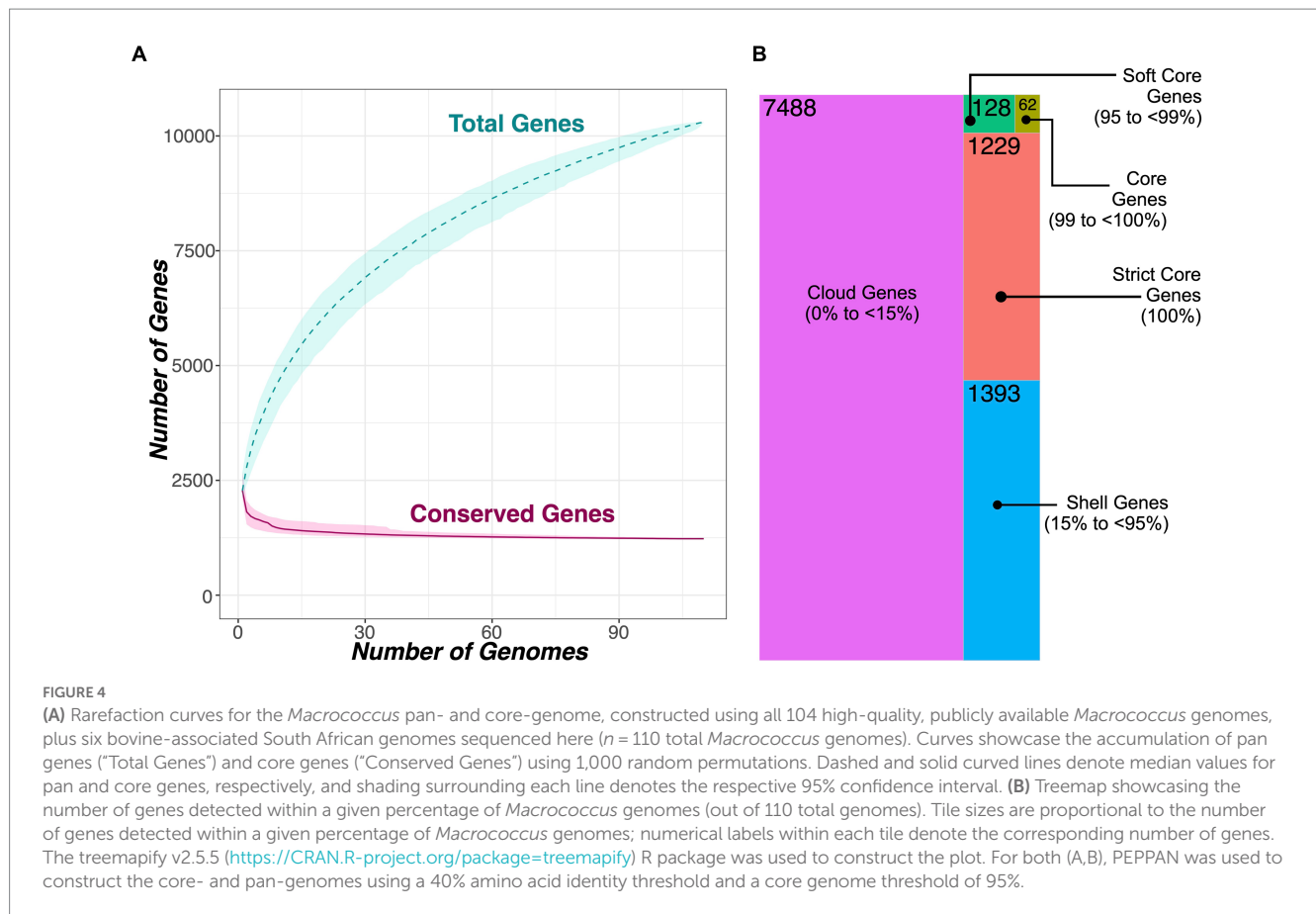
Using PEPPAN and a 40% AA identity threshold, a total of 10,300 genes were detected among the 110 *Macrococcus* genomes aggregated here, 1,229 of which were core genes present in all 110 genomes (11.9% of all *Macrococcus* genes; Figure 4; Supplementary Figures S1, S3, S5). Comparatively, at a 20% AA identity threshold, 9,835 total genes were detected, 1,235 of which were core genes present in all 110 genomes (12.6% of all *Macrococcus* genes; Supplementary Figures S2–S4, S6). Based on trees constructed using pan-genome element presence/absence, *Macrococcus* species (per GTDB-Tk) tended to cluster together based on pan-genome composition, although not exclusively (Supplementary Figures S1, S2). Specifically, the topology of the PEPPAN pan-genome tree differed from that of the PEPPAN Core Genome Allelic Variation (CGAV) tree, as some *Macrococcus* species were polyphyletic based on pan-genome element presence/absence (Supplementary Figures S1, S2). Overall, *Macrococcus* species tend to differ via both core genome phylogeny (Figures 1, 3) and pan-genome composition (Figure 4; Supplementary Figures S1–S6).

3.7. *Macrococcus caseolyticus* and *Macrococcus armenti* are composed of multiple within-species subclusters separated by recent gene flow

PopCOGenT identified 18 “main clusters” (species) across *Macrococcus* in its entirety; within two of these main clusters (i.e., PopCOGenT Main Clusters 0 and 2 in Figure 1), PopCOGenT identified multiple “subclusters” separated by recent gene flow (i.e., populations that were still connected by some gene flow, but had significantly more gene flow within the population than between populations; Figure 1). Specifically, (i) within PopCOGenT Main Cluster 0 (corresponding to GTDB-Tk’s *M. caseolyticus* genomospecies), five subclusters were identified, and (ii) within PopCOGenT Main Cluster 2 (an unknown species via GTDB-Tk, which contains the *M. armenti* type strain and will thus be referred to as *M. armenti* hereafter), two subclusters were identified. As such, we will discuss these two species individually in detail below (Figure 1).

3.7.1. African and European *Macrococcus caseolyticus* strains largely belong to separate lineages

The 58 *M. caseolyticus* genomes (per GTDB-Tk) were divided into five PopCOGenT subclusters and five RhierBAPS clusters, although the composition of those (sub)clusters differed slightly (Figure 5; Supplementary Figure S12; Supplementary Table S8). Notably, the majority of European *M. caseolyticus* genomes ($n=33$ of 42 European *M. caseolyticus* genomes, 78.6%) were assigned to a well-supported clade within the species phylogeny (referred to hereafter as the “*M. caseolyticus* major European lineage,” which is denoted in Figure 5 as RhierBAPS Cluster 1; ultrafast bootstrap support = 100%). Members of the *M. caseolyticus* major European lineage were overwhelmingly of bovine origin (34 of 36 RhierBAPS Cluster 1 genomes, 94.4%), and nearly all genomes within the lineage were reportedly isolated from European countries: 30 from the United Kingdom (83.3%), and two and one genome(s) from Switzerland and Ireland, respectively (5.6 and 2.8%); the only genome isolated from outside of Europe was reportedly isolated from rosy milk in the United States in 1920 (NCBI GenBank Assembly accession GCA_900453015.1; Figure 5; Supplementary Figure S12; Supplementary Table S3). Interestingly, most genomes within the *M. caseolyticus* major European lineage were predicted to be MDR (Figures 3, 5; Supplementary Figure S12). Specifically, (i) all genomes in the *M. caseolyticus* major European lineage (36 of 36 genomes, 100%) were predicted to be resistant to macrolides, largely due to the presence of *abc-f* (35 of 36 *M. caseolyticus* major European lineage genomes, 97.2%; the only genome in which *abc-f* was not detected possessed *erm(B)* and was thus still predicted to be macrolide-resistant via AMRFinderPlus); (ii) nearly all (33 of 36 genomes, 91.7%) were predicted to be resistant to beta-lactams, largely due to the presence of *mecD* in 28 genomes (77.8% of 36 genomes in the lineage; the remaining five genomes that were predicted to be beta-lactam-resistant harbored *bla* and *mecB*); (iii) a majority (21 of 36 genomes in the lineage, 58.3%) were predicted to be resistant to aminoglycosides, due largely to the presence of *str* and/or *aadD1* (detected in 14 and 9 of 36 genomes, 38.9 and 25.0%, respectively; Figures 3, 5; Supplementary Figure S12). Additionally, three genomes possessed genes conferring resistance to bleomycin; these were the



only genomes within the *Macrocooccus* genus, which harbored bleomycin resistance-conferring gene *bleO* (Figures 3, 5; Supplementary Figure S12).

