




# Comparative Genomics Suggests a Taxonomic Revision of the *Staphylococcus cohnii* Species Complex

Anna Lavecchia <sup>1,†</sup>, Matteo Chiara <sup>1,2,†</sup>, Caterina De Virgilio<sup>3</sup>, Caterina Manzari<sup>1</sup>, Carlo Pazzani<sup>4</sup>, David Horner<sup>1,2</sup>, Graziano Pesole <sup>1,3,5</sup>, and Antonio Placido<sup>1,\*</sup>

<sup>1</sup>Institute of Biomembranes, Bioenergetics and Molecular Biotechnologies, National Research Council, Bari, Italy

<sup>2</sup>Department of Biosciences, University of Milan, Milan, Italy

<sup>3</sup>Department of Biosciences, Biotechnology and Biopharmaceutics, University of Bari “Aldo Moro”, Bari, Italy

<sup>4</sup>Department of Biology, University of Bari Aldo Moro, Bari, Italy

<sup>5</sup>Consorzio Interuniversitario Biotecnologie, Trieste, Italy

†These authors contributed equally to this work.

\*Corresponding author: E-mail: a.placido@ibiom.cnr.it.

Accepted: 22 January 2021

## Abstract

*Staphylococcus cohnii* (SC), a coagulase-negative bacterium, was first isolated in 1975 from human skin. Early phenotypic analyses led to the delineation of two subspecies (subsp.), *Staphylococcus cohnii* subsp. *cohnii* (SCC) and *Staphylococcus cohnii* subsp. *urealyticus* (SCU). SCC was considered to be specific to humans, whereas SCU apparently demonstrated a wider host range, from lower primates to humans. The type strains ATCC 29974 and ATCC 49330 have been designated for SCC and SCU, respectively. Comparative analysis of 66 complete genome sequences—including a novel SC isolate—revealed unexpected patterns within the SC complex, both in terms of genomic sequence identity and gene content, highlighting the presence of 3 phylogenetically distinct groups. Based on our observations, and on the current guidelines for taxonomic classification for bacterial species, we propose a revision of the SC species complex. We suggest that SCC and SCU should be regarded as two distinct species: SC and SU (*Staphylococcus urealyticus*), and that two distinct subspecies, SCC and SCB (SC subsp. *barensis*, represented by the novel strain isolated in Bari) should be recognized within SC. Furthermore, since large-scale comparative genomics studies recurrently suggest inconsistencies or conflicts in taxonomic assignments of bacterial species, we believe that the approach proposed here might be considered for more general application.

**Key words:** comparative genomics, genome shotgun sequencing, DNA–DNA hybridization analyses, average nucleotide identity, phylogenetic analyses, *Staphylococcus cohnii*.

## Significance

In recent years, as the extent of its involvement in a multitude of human and animal infections has become evident, much research has been focused on *Staphylococcus cohnii*. Moreover, *S. cohnii* is also widely used as a model in the development of biotechnological applications and antibacterial medical devices. Its relevance notwithstanding, comparative genomic studies of this species complex are lacking, and the current work suggests the need for a major taxonomic revision. More generally, the computational approach presented here can be applied on a larger scale, both for the resolution of complex or conflicting taxonomic assignments and as a tool to contribute to the understanding of the history of biomedically and biotechnologically important traits at species and subspecies levels.

© The Author(s) 2021. Published by Oxford University Press on behalf of the Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact [journals.permissions@oup.com](mailto:journals.permissions@oup.com)

## Introduction

*Staphylococcus cohnii* (SC), a Gram-positive bacterium of the Coagulase-Negative Staphylococci (CoNS) group, was first isolated by Schleifer and Kloos in 1975. The name *cohnii* was adopted in memory of Ferdinand Julius Cohn, a German botanist and bacteriologist (Schleifer and Kloos 1975).

According to the current classification—which is fundamentally based on phenotypic traits—SC includes two subspecies (subsp.): *Staphylococcus cohnii* subsp. *cohnii* (SCC) and *Staphylococcus cohnii* subsp. *urealyticum* (SCU). The SCC ATCC 29974 and SCU ATCC 49330 isolates have been designated as the type strains for SCC and SCU, respectively (Kloos and Wolfshohl 1983, 1991). The original spelling, SC subsp. *urealytium* (sic), was corrected by Sneath to SC subsp. *urealyticus* (SCU) in 1992 (Sneath 1992).