Of the nine European *M. caseolyticus* genomes that were not members of the *M. caseolyticus* major European lineage, seven belonged to a well-supported clade containing 10 genomes (ultrafast bootstrap support = 100%; referred to hereafter as the "*M. caseolyticus* minor European lineage," which is denoted in Figure 5 as RhierBAPS Cluster 2 and PopCOGenT Subcluster 0.1). Aside from two genomes of unknown origin, one genome within the *M. caseolyticus* minor European lineage was reportedly of non-European origin (i.e., strain CCM 3540, reportedly isolated from cow's milk in the Washington, D.C. vicinity of the United States in 1916; NCBI GenBank Assembly accession GCA_003259685.1, Supplementary Table S3; Evans, 1916). Like the *M. caseolyticus* major European lineage, all genomes within the *M. caseolyticus* minor European lineage were predicted to be resistant to macrolides, as all harbored *abc-f* (Figures 3, 5; Supplementary Figure S12). However, a predicted MDR phenotype (i.e., resistant to three or more antimicrobial classes) was less prevalent among genomes within the minor European lineage ($n = 3$ of 10 *M. caseolyticus* minor European lineage genomes, 30%): the MDR genomes were similar on a genomic scale (99.7–99.9 ANI via OrthoANI) and were confined to a single, well-supported clade within the *M. caseolyticus* minor European lineage (ultrafast bootstrap support = 100%; Figure 5; Supplementary Figure S12). Additionally, within the *M. caseolyticus* minor European lineage, PopCOGenT identified six "flexible" genes (i.e., PopCOGenT subcluster-specific

orthologous gene clusters), which were specific to the *M. caseolyticus* minor European lineage (denoted as gene group "C" within the PopCOGenT Flexible Gene heatmap in Figure 5; PopCOGenT $p < 0.05$). All six genes were chromosomal and included (i) large conductance mechanosensitive channel protein MscL, and (ii) genes associated with Y-family DNA polymerases (Figure 5; Supplementary Figure S12; Supplementary Table S9). Compared to all other *M. caseolyticus* genes, numerous biological processes (BPs) and molecular functions (MFs) were enriched in the *M. caseolyticus* minor European lineage flexible genes, including DNA-related BPs/MFs (e.g., DNA biosynthesis, replication, and repair), and those related to ion binding/transport (topGO Fisher's Exact Test [FET] $p < 0.05$; Supplementary Tables S9, S20).

Of the 12 *M. caseolyticus* genomes that were not members of the major and minor European lineages, seven were African in origin, three were North American, and two were European, including the one human-associated *M. caseolyticus* genome (i.e., strain CCM 7927, which was isolated in Pribram, Czech Republic in 2003 from a vaginal swab taken from an acute vaginitis case in a 40–45 year-old patient, NCBI GenBank Assembly accession GCA_002742395.2; Figure 5; Supplementary Figure S12; Supplementary Table S3; Maslanova et al., 2018). Notably, of the five South African *M. caseolyticus* strains isolated and sequenced here, four were assigned to a single PopCOGenT subcluster (i.e., PopCOGenT Subcluster 0.3 in Figure 5). Unlike the major and minor European lineages, members of this subcluster did not possess macrolide resistance genes (Figures 3, 5; Supplementary Figure S12). AMR genes were detected sporadically

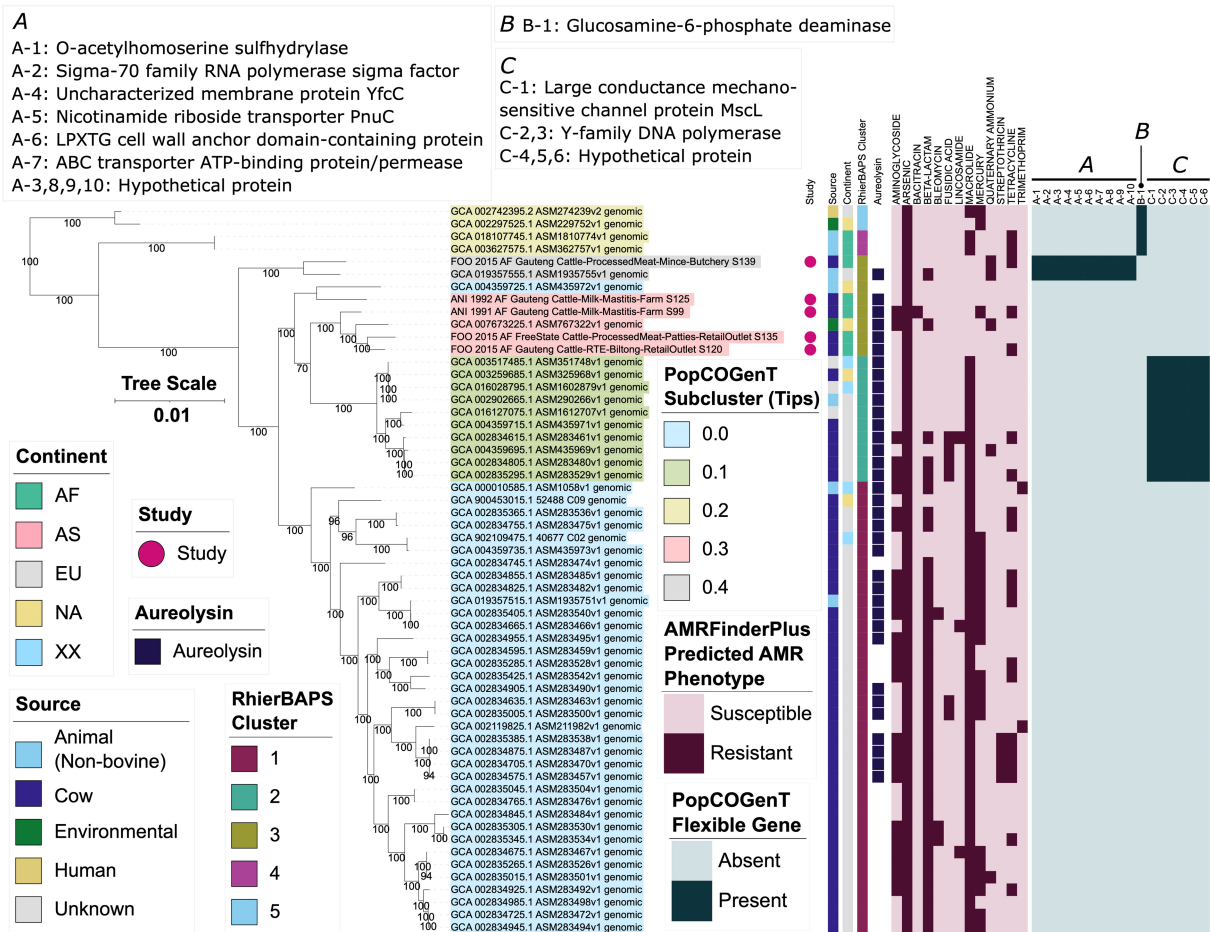


FIGURE 5
 Maximum likelihood (ML) phylogeny of 58 genomes assigned to the Genome Taxonomy Database’s (GTDB) *M. caseolyticus* genomospecies. Tip label colors correspond to subcluster assignments obtained using PopCOGenT (“PopCOGenT Subcluster”). Pink circles denote genomes sequenced in this study (“Study”). Color strips/heatmaps to the right of the phylogeny denote (from left to right): (i) the source from which each strain was reportedly isolated (“Source”; because all food-associated isolates were derived from animal products, food-associated isolates were categorized as “Cow” if they were derived from bovine-associated products or “Animal” if they were derived from non-bovine or unspecified animal products); (ii) the continent from which each strain was reportedly isolated (“Continent”); (iii) cluster assigned using RhierBAPS (“RhierBAPS Cluster”); (iv) presence of gene(s) sharing homology to aureolysin at 40% amino acid identity and 50% coverage (“Aureolysin”); (v) predicted antimicrobial resistance (AMR) and stress response phenotype, obtained using AMR and stress response determinants identified via AMRFinderPlus (“AMRFinderPlus Predicted AMR Phenotype”); (vi) presence and absence of flexible genes identified via PopCOGenT (“PopCOGenT Flexible Gene”), with corresponding gene annotations displayed in the boxes marked “A,” “B,” and “C.” The ML phylogeny was constructed using an alignment of 1,751 core genes identified among all 58 *M. caseolyticus* genomes, plus an outgroup *Macrococcus* spp. genome from bactaxR Cluster 2 (NCBI GenBank Assembly accession GCA_019357535.1; [Figure 1](#)), using Panaroo and a 70% protein family sequence identity threshold. The tree was rooted using the outgroup (omitted for readability), and branch lengths are reported in substitutions per site. Branch labels correspond to branch support percentages obtained using one thousand replicates of the ultrafast bootstrap approximation. AF, Africa; AS, Asia; EU, Europe; NA, North America; XX, unknown/unreported geographic location. For an extended version of this phylogeny, see [Supplementary Figure S12](#).

within these genomes. Specifically, (i) strain S99 possessed genes associated with aminoglycoside (streptomycin), bacitracin, and tetracycline resistance (*str*, *bcrBC*, and *tet(L)*, respectively); (ii) GCA_007673225.1 (an environmental strain isolated in 2018 in Durham, North Carolina, United States) possessed genes associated with aminoglycoside and beta-lactam resistance (i.e., *aph(2'')-IIIa*, *str*, and *mecD*, associated with amikacin/gentamicin/kanamycin/tobramycin, streptomycin, and methicillin resistance, respectively); (iii) S120 possessed tetracycline resistance gene *tet(L)* ([Figure 5](#); [Supplementary Figure S12](#)). Despite most genomes being South African in origin, the five *M. caseolyticus* genomes within this subcluster were considerably diverse, sharing 98.6–99.4 ANI with each other (via OrthoANI; [Figure 5](#); [Supplementary Figure S12](#)).

The remaining South African genome sequenced in this study (i.e., S139), plus GCA_019357555.1 (isolated from a calf nasal swab in Switzerland in 2019), were assigned to a separate subcluster via PopCOGenT (i.e., PopCOGenT Subcluster 0.4 in [Figure 5](#)). Neither genome possessed macrolide resistance genes, although both possessed quaternary ammonium resistance gene *qacH* ([Figure 5](#); [Supplementary Figure S12](#)). S139 additionally possessed tetracycline resistance gene *tet(L)*, while GCA_019357555.1 possessed beta-lactam resistance genes *mecB* (methicillin) and *bla* ([Figure 5](#); [Supplementary Figure S12](#)). Most notably, however, PopCOGenT identified 10 flexible genes within this subcluster (denoted as gene group “A” within the PopCOGenT Flexible Gene heatmap in [Figure 5](#), PopCOGenT $p < 0.05$; [Figure 5](#); [Supplementary Figure S12](#);

Supplementary Table S9); ATP- and transmembrane-associated BPs/MFs were enriched in this subcluster's flexible genes (topGO FET $p < 0.05$; Supplementary Table S21).

Four additional *M. caseolyticus* genomes were assigned to a single subcluster using PopCOGenT (i.e., PopCOGenT Subcluster 0.2 in Figure 5). Interestingly, like the major and minor European clades, three of the four genomes within this subcluster were predicted to be macrolide resistant, as they possessed *abc-f* and *mef(D)* (Figure 5; Supplementary Figure S12). Two highly similar genomes derived from strains isolated in 2016 from wounded animals in Sudan additionally possessed tetracycline resistance gene *tet(L)* (100.0 ANI and 0 SNPs via OrthoANI and Snippy, respectively, NCBI GenBank Assembly accessions GCA_018107745.1 and GCA_003627575.1; Figure 5; Supplementary Figure S12). Additionally, unlike the other *M. caseolyticus* subclusters described above, none of the genomes within this PopCOGenT subcluster possessed genes sharing homology to aureolysin-encoding genes (Figure 5; Supplementary Figure S12). Further, PopCOGenT identified one flexible gene within this subcluster (denoted as gene group "B" within the PopCOGenT Flexible Gene heatmap in Figure 5, PopCOGenT $p < 0.05$; Supplementary Table S9); glucosamine-6-phosphate deaminase, which was associated with the enrichment of several GO terms, including antibiotic catabolic process, carbohydrate metabolic process, and N-acetylglucosamine-associated processes (topGO FET $p < 0.05$; Figure 5; Supplementary Figure S12; Supplementary Tables S9, S22).

Overall, these results indicate that *M. caseolyticus* genomes from geographic regions outside of Europe, particularly Africa, belong to separate lineages within the species. However, future genomic sequencing efforts are needed to provide further evidence of lineage-geography associations.

3.7.2. Putative virulence factors are differentially associated with *Macrococcus armentii* lineages

Like *M. caseolyticus*, *M. armentii* could be differentiated into subclusters via PopCOGenT (Figure 6A; Supplementary Figure S13; Supplementary Table S8). Specifically, (i) PopCOGenT Subcluster 2.0 contained five genomes from animals in Switzerland (two from strains isolated from the nasal cavities of calves in 2019, and three from the skins of pigs in 2021); and (ii) PopCOGenT Subcluster 2.1 contained two genomes from pigs in Switzerland (one from the nasal cavity of a pig in 2017, and another from the skin of a pig in 2021; Figure 6A; Supplementary Figure S13; Supplementary Tables S3, S8). An additional genome, derived from a pig-associated strain isolated in the United Kingdom in 1963 (NCBI GenBank Assembly accession GCA_022808015.1) was additionally assigned to the *M. armentii* species via ANI-based methods (i.e., using OrthoANI, it shared 96.5–97.7 ANI with all other *M. armentii* genomes; Figure 2; Supplementary Table S5); however, PopCOGenT assigned this genome to a separate main cluster (i.e., "species"), and it was thus not included in the subsequent within-main cluster flexible gene analyses (Figure 6A; Supplementary Figure S13; Supplementary Table S8).

Within PopCOGenT Subcluster 2.0, PopCOGenT identified 43 flexible genes (denoted as gene group "A" within the PopCOGenT Flexible Gene heatmap in Figure 6A; PopCOGenT $p < 0.05$; Supplementary Table S10), which together were associated with the enrichment of eight GO terms (topGO FET $p < 0.05$; Supplementary Table S23). The most highly enriched GO terms were by far "diaminopimelate biosynthetic process" (GO:0019877) and

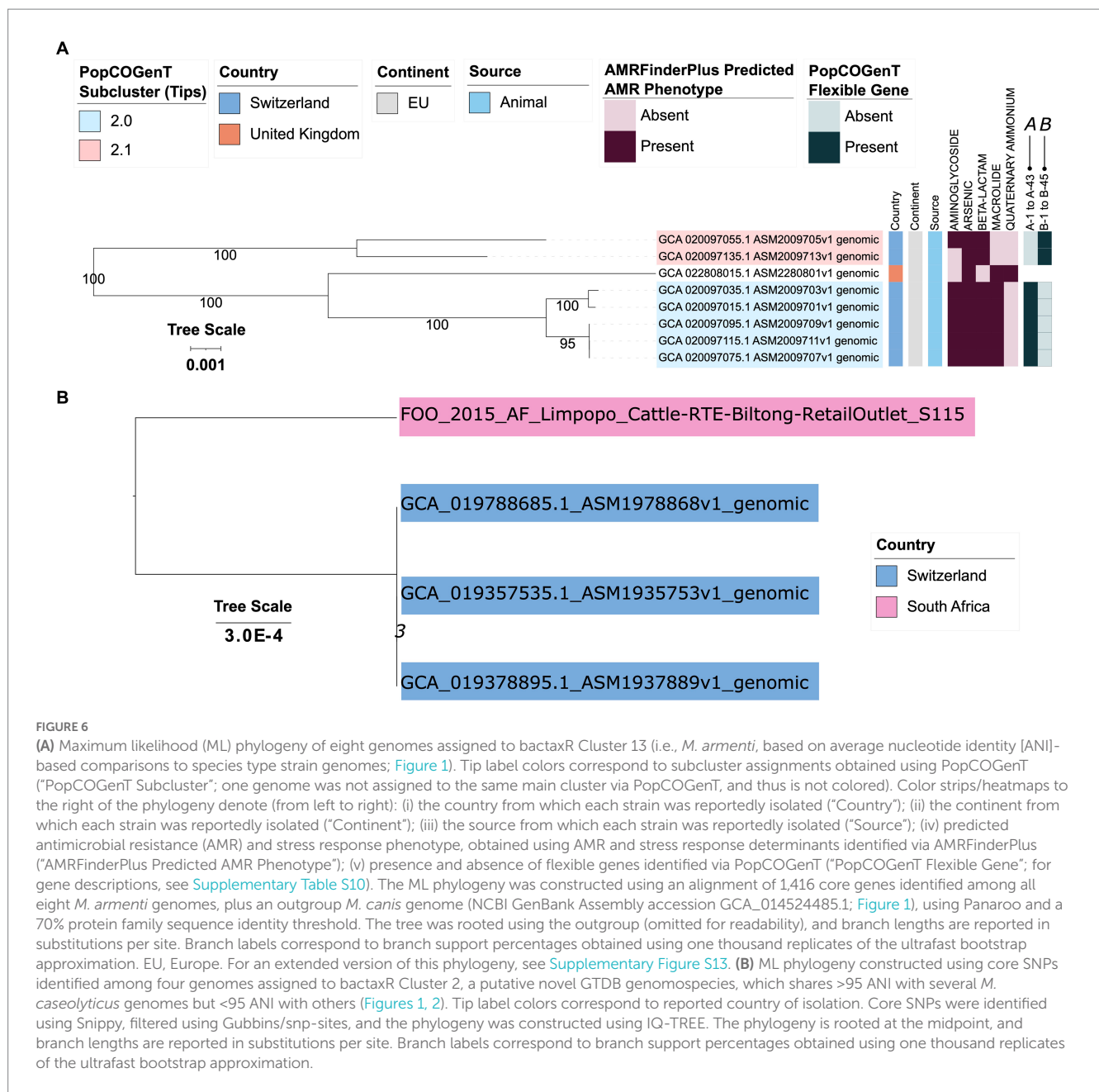
"lysine biosynthetic process via diaminopimelate" (GO:0009089, topGO FET $p < 1.0 \times 10^{-30}$; Supplementary Table S23), which were assigned to a cluster of three consecutive flexible genes (PopCOGenT $p < 0.05$): (i) 4-hydroxy-tetrahydrodipicolinate reductase/dihydrodipicolinate reductase *dapB* (NCBI Protein accession UBH07557.1); (ii) 2,3,4,5-tetrahydropyridine-2,6-dicarboxylate N-acetyltransferase *dapD* (NCBI Protein accession UBH09720.1); (iii) an amidohydrolase (NCBI Protein accession UBH07558.1; Supplementary Table S10).

Most notably, genes sharing homology to *S. aureus* Type VII secretion system proteins were among the flexible genes within Subcluster 2.0 (PopCOGenT $p < 0.05$), including genes sharing homology to extracellular protein EsxD (VFDB ID VFG049714), chaperone protein EsaE (VFDB ID VFG049701), secreted protein EsxB (VFDB ID VFG002411), secretion substrate EsaC (at 97% query coverage and 38% AA identity; NCBI Protein accession HCD1544785.1), EssB (NCBI Protein accession UBH08107.1), EsaA (NCBI Protein accession UBH08110.1), and secreted protein EsxA (VFDB ID VFG002405; Figure 6A; Supplementary Figure S13; Supplementary Table S10).