A larger colony size, distinct pigmentation, differences in fatty acid profile, and the presence of metabolic activities, including  $\beta$ -glucuronidase and  $\beta$ -galactosidase activities, delayed alkaline phosphatase activity, and the ability to produce acid aerobically from  $\alpha$ -lactose discriminate SCU from SCC. Moreover, SCC was originally reported to colonize only humans, whereas SCU can also colonize other primates (Kloos and Wolfshohl 1991). More recently, SCU has also been isolated from healthy dogs (Bean et al. 2017) and goats (Seni et al. 2019).

Similar to most CoNS, SCC and SCU are typically commensal bacteria of the skin and mucous membranes (Waldon et al. 2002; Crossley et al. 2009). However, several opportunistically pathogenic strains have been described and implicated in nosocomial infections, including meningitis, primary septic arthritis, septicemia, brain abscess, and catheter invasion (Okudera et al. 1991; Mastroianni et al. 1996; Basaglia et al. 2003; Yamashita et al. 2005; Adeyemi et al. 2010; Mendoza-Olazarán et al. 2017). The ability to form biofilms appears to play an important role in staphylococcal virulence (Yong et al. 2019), and biofilm-associated infections are of particular concern because they are often difficult to resolve with antibiotics. Recent studies suggest that similar to other CoNS, SC can adhere to and invade human HeLa cells through the formation of biofilms (Szczyka et al. 2016), whereas strains of SC and especially those isolated from hospital environments, including pediatric wards and intensive-care units, have been reported to be resistant to several antibiotics (Szewczyk et al. 2000, 2004; Song et al. 2017). Indeed, multidrug-resistant CoNS bacteria constitute an emerging source of concern for Public Health Organizations (David and Elliott 2015; Moawad et al. 2019) as they have been associated with an increasing proportion of nosocomial infections, and because they can act as a reservoir of resistance determinants for *Staphylococcus aureus* through

horizontal gene transfer (Otto 2013; Winstel et al. 2013; Larsen et al. 2017; Argemi et al. 2019).

Comparative genomic and phylogenetic studies allow the characterization of evolutionary dynamics and the identification of genes and pathways potentially involved in pathogenesis and/or antibiotic resistance. Here we present comprehensive analyses of the complete collection of 65 publicly available SC genomes as well as that of a novel strain. We uncover striking patterns of genomic evolution, including high levels of genomic diversity and differential gene acquisition and loss, which suggest a taxonomic revision of the SC species complex. We propose that SC should be divided into two species, SC and SU (*Staphylococcus urealyticus*). Moreover, two subspecies SCC and SCB (SC subsp. *barensis*, exemplified by our novel strain isolated in Bari) should be distinguished within SC. Of more general interest, we describe an approach based on the integration of different types of phylogenetic, genomic and gene content analyses that could be applied on a larger scale for the resolution of complex or conflicting taxonomic assignments.

## Materials and Methods

### Isolation of SC5 and Preliminary Taxonomic Detection

A Nunc bioassay plate containing Luria Bertani (LB) agar supplemented with 0.2 mM Potassium Chromate and Chloramphenicol (12.5  $\mu$ g/ml) was prepared and placed on the worktop of a class II biological safety cabinet. Airflow was switched on for 1 h, and the plate was subsequently incubated for 16 h at 37°C. Sixty-two colonies were grown.

All colonies were replicated on LB agar containing Cr<sup>+6</sup> 150 mM. Only five colonies survived (Lavecchia et al. 2018) and were subjected to a preliminary taxonomic characterization through partial 16S rDNA amplification and Sanger sequencing. One colony preliminarily assigned to *Staphylococcus cohnii* and labeled SC5 was then subjected to whole-genome sequencing.