Several additional clusters of genes were among the flexible genes within Subcluster 2.0 (PopCOGenT $p < 0.05$; Figure 6A; Supplementary Figure S13; Supplementary Table S10), including: (i) a cluster of genes involved in nitrous oxide reduction, e.g., c-type cytochrome (NCBI Protein accession UBH08788.1), a Sec-dependent nitrous-oxide reductase (NCBI Protein accession UBH08789.1), nitrous oxide reductase family maturation protein NosD (NCBI Protein accession UBH08791.1); (ii) a cluster of genes that included an ImmA/IrrE family metallo-endopeptidase (NCBI Protein accession UBH09010.1), a LaCl family DNA-binding transcriptional regulator (NCBI Protein accession UBH09033.1), a sucrose-6-phosphate hydrolase (NCBI Protein accession UBH09034.1), a carbohydrate kinase (NCBI Protein accession UBH09035.1), and sucrose-specific PTS transporter subunit IIBC (NCBI Protein accession UBH09036.1); (iii) a cluster of genes that included a pathogenicity island protein (NCBI Protein accession UBH09209.1; Figure 6A; Supplementary Figure S13; Supplementary Table S10).

Interestingly, a protein most closely resembling immune inhibitor A was also identified by PopCOGenT as a flexible gene (at 98% query coverage and 97.65% AA identity, NCBI Protein accession WP_224185801.1, PopCOGenT $p < 0.05$; Figure 6A; Supplementary Figure S13; Supplementary Table S10). The "immune inhibitor A peptidase M6" protein domain identified in this protein (PFAM ID 05547) has previously been identified in virulence factors secreted by members of the *Bacillus cereus* group (immune inhibitor A; InhA) and *Vibrio cholerae* (secreted metalloprotease PrtV; Vaitkevicius et al., 2008).

Comparatively, within Subcluster 2.1, PopCOGenT identified 45 flexible genes (denoted as gene group "B" within the PopCOGenT Flexible Gene heatmap in Figure 6A; PopCOGenT $p < 0.05$) associated with 22 enriched GO terms (topGO FET $p < 0.05$; Figure 6A; Supplementary Figure S13; Supplementary Tables S10, S24). By far the most highly enriched GO term within this subcluster corresponded to BP "lipoteichoic acid biosynthetic process" (GO:0070395, topGO FET $p = 2.2 \times 10^{-18}$; Supplementary Table S24). Notably, a cluster of five consecutive, chromosomally encoded flexible genes within PopCOGenT Subcluster 2.1 were associated with (lipo) teichoic acid synthesis (PopCOGenT $p < 0.05$;



Supplementary Table S10): teichoic acid D-Ala incorporation-associated protein DltX (NCBI Protein accession UBH13741.1), D-alanine--poly(phosphoribitol) ligase subunit DltA (NCBI Protein accession UBH13742.1), D-alanyl-lipoteichoic acid biosynthesis protein DltB (NCBI Protein accession UBH13743.1), D-alanine--poly(phosphoribitol) ligase subunit 2 DltC (NCBI Protein accession UBH13744.1), and D-alanyl-lipoteichoic acid biosynthesis protein DltD (NCBI Protein accession UBH13745.1). Interestingly, this cluster of five genes was located several genes downstream of two consecutive, chromosomally encoded beta-lactamase family proteins, which were also both identified as being flexible genes (PopCOGenT $p < 0.05$). Both beta-lactamase family proteins were annotated via eggNOG-mapper as "autolysis and methicillin resistant-related protein PbpX" (NCBI Protein accessions UBH13736.1 and UBH13737.1) and were associated with "response to antibiotic"

(GO:0046677), a BP that was also enriched in PopCOGenT Subcluster 2.1 (topGO FET $p = 2.3 \times 10^{-3}$; Supplementary Tables S10, S24).

Several GO terms associated with transporter activity were also enriched in PopCOGenT Subcluster 2.1 (topGO FET $p < 0.05$), including MF "ABC-type transporter activity" (GO:0140359; Supplementary Table S24). Congruently, four separate clusters of genes containing regions annotated as ABC transporter components were included among PopCOGenT's set of flexible genes (PopCOGenT $p < 0.05$; Supplementary Tables S10, S24).

Interestingly, a protein annotated as immune inhibitor A was also among the flexible genes detected within PopCOGenT Subcluster 2.1 (PopCOGenT $p < 0.05$, NCBI Protein accession UBH13622.1; Figure 6A; Supplementary Figure S13; Supplementary Table S10). Further, genes encoding a type II toxin-antitoxin system were among the flexible genes identified by PopCOGenT within this PopCOGenT

subcluster (PopCOGenT $p < 0.05$; Figure 6A; Supplementary Figure S13; Supplementary Table S10), specifically: (i) a type II toxin-antitoxin system RelE/ParE family toxin, which was immediately upstream of (ii) a type II toxin-antitoxin system Phd/YefM family antitoxin (NCBI Protein accessions UBH12746.1 and UBH12747.1, respectively).

Overall, (i) *M. armentis* boasts two major subclusters, which are largely separated by recent gene flow; and (ii) flexible genes differentially present within these major subclusters (e.g., a type VII secretion system, toxin-antitoxin genes, beta-lactamase family genes) indicate that these two subclusters may differ phenotypically, although future experiments will be necessary to confirm this.

3.8. A novel GTDB genomospecies encompasses *Macrococcus* genomes from Switzerland and South Africa

Of the six *Macrococcus* spp. genomes sequenced in this study, five were assigned to *M. caseolyticus* (per GTDB-Tk; Figure 1; Supplementary Table S4). The genome of S115, however, could not be assigned to a known species via GTDB-Tk (Figure 1; Supplementary Table S4). Using *bactaxR* and a 95 ANI threshold (i.e., an approach similar to that of GTDB-Tk), three additional, publicly available genomes belonged to this putative novel GTDB genomospecies (i.e., *bactaxR* Cluster 2, $n = 4$ total genomes; Figures 1, 6B). In addition to (i) S115, a food-associated strain isolated in 2015 from beef biltong sold in South Africa's Limpopo province, this genomospecies included three strains isolated in Switzerland in 2019: (ii) 19Msa1099, isolated from pork meat (NCBI GenBank Assembly accession GCA_019357535.1), plus (iii) 19Msa1047 and (iv) 19Msa0499, each isolated from calf nasal swab samples (NCBI GenBank Assembly accessions GCA_019378895.1 and GCA_019788685.1, respectively; Supplementary Tables S1, S3). Notably, the South African genome was relatively distantly related to the Swiss genomes, sharing 99.2 ANI with each via OrthoANI and differing by 8,614–8,637 SNPs (identified via Snippy relative to each individual Swiss genome).

Comparatively, the three Swiss genomes shared >99.99 ANI with each other via OrthoANI and differed by 1–34 core SNPs (calculated via Snippy with the South African S115 strain excluded): strains 19Msa1047 (from a calf nasal swab) and 19Msa1099 (from pork meat) differed by a single core SNP identified in a gene annotated as a CBS domain-containing protein (NCBI Protein accession WP_219491817.1, corresponding to locus tag KYI07_RS05750 within the *M. caseolyticus* str. 19Msa0499 reference chromosome with NCBI Nucleotide accession NZ_CP079969.1). These two strains differed from strain 19Msa0499 (from a calf nasal swab) by 33 and 34 core SNPs, all of which fell within two regions of the *M. caseolyticus* str. 19Msa0499 reference chromosome: (i) 13 core SNPs within positions 312,553–367,746 bp, and (ii) 20 core SNPs within positions 1,778,236–1,778,444 bp, indicating that genetic differences within these regions may be due to recombination.

4. Discussion

In this study, WGS was used to characterize *Macrococcus* spp. strains isolated from South African cattle (i.e., two strains from bovine

clinical mastitis cases) and beef products (i.e., two stains from RTE beef biltong and two from minced/processed beef products). Using these genomes in combination with all publicly available, high quality *Macrococcus* spp. genomes, insight is provided into the evolution, population structure, and functional potential of the *Macrococcus* genus as a whole. Importantly, we observed (i) differences in functional potential (e.g., AMR potential, virulence potential) between and within *Macrococcus* spp., and (ii) that some *Macrococcus* species lack clear boundaries at conventional genomospecies delineation thresholds (i.e., 95 ANI), which may cause taxonomic issues in the future. Below, we discuss these findings in detail, as well as (iii) future opportunities in the *Macrococcus* genomics space.

4.1. Differences in functional potential can be observed between and within *Macrococcus* species

Bacteria can adapt to stressors and stimuli in their respective environments through the acquisition of genomic material in the “flexible” gene pool (Arevalo et al., 2019). Thus, intraspecies differences in genomic content can be observed for many bacterial species (Tonkin-Hill et al., 2020), and differences resulting from recent gene flow (i.e., genomic elements acquired post-speciation) can be used to delineate populations within those species (Arevalo et al., 2019). Here, we queried all *Macrococcus* spp. genomes and identified genomic determinants variably present within species, indicative of within-species differences in functional potential. For example, of the 50 putative AMR and stress response determinants identified across *Macrococcus* in its entirety, nearly half (24 of 50, 48%) were species-specific (based on GTDB-Tk species assignments); of these species-specific AMR and stress response determinants, all (24 of 24, 100%) were variably present within their given species, indicating that AMR potential can vary within *Macrococcus* species. Antimicrobial exposure can select for AMR (Hendriksen et al., 2019; Olesen et al., 2020), and reducing exposure (e.g., limiting antimicrobial use outside of treating human disease, minimizing unnecessary antibiotic use for human illness cases) can reduce the risk of AMR (Murray et al., 2022). Thus, it is not particularly surprising that intraspecies differences in AMR potential exist within *Macrococcus*; the genomes aggregated here were derived from *Macrococcus* strains isolated from a range of sources (e.g., humans, animal hosts, animal products, environmental samples), geographic locations (i.e., four continents), and timeframes (i.e., between the years of 1916 and 2021) and thus have likely been exposed to different selective pressures.