Polymerase chain reaction (PCR) amplification of the SC5 16S rDNA was performed using 200 ng of DNA (see DNA Isolation section) and the following primers: 1  $\mu$ M Forward 5'-TACGGGAGGCAGCAGTAG-3' (16S rDNA position 369-386), 1  $\mu$ M Reverse 5'-CATGGTGTGACGGGCG GT-3' (position 1424-1441). The reaction mixture (final volume 50  $\mu$ l) was completed using: each NTP 200  $\mu$ M, 2 mM MgCl<sub>2</sub>, and 2 U Taq DNA polymerase (Thermo Fisher Scientific, Waltham, MA). Taq DNA polymerase was activated as follows: 94°C for 5 min; 30 cycles at 94°C for 30 s, 55°C for 15 s, and 72°C for 1 min.

Sanger sequencing was performed by Macrogen (Amsterdam, Netherlands), and the preliminary taxonomic assessment was made by probing the 16S rRNA (Bacteria and Archaea) database available at the NCBI (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>), using the BlastN algorithm.

### DNA Isolation, Library Preparation, and Sequencing

Genomic DNA was extracted with the DNeasy Blood and Tissue kit (Qiagen, Hilden, Germany). The library was prepared using the Nextera XT library prep workflow (Illumina, Illumina, San Diego, CA) and 2 × 250 nt paired-end reads generated on an Illumina MiSeq instrument.

### Genome Assembly and Annotation

Raw data were processed using a modified version of the “Fosmid1” pipeline in the A-GAME (Chiara et al. 2018) Galaxy framework (Afgan et al. 2018). Quality trimming was executed using the sliding-window operation in Trimmomatic with default parameters (Bolger et al. 2014). Overlapping reads were merged using PEAR with standard parameters (Zhang et al. 2014). The final assembly was performed using the SPAdes assembler (version 3.50) using kmers of 33, 55, 77, 99, and 121 nt (Bankevich et al. 2012). Annotation was performed with PROKKA using default parameters (Seemann 2014).

### *Staphylococcus cohnii* Genomes Used in This Study and Annotation of Protein-Coding Genes

The complete collection of 65 *S. cohnii* (SC) genome assemblies (including SC subsp. *cohnii* ATCC 29974 and SC subsp. *urealyticus* ATCC 49330), as available in GenBank on July 1, 2020, was downloaded from the NCBI assembly database, directly from the “Download Assemblies” link, as available from the web interface. To avoid possible ascertainment biases, all the genomes were reannotated using the procedure described above. Annotations of protein-coding genes, as obtained from Prokka, were used in all the subsequent analyses. A complete list of the accession number of the genomes used in this study is provided in [supplementary table S1](#), [Supplementary Material](#) online.

### Calculation of Average Nucleotide Identity and *In Silico* DNA–DNA Hybridization

Average Nucleotide Identity based on BLAST (ANIb) between all 66 genomes (SC5 included) was computed according to the method described by Rossello-Mora (Richter and Rosselló-Móra 2009), as implemented by a custom script, which is available at [https://github.com/cvulpispaper/compute\\_anib](https://github.com/cvulpispaper/compute_anib). *In silico* DNA–DNA hybridization (DDH) was computed using the GGDC (Genome-to-Genome Distance Calculator 2.1) available from <https://ggdc.dsmz.de/ggdc.php#> (Meier-Kolthoff et al. 2013). As recommended in Auch et al. (2010), all the comparisons performed in this study were based on formula 2.

### Clusters of Orthologous Genes, Core, and Accessory Genome

The makeblastdb utility, as incorporated in the blast+ software package (Camacho et al. 2009) was used to prepare a Blast protein database, containing all of the protein-coding genes as predicted by Prokka in the 66 SC genomes included in this study. All-against-all BLASTp (Altschul et al. 1990) was performed using the BLOSUM80 matrix and accepting only best reciprocal hits with  $e\text{-value} \leq 1e-5$  and where “second-best” hits from the same genome produce bit scores <90% of that associated with the best match. Putative clusters of orthologous genes (COGs) were established as groups of best reciprocal BLAST hits. Core genes were defined as COGs containing single representatives from all genomes included in our analyses (or all genomes within major groups) and accessory genes as COGs with incomplete representation. The program used for the identification of COGs is available at [https://github.com/cvulpispaper/compute\\_anib](https://github.com/cvulpispaper/compute_anib).

### Estimation of Completeness of Core and Accessory Genomes

Size of core and accessory genomes were established by rarefaction analyses based on random resampling of genomic sequences of each major group. For each number of strains considered (2–28 for B, 2–22 for A1, 2–16 for A2, and 2–66 for SC) the inferred sizes of core and accessory genomes were recorded for 10,000 replicates of randomly selected combinations of genomes. Plots were prepared showing mean and standard deviation of these statistics.

### Phyletic Patterns and Clustering of Gene Presence/Absence Profiles

The phyletic pattern of genes presence/absence in the genomes of the 66 SC isolates was inferred directly by comparison of clusters of orthologous genes. Only COGs containing ten or more genes were considered in this analysis. A matrix of gene presence/absence was compiled, with genes on the rows and isolates on the columns. A value of 0 was used to indicate the absence of a gene, a value of 1 its presence. A correlation-based distance matrix of gene presence/absence profiles was obtained by applying the cor and the dist functions, from the stat library of the R programming language with default parameters (Pearson correlation and Euclidean distances, respectively). Clustering was performed by applying the hclust function with median linkage, from the same software package.

### Phylogenetic Analyses

The conceptually translated sequences of the 1468 SC core genes were independently aligned using Muscle (Edgar 2004) and ambiguously aligned regions were excluded using the GBlocks software (Castresana 2000). Maximum-likelihood

**Table 1.**

Main Genome Assembly Features of SC5, SCC ATCC 29974 and SCU ATCC 49330 Strains

Strain	Size (Mb)	GC (%)	Contigs	N50 (kb)	Proteins	rRNAs	tRNAs	DDH (%) <sup>a</sup>			ANiB (%) <sup>b</sup>		
								SC5	SCC	SCU	SC5	SCC	SCU
SC5	2.62	32.2	34	806	2,510	10	61	—	67.5	41.9	—	95.4	91.0
SCC ATCC 29974	2.71	32.6	83	114	2,422	9	58	67.5	—	—	95.4	—	91.5
SCU ATCC 49330	2.67	32.5	223	26	2,457	13	61	41.9	—	—	91.0	91.5	—

NOTE.—DDH and ANiB values between strains.

<sup>a</sup> DDH (cut-off for species affiliation > 70%); <sup>b</sup> ANiB (cut-off for species affiliation > 96%). SC5, our isolate; SCC ATCC 29974, *Staphylococcus cohnii* subsp. *cohnii* (type strain) and SCU ATCC 49330, *Staphylococcus cohnii* subsp. *urealyticus* (type strain).

phylogenetic reconstruction and bootstrap analyses of concatenated alignments were performed using the software PHYML (Guindon et al. 2009) under the WAG (Whelan and Goldman 2001) substitution model, suggested by the software ProtTest (Darriba et al. 2011) to best fit the data, with invariable and four gamma-distributed substitution rate categories.

### Statistical Analyses

Welch *t*-test *P*-values for the comparison of ANiB distributions and the size of the core and accessory genomes were computed by means of the *t*-test function as implemented in the stats R package.

## Results

### Isolation and Whole-Genome Shotgun Sequencing of a Novel Strain of *Staphylococcus cohnii*

Five strains of staphylococci were isolated from a disused class II biological safety cabinet during a study aiming to identify bacterial strains resistant to hexavalent chromate (Lavecchia et al. 2018). Preliminary taxonomic analyses based on partial 16S rDNA Sanger sequencing identified four of these strains as *Staphylococcus arlettae* (Lavecchia et al. 2018), whereas one isolate showed high levels of similarity (99.1% and 98.8%), respectively, with the type strains of SCC (*Staphylococcus cohnii* subsp. *cohnii*) and SCU (*Staphylococcus cohnii* subsp. *urealyticus*).