Comparatively, some genomic elements identified here were present across multiple *Macrococcus* spp., indicating shared inter-species functional potential for some phenotypes. Methicillin resistance genes *mecB* and *mecD*, for example, were variably present within multiple *Macrococcus* species (via GTDB-Tk; Figure 3), mirroring previous studies, which have reported *mecB* and/or *mecD* in various *Macrococcus* spp., including *M. caseolyticus* (Schwendener et al., 2017; MacFadyen et al., 2018; Zhang et al., 2022), *M. bohemicus* (Foster and Paterson, 2020), *M. canis* (Chancaithong et al., 2019), and *M. goetzii* (Maslanova et al., 2018). Outside of the AMR space, we further identified proteins that shared homology with virulence

factors in other species. Perhaps most notably, we detected homologues of aureolysin in multiple *Macrococcus* species (Figure 3). Aureolysin is an extracellular zinc-dependent metalloprotease secreted by *S. aureus*, which plays a crucial role in host immune system evasion (Thammavongsa et al., 2015; Pietrocola et al., 2017). While others have detected aureolysin homologues in *Macrococcus* genomes (Mazhar et al., 2019a,b; Zhang et al., 2022), the roles this protein plays in *Macrococcus* interactions with human or animal hosts (if any) are unknown.

Finally, for *M. caseolyticus* and *M. armenti*, which were composed of multiple populations (subclusters) separated by recent gene flow, some variably present genomic elements were subcluster-specific genes, which had been acquired post-speciation and differentially swept through these subclusters (i.e., flexible genes identified via PopCOGenT). Similar to results observed for *Ruminococcus gnavus* (Arevalo et al., 2019), transporter functions (e.g., ABC-type transporters, genes involved in ion transport) were enriched in subcluster-specific flexible gene sets within both *M. caseolyticus* and *M. armenti*. Notably, within *M. armenti*, we further identified two distinct subclusters with different flexible genes in each, including (i) one subcluster with a type VII secretion system, *S. aureus*-like virulence factors, and a putative pathogenicity island (Subcluster 2.0), and (ii) another with beta-lactamase family proteins and a type II toxin-antitoxin system (Subcluster 2.1). Taken together, these results indicate that there may be differences in the functional potential of these two *M. armenti* subclusters; however future experimental work will be needed to confirm the roles of these subcluster-specific flexible genes in each subcluster, as there are no clear differences in terms of each subcluster's ecological niche (strains in both subclusters were isolated from livestock in Switzerland).

Overall, proteins with potential virulence- and AMR-related functions, which were differentially present within and across *Macrococcus* species were identified. This indicates that there are potential within- and between-species differences in *Macrococcus* virulence and AMR potential. Future experimental efforts will thus be needed to investigate these differences further, as the study conducted here did not consider phenotypic data.

4.2. The lack of clear genomospecies boundaries between some *Macrococcus* species may cause taxonomic issues in the future

The delineation of prokaryotes into species-level taxonomic units is notoriously challenging, as horizontal gene transfer can obscure prokaryotic population boundaries (Jain et al., 2018; Arevalo et al., 2019). With the increasing availability of WGS, taxonomic assignment has largely shifted to *in silico* methods; however, numerous approaches exist for this purpose and may produce conflicting results (e.g., various implementations of ANI-based methods, marker gene-based methods, metrics using recent gene flow, *in silico* DNA–DNA hybridization; Meier-Kolthoff et al., 2013; Mende et al., 2013; Lee et al., 2016; Yoon et al., 2017; Jain et al., 2018; Arevalo et al., 2019; Chaumeil et al., 2019; Meier-Kolthoff et al., 2022; Parks et al., 2022). Here, we applied multiple species-level taxonomic assignment methods

to all publicly available *Macrococcus* genomes, specifically ANI-based approaches (i.e., OrthoANI/bactaxR and GTDB-Tk), an approach that uses a metric based on recent gene flow (i.e., PopCOGenT), and a marker gene-based method (i.e., specI; Figure 1). Overall, we observed similar results for three of four approaches; the marker gene-based approach only recovered two species, likely due to a lack of *Macrococcus* genomes of species other than *M. caseolyticus* and *M. canis* during species cluster database construction (this will likely be remedied in future specI database versions). However, even among the approaches that produced highly similar results, no two methods produced identical results.

Furthermore, at the conventional 95 ANI genomospecies threshold, several *Macrococcus* species were found to overlap (i.e., members of one species shared ≥ 95 ANI with members of a different species; Figure 2). We have previously observed a similar phenomenon among members of the *Bacillus cereus* group (Carroll et al., 2020), as others have done for *Escherichia/Shigella* spp., *Mycobacterium* spp., and *Neisseria gonorrhoeae/Neisseria meningitidis* (Jain et al., 2018). For *Macrococcus*, ambiguous species boundaries may not seem immediately concerning, as members of the genus are often viewed as animal commensals (Mazhar et al., 2018); thus, taxonomic misidentifications may not be viewed as “high consequence” compared to other organisms plagued by taxonomic issues (e.g., anthrax-causing organisms within the *Bacillus cereus* group, botulinum neurotoxin-producing members of *Clostridium*; Smith et al., 2018; Bower et al., 2022). However, as more *Macrococcus* strains undergo WGS and more is learned about the pathogenic potential of these organisms in animals and humans, there may be a greater need to ensure that species are clearly defined (e.g., in clinical laboratory, diagnostic, or regulatory settings). While there is some evidence that one of the South African genomes sequenced here belongs to a putative novel species (i.e., S115), we do not advocate for any changes to the taxonomy at this time, due to the limited number of genomes available. However, we encourage readers to be aware of ambiguous species boundaries for some *Macrococcus* spp., which may cause taxonomic issues in the future.

4.3. Future genomic sequencing, metadata collection, and phenotypic characterization efforts are needed to gain insight into *Macrococcus* population structure, antimicrobial resistance, and virulence potential

WGS has proven to be revolutionary in the food, veterinary, and human clinical microbiology spaces and is being used for—among other applications—pathogen surveillance, outbreak and cluster detection, source tracking, and diagnostics (Rossen et al., 2018; Brown et al., 2021; Ferdinand et al., 2021; Forde et al., 2023). Massive WGS efforts are being undertaken to query bacterial pathogens such as *Salmonella enterica*, *Escherichia coli*, and *Listeria monocytogenes* (Allard et al., 2016; Stevens et al., 2017; Brown et al., 2019), and large amounts of genomic data and metadata are publicly available for these organisms. As of 5 February 2023, (i) 455,330 genomes had been submitted to NCBI's GenBank Assembly database as *Salmonella*

enterica, (ii) 200,204 as *Escherichia coli*, and (iii) 51,579 as *Listeria monocytogenes*. *S. aureus* is a close relative of *Macrococcus* and boasts a total of 68,631 publicly available, assembled genomes (per NCBI's GenBank Assembly database, accessed 5 February 2023). These numbers dwarf those of *Macrococcus*, with 110 available, high-quality genomes for the entire genus at the time of this study, including the genomes generated here.

While our study provides insight into the evolution, population structure, and functional potential of *Macrococcus*, much more needs to be done to understand the role that *Macrococcus* spp. play as animal commensals, in animal-associated foodstuffs, and as opportunistic pathogens in animals and humans. First and foremost, future WGS efforts are needed to characterize these organisms, as increased availability of genomes will provide further insight into *Macrococcus* evolution (e.g., facilitating the identification of novel species, novel lineages within species). It is equally important that future WGS efforts are complemented with publicly available metadata (e.g., information conveying when and where a given strain was isolated) as this information can be used to identify potential host or geographic associations or potential migration or transmission events (e.g., between hosts or geographic regions). Finally, phenotypic data will be essential to confirm or invalidate the preliminary findings posited here regarding *Macrococcus* functional potential. Genomic AMR prediction, for example, does not necessarily translate to phenotypic AMR (Ransom et al., 2020). Similarly, any genomic determinants identified here based on their homology to known virulence factors (e.g., aureolysin, PVL, immune inhibitor A, the type VII secretion system identified in one *M. armentis* subcluster) must be evaluated experimentally. Thus, we hope that the results provided here can serve as a guide for further studies of the AMR and virulence potential of *Macrococcus* spp.

Data availability statement

The datasets presented in this study can be found in online repositories. The names of the repository/repositories and accession number(s) can be found in the article/Supplementary material.

References

- Alexa, A., Rahnenfuhrer, J., and Lengauer, T. (2006). Improved scoring of functional groups from gene expression data by decorrelating GO graph structure. *Bioinformatics* 22, 1600–1607. doi: 10.1093/bioinformatics/btl140
- Ali, D. E., Allam, M., Altayb, H. N., Mursi, D., Adalla, M. A., Mohammed, N. O., et al. (2022). A prevalence and molecular characterization of novel pathogenic strains of *Macrococcus caseolyticus* isolated from external wounds of donkeys in Khartoum state-Sudan. *BMC Vet. Res.* 18:197. doi: 10.1186/s12917-022-03297-2
- Allard, M. W., Strain, E., Melka, D., Bunning, K., Musser, S. M., Brown, E. W., et al. (2016). Practical value of food pathogen traceability through building a whole-genome sequencing network and database. *J. Clin. Microbiol.* 54, 1975–1983. doi: 10.1128/JCM.00081-16
- Arevalo, P., Vaninsberghe, D., Elsherbini, J., Gore, J., and Polz, M. F. (2019). A reverse ecology approach based on a biological definition of microbial populations. *Cell* 178:e814, 820–834.e14. doi: 10.1016/j.cell.2019.06.033
- Ashburner, M., Ball, C. A., Blake, J. A., Botstein, D., Butler, H., Cherry, J. M., et al. (2000). Gene Ontology: tool for the unification of biology. *Nat. Genet.* 25, 25–29. doi: 10.1038/75556
- Blin, K., Shaw, S., Kloosterman, A. M., Charlop-Powers, Z., Van Wezel, G. P., Medema, M. H., et al. (2021). antiSMASH 6.0: improving cluster detection and comparison capabilities. *Nucleic Acids Res.* 49, W29–W35. doi: 10.1093/nar/gkab335
- Bolger, A. M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 30, 2114–2120. doi: 10.1093/bioinformatics/btu170
- Bower, W. A., Hendricks, K. A., Vieira, A. R., Traxler, R. M., Weiner, Z., Lynfield, R., et al. (2022). What is Anthrax? *Pathogens* 11:690. doi: 10.3390/pathogens11060690
- Brown, B., Allard, M., Bazaco, M. C., Blankenship, J., and Minor, T. (2021). An economic evaluation of the whole genome sequencing source tracking program in the U.S. *PLoS One* 16:e0258262. doi: 10.1371/journal.pone.0258262
- Brown, E., Dessai, U., McGarry, S., and Gerner-Smidt, P. (2019). Use of whole-genome sequencing for food safety and public health in the United States. *Foodborne Pathog. Dis.* 16, 441–450. doi: 10.1089/fpd.2019.2662
- Buchfink, B., Xie, C., and Huson, D. H. (2015). Fast and sensitive protein alignment using DIAMOND. *Nat. Methods* 12, 59–60. doi: 10.1038/nmeth.3176
- Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., et al. (2009). BLAST+: architecture and applications. *BMC Bioinform.* 10:421. doi: 10.1186/1471-2105-10-421
- Cantalapiedra, C. P., Hernandez-Plaza, A., Letunic, I., Bork, P., and Huerta-Cepas, J. (2021). eggNOG-mapper v2: Functional Annotation, Orthology Assignments, and Domain Prediction at the Metagenomic Scale. *Mol. Biol. Evol.* 38, 5825–5829. doi: 10.1093/molbev/msab293
- Carroll, L. M., Larralde, M., Fleck, J. S., Ponnudurai, R., Milanese, A., Cappio, E., et al. (2021). Accurate *de novo* identification of biosynthetic gene clusters with GECCO. *bioRxiv:2021.2005.2003.442509*.

Author contributions

LC performed all computational analyses. IM performed bacterial isolation and identification as well as DNA extraction. RP supervised the sequencing of the isolates. IM and KM sourced the funding for sequencing of the isolates. All authors contributed to the article and approved the submitted version.

Funding

LC was supported by the SciLifeLab & Wallenberg Data Driven Life Science Program (grant: KAW 2020.0239). Additional funding was provided by the Gauteng Department of Agriculture and Rural Development.

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fmicb.2023.1181376/full#supplementary-material>