The latter, preliminarily named SC5, was subsequently subjected to Whole-Genome Shotgun Sequencing using an Illumina MiSeq instrument. A total of 2,402,324 paired-end reads were obtained, with an average insert size of 254.74 bp, providing a theoretical 230× coverage of the genome. Raw reads were subjected to quality trimming and assembly by means of a modified version of the Fosmid1 pipeline as incorporated in A-GAME (Chiara et al. 2018). Salient features of the SC5 genome assembly are summarized in table 1. An overall good level of contiguity was observed with more than 90% of the assembly incorporated in contigs >100 kb in size (N90 108 kb). A total of 10 rRNA, 61 tRNA, and 2,510 protein-coding genes were predicted by in silico

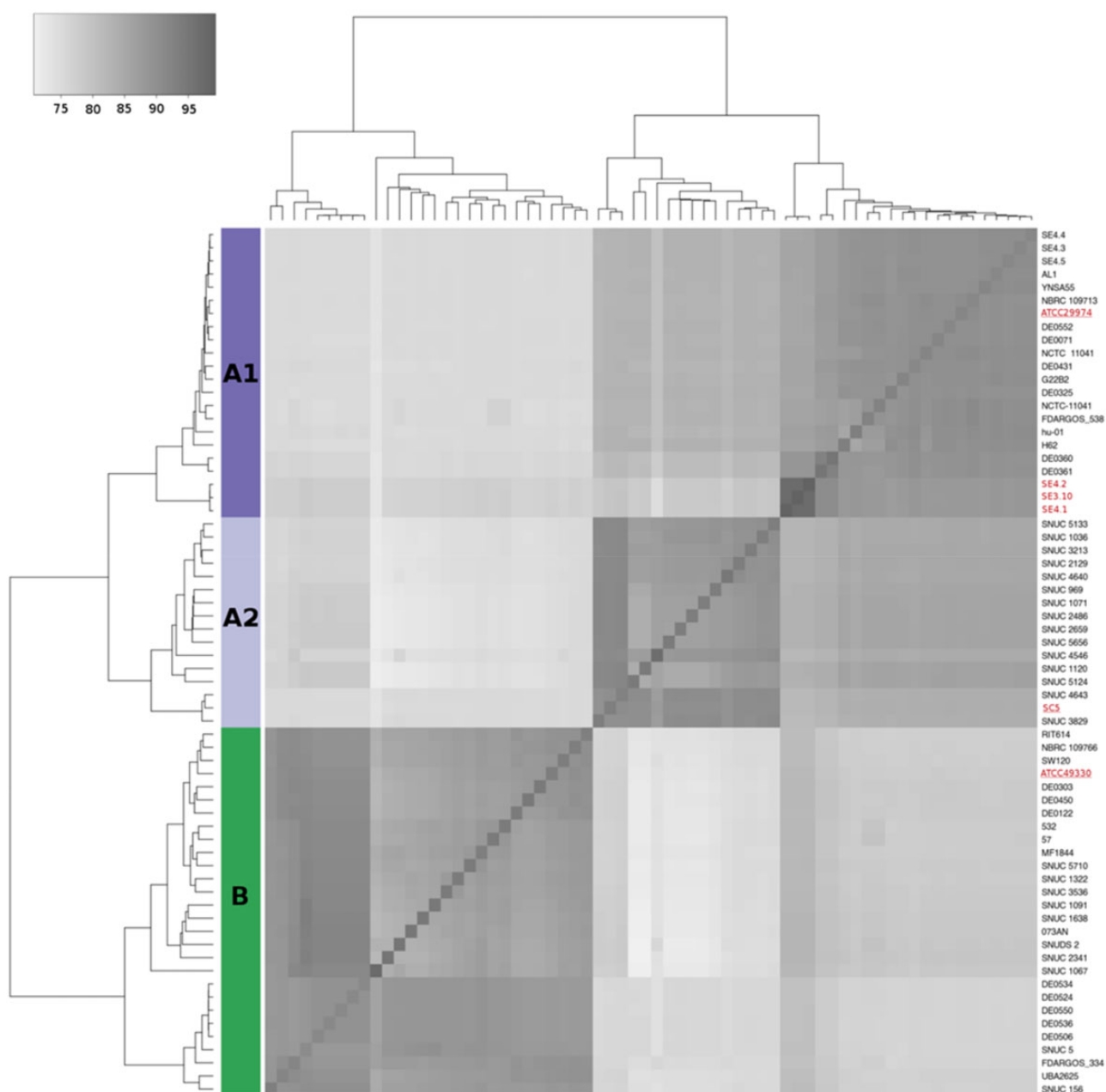
annotation of the genome. Of note, the *emrA* and *emrB* genes, implicated in chromate and ampicillin co-resistance in *Staphylococcus aureus* LZ 01 (Zhang et al. 2016), were identified in the genome of SC5. These genes are also observed in the genomes of the other four *S. arlettae* chromium-resistant strains isolated from the same environment. The draft genome sequence of SC5 was deposited in NCBI under the accession number JAALCY000000000, BioSample accession number SAMN14142771, and BioProject number ID PRJNA607668.

### In Silico DNA-DNA Hybridization Analyses

In silico DNA-DNA hybridization (DDH) analyses were performed to refine the taxonomic delineation of SC5. Strikingly, although 16S rDNA taxonomic assignment suggested that SC5 was closely related to SCC, in silico hybridization assays against the SCC ATCC 29974 and SCU ATCC 49330 type strains recovered somewhat unexpected patterns. Indeed (table 1), the observed DDH values, 67.5% and 41.9%, respectively, for SCC and SCU, were borderline or well below the cut-off value of 70%, that is, normally used to delineate species by this method (Meier-Kolthoff et al. 2013; Colston et al. 2014; Garrido-Sanz et al. 2016). A systematic comparison of in silico hybridization profiles of SC5 against the complete collection of the other 65 SC draft genomes considered in this study (supplementary table S1, Supplementary Material online) was performed. Notably, contrasting patterns of sequence similarity profiles (supplementary table S2, Supplementary Material online) were observed. SC5 showed DDH levels > 90% with 15 isolates (supplementary fig.S1, Supplementary Material online), but lower than 70% with 22 (including the SCC ATCC 29974 type strain) and around 45% with the remaining 28 isolates (including the SCU ATCC 49330 type strain).

### Analysis of Genomic Identity

Levels of pairwise genome identity between all currently available SC genomes were established by Average Nucleotide Identity on BLAST (ANiB). Hierarchical clustering of ANiB profiles was applied to identify clades/groups of SC genomes with similar levels of genome identity. As shown in figure 1