- Carroll, L. M., Wiedmann, M., and Kovac, J. (2020). Proposal of a taxonomic nomenclature for the *Bacillus cereus* group which reconciles genomic definitions of bacterial species with clinical and industrial phenotypes. *MBio* 11, e00034–e00020. doi: 10.1128/mBio.00034-20
- Chanchaithong, P., Perreten, V., and Schwendener, S. (2019). *Macrocooccus canis* contains recombinogenic methicillin resistance elements and the *mecB* plasmid found in *Staphylococcus aureus*. *J. Antimicrob. Chemother.* 74, 2531–2536. doi: 10.1093/jac/dkz260
- Chaumeil, P. A., Mussig, A. J., Hugenholtz, P., and Parks, D. H. (2019). GTDB-Tk: a toolkit to classify genomes with the Genome Taxonomy Database. *Bioinformatics*. doi: 10.1093/bioinformatics/btz848
- Cingolani, P., Platts, A., Wang, L. L., Coon, M., Nguyen, T., Wang, L., et al. (2012). A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly* 6, 80–92. doi: 10.4161/fly.19695
- Cotting, K., Strauss, C., Rodriguez-Campos, S., Rostaher, A., Fischer, N. M., Roosje, P. J., et al. (2017). *Macrocooccus canis* and *M. caseolyticus* in dogs: occurrence, genetic diversity and antibiotic resistance. *Vet. Dermatol.* 28, 559–e133. doi: 10.1111/vde.12474
- Croucher, N. J., Page, A. J., Connor, T. R., Delaney, A. J., Keane, J. A., Bentley, S. D., et al. (2015). Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res.* 43:e15. doi: 10.1093/nar/gku1196
- Evans, A. C. (1916). The Bacteria of Milk Freshly Drawn From Normal Udders. *J. Infect. Dis.* 18, 437–476. doi: 10.1093/infdis/18.5.437
- Ewels, P., Magnusson, M., Lundin, S., and Kaller, M. (2016). MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics* 32, 3047–3048. doi: 10.1093/bioinformatics/btw354
- Feldgarden, M., Brover, V., Haft, D. H., Prasad, A. B., Slotta, D. J., Tolstoy, I., et al. (2019). Validating the AMRFinder tool and resistance gene database by using antimicrobial resistance genotype-phenotype correlations in a collection of isolates. *Antimicrob. Agents Chemother.* 63, e00483–19. doi: 10.1128/AAC.00483-19
- Ferdinand, A. S., Kelaher, M., Lane, C. R., Da Silva, A. G., Sherry, N. L., Ballard, S. A., et al. (2021). An implementation science approach to evaluating pathogen whole genome sequencing in public health. *Genome Med.* 13:121. doi: 10.1186/s13073-021-00934-7
- Forde, B. M., Bergh, H., Cuddihy, T., Hajkowicz, K., Hurst, T., Playford, E. G., et al. (2023). Clinical implementation of routine whole-genome sequencing for hospital infection control of multi-drug resistant pathogens. *Clin. Infect. Dis.* 76, e1277–e1284. doi: 10.1093/cid/ciac726
- Foster, G., and Paterson, G. K. (2020). Methicillin-resistant *Macrocooccus bohemicus* encoding a divergent SCC*mecB* element. *Antibiotics* 9:590. doi: 10.3390/antibiotics9090590
- Fu, L., Niu, B., Zhu, Z., Wu, S., and Li, W. (2012). CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* 28, 3150–3152. doi: 10.1093/bioinformatics/bts565
- Garrison, E., and Marth, G. (2012). Haplotype-based variant detection from short-read sequencing. *arXiv:1207.3907*.
- Gobeli Brawand, S., Cotting, K., Gomez-Sanz, E., Collaud, A., Thomann, A., Brodard, I., et al. (2017). *Macrocooccus canis* sp. nov., a skin bacterium associated with infections in dogs. *Int. J. Syst. Evol. Microbiol.* 67, 621–626. doi: 10.1099/ijsem.0.001673
- Gomez-Sanz, E., Schwendener, S., Thomann, A., Gobeli Brawand, S., and Perreten, V. (2015). First staphylococcal cassette chromosome *mec* containing a *mecB*-carrying gene complex independent of transposon Tn6045 in a *Macrocooccus canis* isolate from a canine infection. *Antimicrob. Agents Chemother.* 59, 4577–4583. doi: 10.1128/AAC.05064-14
- Gurevich, A., Saveliev, V., Vyahhi, N., and Tesler, G. (2013). QUAST: quality assessment tool for genome assemblies. *Bioinformatics* 29, 1072–1075. doi: 10.1093/bioinformatics/btt086
- Hendriksen, R. S., Munk, P., Njage, P., van Bunnik, B., McNally, L., Lukjancenko, O., et al. (2019). Global monitoring of antimicrobial resistance based on metagenomics analyses of urban sewage. *Nat. Commun.* 10:1124. doi: 10.1038/s41467-019-08853-3
- Huerta-Cepas, J., Szklarczyk, D., Heller, D., Hernandez-Plaza, A., Forslund, S. K., Cook, H., et al. (2019). eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Res.* 47, D309–D314. doi: 10.1093/nar/gky1085
- Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T., and Aluru, S. (2018). High throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat. Commun.* 9:5114. doi: 10.1038/s41467-018-07641-9
- Jolley, K. A., Bray, J. E., and Maiden, M. C. J. (2018). Open-access bacterial population genomics: BIGSdb software, the PubMLST.org website and their applications. *Wellcome Open Res.* 3:124. doi: 10.12688/wellcomeopenres.14826.1
- Jolley, K. A., and Maiden, M. C. (2010). BIGSdb: scalable analysis of bacterial genome variation at the population level. *BMC Bioinform.* 11:595. doi: 10.1186/1471-2105-11-595
- Jost, G., Schwendener, S., Liassine, N., and Perreten, V. (2021). Methicillin-resistant *Macrocooccus canis* in a human wound. *Infect. Genet. Evol.* 96:105125. doi: 10.1016/j.meegid.2021.105125
- Kalyanamoorthy, S., Minh, B. Q., Wong, T. K. F., Von Haeseler, A., and Jeremiin, L. S. (2017). ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods* 14, 587–589. doi: 10.1038/nmeth.4285
- Katoh, K., and Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* 30, 772–780. doi: 10.1093/molbev/mst010
- Kautsar, S. A., Blin, K., Shaw, S., Navarro-Munoz, J. C., Terlouw, B. R., Van Der Hoof, J. J., et al. (2020). MIBiG 2.0: a repository for biosynthetic gene clusters of known function. *Nucleic Acids Res.* 48, D454–D458. doi: 10.1093/nar/gkz882
- Keller, J. E., Schwendener, S., Overesch, G., and Perreten, V. (2022). *Macrocooccus armenti* sp. nov., a novel bacterium isolated from the skin and nasal cavities of healthy pigs and calves. *Int. J. Syst. Evol. Microbiol.* 72. doi: 10.1099/ijsem.0.005245
- Kimura, M. (1981). Estimation of evolutionary distances between homologous nucleotide sequences. *Proc. Natl. Acad. Sci. U. S. A.* 78, 454–458. doi: 10.1073/pnas.78.1.454
- Kitts, P. A., Church, D. M., Thibaud-Nissen, F., Choi, J., Hem, V., Sapojnikov, V., et al. (2016). Assembly: a resource for assembled genomes at NCBI. *Nucleic Acids Res.* 44, D73–D80. doi: 10.1093/nar/gkv1226
- Kloos, W.E., Ballard, D.N., George, C.G., Webster, J.A., Hubner, R.J., Ludwig, W., Schleifer, K.H., Fiedler, F., and Schubert, K. (1998). Delimiting the genus *Staphylococcus* through description of *Macrocooccus caseolyticus* gen. nov., comb. nov. and *Macrocooccus equiperficus* sp. nov., and *Macrocooccus bovicus* sp. no. and *Macrocooccus carouselicus* sp. nov. *Int. J. Syst. Bacteriol.* 48 Pt, 859–877. doi: 10.1099/00207713-48-3-859
- Le, S. Q., and Gascuel, O. (2008). An improved general amino acid replacement matrix. *Mol. Biol. Evol.* 25, 1307–1320. doi: 10.1093/molbev/msn067
- Lee, I., Ouk Kim, Y., Park, S. C., and Chun, J. (2016). OrthoANI: an improved algorithm and software for calculating average nucleotide identity. *Int. J. Syst. Evol. Microbiol.* 66, 1100–1103. doi: 10.1099/ijsem.0.000760
- Letunic, I., and Bork, P. (2021). Interactive Tree of Life (iTOL) v5: an online tool for phylogenetic tree display and annotation. *Nucleic Acids Res.* 49, W293–W296. doi: 10.1093/nar/gkab301
- Li, H. (2011). A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* 27, 2987–2993. doi: 10.1093/bioinformatics/btr509
- Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv:1303.3997*.
- Li, H. (2019). "Seqtk: a fast and lightweight tool for processing sequences in the FASTA or FASTQ format". 1.2-r102-dirty ed.).
- Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25, 1754–1760. doi: 10.1093/bioinformatics/btp324
- Li, W., and Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* 22, 1658–1659. doi: 10.1093/bioinformatics/btl158
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., et al. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078–2079. doi: 10.1093/bioinformatics/btp352
- Liu, B., Zheng, D., Jin, Q., Chen, L., and Yang, J. (2019). VFDB 2019: a comparative pathogenomic platform with an interactive web interface. *Nucleic Acids Res.* 47, D687–D692. doi: 10.1093/nar/gky1080
- Löffler, B., Hussain, M., Grundmeier, M., Bruck, M., Holzinger, D., Varga, G., et al. (2010). *Staphylococcus aureus* Panton-Valentine leukocidin is a very potent cytotoxic factor for human neutrophils. *PLoS Pathog.* 6:e1000715. doi: 10.1371/journal.ppat.1000715
- Macfadyen, A. C., Fisher, E. A., Costa, B., Cullen, C., and Paterson, G. K. (2018). Genome analysis of methicillin resistance in *Macrocooccus caseolyticus* from dairy cattle in England and Wales. *Microb. Genom.* 4:e000191. doi: 10.1099/mgen.0.000191
- Mannerova, S., Pantucek, R., Doskar, J., Svec, P., Snauwaert, C., Vancanneyt, M., et al. (2003). *Macrocooccus brunensis* sp. nov., *Macrocooccus hajekii* sp. nov. and *Macrocooccus lamae* sp. nov., from the skin of llamas. *Int. J. Syst. Evol. Microbiol.* 53, 1647–1654. doi: 10.1099/ijms.0.02683-0
- Maslanova, I., Wertheimer, Z., Sedlacek, I., Svec, P., Indrakova, A., Kovarovic, V., et al. (2018). Description and comparative genomics of *Macrocooccus caseolyticus* subsp. *hominis* subsp. nov., *Macrocooccus goetzii* sp. nov., *Macrocooccus epidermidis* sp. nov., and *Macrocooccus bohemicus* sp. nov., novel macrococci from human clinical material with virulence potential and suspected uptake of foreign DNA by natural transformation. *Front. Microbiol.* 9:1178. doi: 10.3389/fmicb.2018.01178
- Mazhar, S., Altermann, E., Hill, C., and Mcauliffe, O. (2019a). Draft genome sequences of *Macrocooccus caseolyticus*, *Macrocooccus canis*, *Macrocooccus bohemicus*, and *Macrocooccus goetzii*. *Microbiol. Resour. Announc.* 8, e00343–19. doi: 10.1128/MRA.00343-19
- Mazhar, S., Altermann, E., Hill, C., and Mcauliffe, O. (2019b). Draft genome sequences of the type strains of six *Macrocooccus* species. *Microbiol. Resour. Announc.* 8, e00343–19. doi: 10.1128/MRA.00344-19
- Mazhar, S., Hill, C., and Mcauliffe, O. (2018). The genus *Macrocooccus*: an insight into its biology, evolution, and relationship with *Staphylococcus*. *Adv. Appl. Microbiol.* 105, 1–50. doi: 10.1016/bs.aambs.2018.05.002
- Meier-Kolthoff, J. P., Auch, A. F., Klenk, H.-P., and Göker, M. (2013). Genome sequence-based species delimitation with confidence intervals and improved distance functions. *BMC Bioinform.* 14:60. doi: 10.1186/1471-2105-14-60