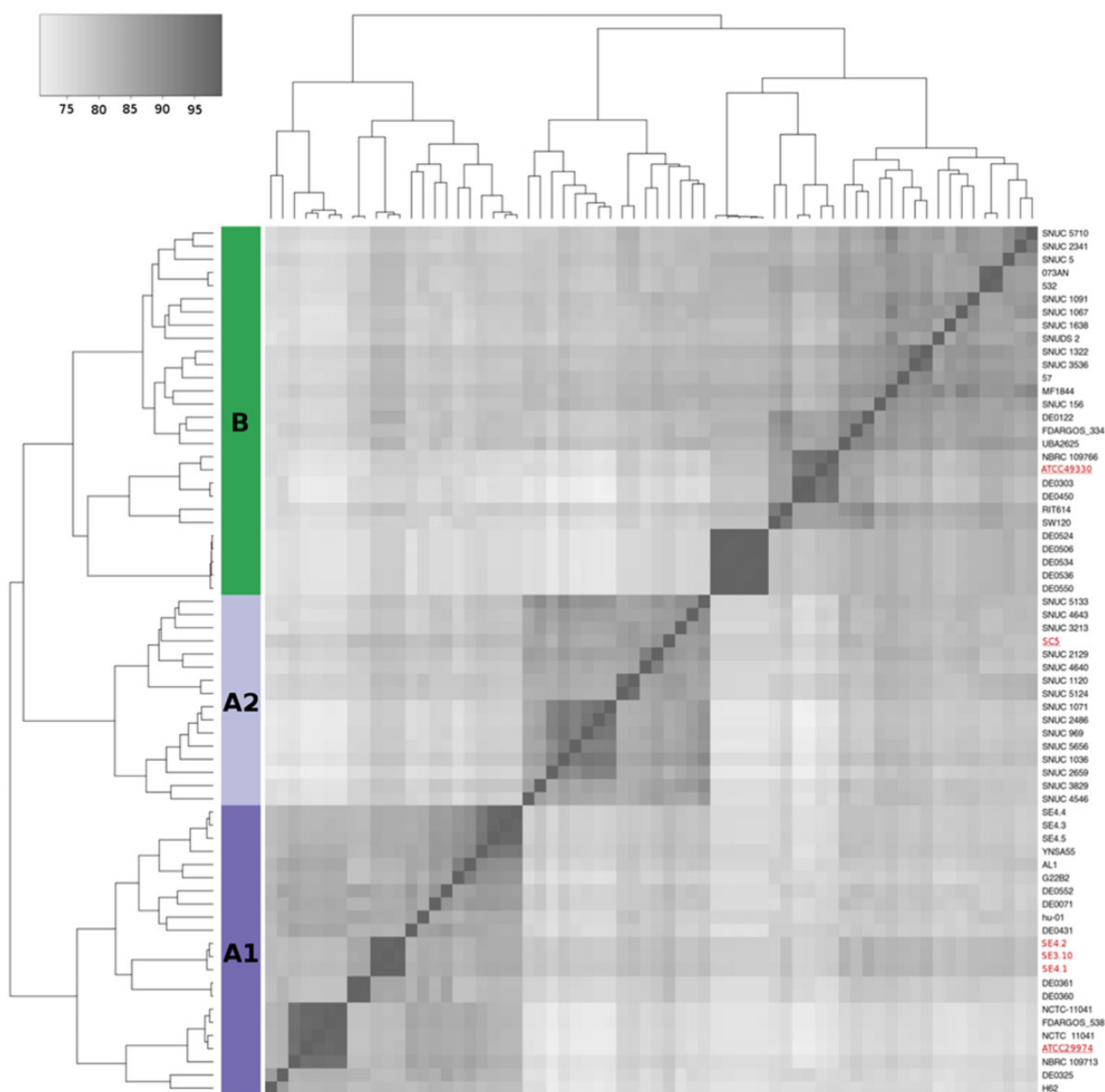


**Fig. 1.**—Heatmap of ANIb between genomes of *Staphylococcus cohnii* isolates. ANIb values are represented using a gray scale color map, with darker colors indicating higher levels of identity, according to the scale represented on the top. Strain identifiers are indicated on the rows. The panel on the left indicates cluster memberships, according to the following color codes: green = B, dark purple = A1, and light purple = A2. Columns and row dendrograms are used to group SC strains based on patterns of genome identity profiles. The novel SC5 isolate and the two type strains SCC ATCC 29974 and SCU ATCC 49330 are highlighted in red and underlined. The SE4.1, SE 4.2, and SE3.10 strains that are also discussed in the text are highlighted in red.

and [supplementary figure S2, Supplementary Material](#) online, and consistent with patterns of in silico DDH profiles, the results of these analyses suggested the presence of three distinct clusters with different mutual levels of ANIb within the SC species complex. The first two groups of isolates, referred to hereafter as A1 and A2, correspond to isolates SCC and are more closely related, whereas the third group (referred as B) is

more distantly related to A1 and A2 and is composed exclusively of SCU isolates. A1 incorporates 22 strains, including the SCC type strain ATCC 29974, A2 is composed of 16 strains and includes SC5. Finally, group B contains 28 isolates, including the type strain SCU ATCC 49330.