- Meier-Kolthoff, J. P., Carbasse, J. S., Peinado-Olarte, R. L., and Goker, M. (2022). TYGS and LPSN: a database tandem for fast and reliable genome-based classification and nomenclature of prokaryotes. *Nucleic Acids Res.* 50, D801–D807. doi: 10.1093/nar/gkab902
- Mende, D. R., Sunagawa, S., Zeller, G., and Bork, P. (2013). Accurate and universal delineation of prokaryotic species. *Nat. Methods* 10, 881–884. doi: 10.1038/nmeth.2575
- Milanese, A., Mende, D. R., Paoli, L., Salazar, G., Ruscheweyh, H. J., Cuenca, M., et al. (2019). Microbial abundance, activity and population genomic profiling with mOTUs2. *Nat. Commun.* 10:1014. doi: 10.1038/s41467-019-08844-4
- Minh, B. Q., Nguyen, M. A., and Von Haeseler, A. (2013). Ultrafast approximation for phylogenetic bootstrap. *Mol. Biol. Evol.* 30, 1188–1195. doi: 10.1093/molbev/mst024
- Murray, C. J. L., Ikuta, K. S., Sharara, F., Swetschinski, L., Robles Aguilar, G., Gray, A., et al. (2022). Global burden of bacterial antimicrobial resistance in 2019: a systematic analysis. *Lancet* 399, 629–655. doi: 10.1016/S0140-6736(21)02724-0
- Navarro-Munoz, J. C., Selem-Mojica, N., Mullowney, M. W., Kautsar, S. A., Tryon, J. H., Parkinson, E. I., et al. (2020). A computational framework to explore large-scale biosynthetic diversity. *Nat. Chem. Biol.* 16, 60–68. doi: 10.1038/s41589-019-0400-9
- Nguyen, L. T., Schmidt, H. A., Von Haeseler, A., and Minh, B. Q. (2015). IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* 32, 268–274. doi: 10.1093/molbev/msu300
- Olesen, S. W., Lipsitch, M., and Grad, Y. H. (2020). The role of "spillover" in antibiotic resistance. *Proc. Natl. Acad. Sci. U. S. A.* 117, 29063–29068. doi: 10.1073/pnas.2013694117
- Ouoba, L. I. I., Vouidibio Mbozo, A. B., Anyogu, A., Obioha, P. I., Lingani-Sawadogo, H., Sutherland, J. P., et al. (2019). Environmental heterogeneity of *Staphylococcus* species from alkaline fermented foods and associated toxins and antimicrobial resistance genetic elements. *Int. J. Food Microbiol.* 311:108356. doi: 10.1016/j.ijfoodmicro.2019.108356
- Page, A. J., Taylor, B., Delaney, A. J., Soares, J., Seemann, T., Keane, J. A., et al. (2016). SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments. *Microb. Genom.* 2:e000056. doi: 10.1099/mgen.0.000083
- Paradis, E., Claude, J., and Strimmer, K. (2004). APE: Analyses of Phylogenetics and Evolution in R language. *Bioinformatics* 20, 289–290. doi: 10.1093/bioinformatics/btg412
- Paradis, E., and Schliep, K. (2019). ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R. *Bioinformatics* 35, 526–528. doi: 10.1093/bioinformatics/bty633
- Parks, D. H., Chuvochina, M., Rinke, C., Mussig, A. J., Chaumeil, P. A., and Hugenholtz, P. (2022). GTDB: an ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and complete genome-based taxonomy. *Nucleic Acids Res.* 50, D785–D794. doi: 10.1093/nar/gkab776
- Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P., and Tyson, G. W. (2015). CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.* 25, 1043–1055. doi: 10.1101/gr.186072.114
- Parte, A. C., Sarda Carbasse, J., Meier-Kolthoff, J. P., Reimer, L. C., and Goker, M. (2020). List of prokaryotic names with standing in nomenclature (LPSN) moves to the DSMZ. *Int. J. Syst. Evol. Microbiol.* 70, 5607–5612. doi: 10.1099/ijsem.0.004332
- Pietrocola, G., Nobile, G., Rindi, S., and Speziale, P. (2017). *Staphylococcus aureus* manipulates innate immunity through own and host-expressed proteases. *Front. Cell. Infect. Microbiol.* 7:166. doi: 10.3389/fcimb.2017.00166
- Poyart, C., Quesne, G., Boumaila, C., and Trieu-Cuot, P. (2001). Rapid and accurate species-level identification of coagulase-negative staphylococci by using the *sodA* gene as a target. *J. Clin. Microbiol.* 39, 4296–4301. doi: 10.1128/JCM.39.12.4296-4301.2001
- Price, M. N., Dehal, P. S., and Arkin, A. P. (2010). FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS One* 5:e9490. doi: 10.1371/journal.pone.0009490
- Quinlan, A. R., and Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 26, 841–842. doi: 10.1093/bioinformatics/btq033
- R Core Team (2021). "R: A language and environment for statistical computing". 4.1.2 ed. (Vienna, Austria: R Foundation for Statistical Computing).
- Ramos, G., Vigoder, H. C., and Nascimento, J. S. (2021). Technological applications of *Macrococcus caseolyticus* and its impact on food safety. *Curr. Microbiol.* 78, 11–16. doi: 10.1007/s00284-020-02281-z
- Ransom, E. M., Potter, R. F., Dantas, G., and Burnham, C. D. (2020). Genomic prediction of antimicrobial resistance: ready or not, here it comes! *Clin. Chem.* 66, 1278–1289. doi: 10.1093/clinchem/hvaa172
- Rossen, J. W. A., Friedrich, A. W., Moran-Gilad, J., Genomic, E. S. G. F., and Molecular, D. (2018). Practical issues in implementing whole-genome-sequencing in routine diagnostic microbiology. *Clin. Microbiol. Infect.* 24, 355–360. doi: 10.1016/j.cmi.2017.11.001
- Schoch, C. L., Ciufo, S., Domrachev, M., Hotton, C. L., Kannan, S., Khovanskaya, R., et al. (2020). NCBI Taxonomy: a comprehensive update on curation, resources and tools. *Database* 2020: baaa062. doi: 10.1093/database/baaa062
- Schwendener, S., Cotting, K., and Perreten, V. (2017). Novel methicillin resistance gene *mecD* in clinical *Macrococcus caseolyticus* strains from bovine and canine sources. *Sci. Rep.* 7:43797. doi: 10.1038/srep43797
- Schwendener, S., and Perreten, V. (2022). The *bla* and *mec* families of beta-lactam resistance genes in the genera *Macrococcus*, *Mammaliococcus* and *Staphylococcus*: an in-depth analysis with emphasis on *Macrococcus*. *J. Antimicrob. Chemother.* 77, 1796–1827. doi: 10.1093/jac/dkac107
- Seemann, T. (2014). Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30, 2068–2069. doi: 10.1093/bioinformatics/btu153
- Seemann, T. (2019). "samclip: filter SAM file for soft and hard clipped alignments". 0.2 ed.
- Shallcross, L. J., Fragaszy, E., Johnson, A. M., and Hayward, A. C. (2013). The role of the Pantone-Valentine leucocidin toxin in staphylococcal disease: a systematic review and meta-analysis. *Lancet Infect. Dis.* 13, 43–54. doi: 10.1016/S1473-3099(12)70238-4
- Simonsen, M., Mailund, T., and Pedersen, C. N. S. (2008). "Rapid neighbour-joining" in *Algorithms in bioinformatics. Lecture Notes in Computer Science*(vol. 5251, eds. K. A. Crandall and J. Lagergren (Springer Berlin Heidelberg), 113–122.
- Smith, T., Williamson Charles, H. D., Hill, K., Sahl, J., and Keim, P. (2018). Botulinum Neurotoxin-Producing Bacteria. Isn't it Time That We Called a Species a Species? *MBio* 9, e01469–e01418. doi: 10.1128/mBio.01469-18
- Soubrier, J., Steel, M., Lee, M. S. Y., Der Sarkissian, C., Guindon, S., Ho, S. Y. W., et al. (2012). The Influence of Rate Heterogeneity Among Sites on the Time Dependence of Molecular Rates. *Mol. Biol. Evol.* 29, 3345–3358. doi: 10.1093/molbev/mss140
- Souvorov, A., Agarwala, R., and Lipman, D. J. (2018). SKESA: strategic k-mer extension for scrupulous assemblies. *Genome Biol.* 19:153. doi: 10.1186/s13059-018-1540-z
- Steinberger, M., and Soding, J. (2017). MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat. Biotechnol.* 35, 1026–1028. doi: 10.1038/nbt.3988
- Stevens, E. L., Timme, R., Brown, E. W., Allard, M. W., Strain, E., Bunning, K., et al. (2017). The public health impact of a publicly available, environmental database of microbial genomes. *Front. Microbiol.* 8:808. doi: 10.3389/fmicb.2017.00808
- Tan, A., Abecasis, G. R., and Kang, H. M. (2015). Unified representation of genetic variants. *Bioinformatics* 31, 2202–2204. doi: 10.1093/bioinformatics/btv112
- Tavaré, S. (1986). Some probabilistic and statistical problems in the analysis of DNA sequences. *Lect. Math. Life Sci.* 17, 57–86.
- Tettelin, H., Riley, D., Cattuto, C., and Medini, D. (2008). Comparative genomics: the bacterial pan-genome. *Curr. Opin. Microbiol.* 11, 472–477. doi: 10.1016/j.mib.2008.09.006
- Thammavongsa, V., Kim, H. K., Missiakas, D., and Schneewind, O. (2015). Staphylococcal manipulation of host immune responses. *Nat. Rev. Microbiol.* 13, 529–543. doi: 10.1038/nrmicro3521
- The Gene Ontology Consortium (2018). The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic Acids Res.* 47, D330–D338. doi: 10.1093/nar/gky1055
- Tonkin-Hill, G., Lees, J. A., Bentley, S. D., Frost, S. D. W., and Corander, J. (2018). RhierBAPS: an R implementation of the population clustering algorithm hierBAPS. *Wellcome Open Res.* 3:93. doi: 10.12688/wellcomeopenres.14694.1
- Tonkin-Hill, G., Macalasdair, N., Ruis, C., Weimann, A., Horesch, G., Lees, J. A., et al. (2020). Producing polished prokaryotic pangenomes with the Panaroo pipeline. *Genome Biol.* 21:180. doi: 10.1186/s13059-020-02090-4
- Tshipamba, M. E., Lubanza, N., Adetunji, M. C., and Mwanza, M. (2018). Molecular characterization and antibiotic resistance of foodborne pathogens in street-vended ready-to-eat meat sold in South Africa. *J. Food Prot.* 81, 1963–1972. doi: 10.4315/0362-028X.JFP-18-069
- Vaitkevicius, K., Rompikuntal, P. K., Lindmark, B., Vaitkevicius, R., Song, T., and Wai, S. N. (2008). The metalloprotease PrtV from *Vibrio cholerae*. *FEBS J.* 275, 3167–3177. doi: 10.1111/j.1742-4658.2008.06470.x
- Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
- Yang, Z. (1995). A space-time process model for the evolution of DNA sequences. *Genetics* 139, 993–1005. doi: 10.1093/genetics/139.2.993
- Yoon, S. H., Ha, S. M., Lim, J., Kwon, S., and Chun, J. (2017). A large-scale evaluation of algorithms to calculate average nucleotide identity. *Antonie Van Leeuwenhoek* 110, 1281–1286. doi: 10.1007/s10482-017-0844-4
- Zhang, Y., Min, S., Sun, Y., Ye, J., Zhou, Z., and Li, H. (2022). Characteristics of population structure, antimicrobial resistance, virulence factors, and morphology of methicillin-resistant *Macrococcus caseolyticus* in global clades. *BMC Microbiol.* 22:266. doi: 10.1186/s12866-022-02679-8
- Zhou, Z., Charlesworth, J., and Achtman, M. (2020). Accurate reconstruction of bacterial pan- and core genomes with PEPPAN. *Genome Res.* 30, 1667–1679. doi: 10.1101/gr.260828.120