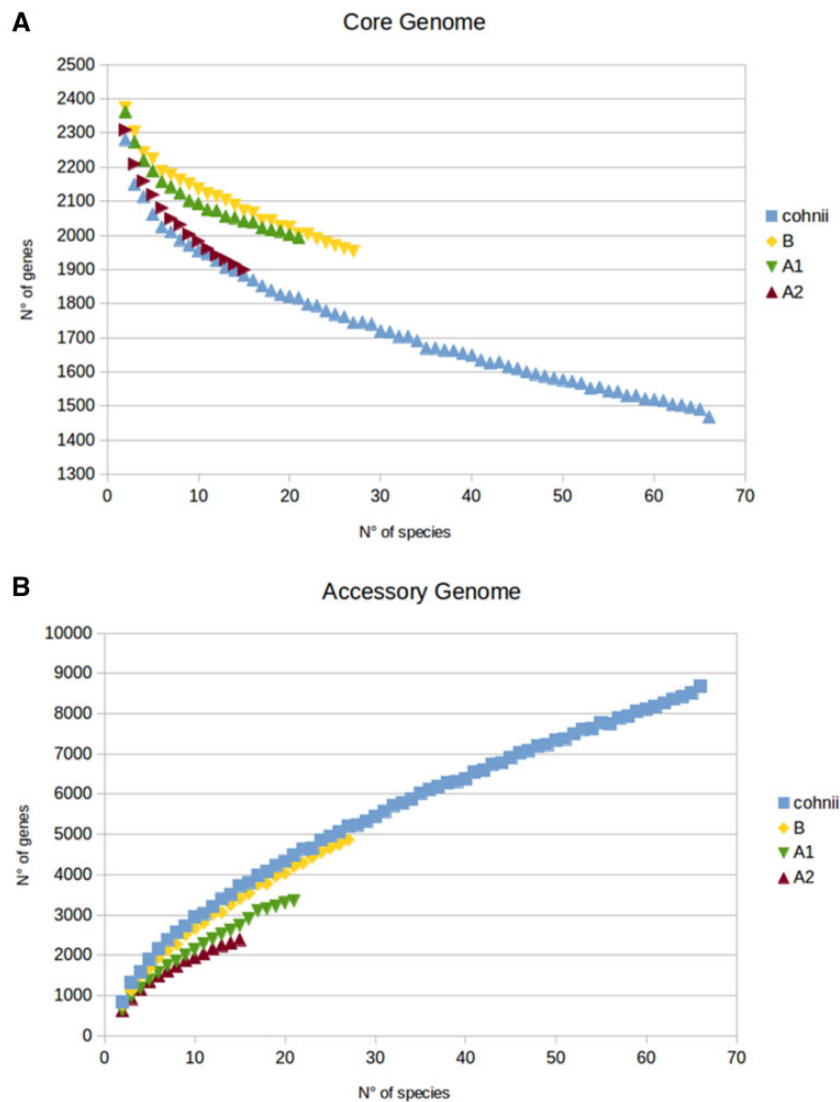
Comparisons of genomic identity levels between the draft genomes of SCC and SCU show an average ANIb of 89.43%,



**Fig. 2.**—Heatmap of gene presence/absence profiles of *Staphylococcus cohnii* isolates. Similarity of gene presence/absence profiles were estimated by computing pairwise Pearson correlation values between all the 66 genomes considered in the study. Pearson correlation coefficients are represented using a gray scale color map. Darker colors indicated higher correlation (similarity) of gene presence/absence profiles. Strain identifiers are indicated on the rows. The panel on the left is used to indicate cluster memberships, with the color codes defined in figure 1. Similar to figure 1, dendrograms are applied to the columns and rows to delineate groups of isolates with similar gene absence prevalence profiles. SC5 and the two type strains SCC ATCC 29974 and SCU ATCC 49330 are highlighted in red and underlined. The SE4.1, SE 4.2, and SE3.10 strains are highlighted in red.

a value that is well below the cut-off normally considered for inclusion in the same bacterial species (Otto 2008; Varghese et al. 2015). This approach also sustains the presence of two distinct clusters within SCC, with significantly different levels of ANIb ( $t$ -test  $P$ -value  $\leq 1e-16$ ) and an average genome sequence identity of 95.82% (supplementary fig.S2,

Supplementary Material online), a value that is normally considered borderline for the identification of bacterial species (Otto 2008; Varghese et al. 2015). Taken together, our analyses of genomic similarity profiles by means of two independent methods strongly suggest that according to the current guidelines for the delineation of bacterial species, SCC and



**FIG. 3.**—Core and accessory genome. (A) Plot of core genome size in the *Staphylococcus cohnii* species complex and in the three distinct groups of SC (A1, A2, and B) delineated in this study. X axis = number of genomes and Y axis = number of genes. (B) Plot of accessory genome size in *S. cohnii* and in the three distinct groups of SC identified by this study. X axis = number of genomes and Y axis = number of genes.

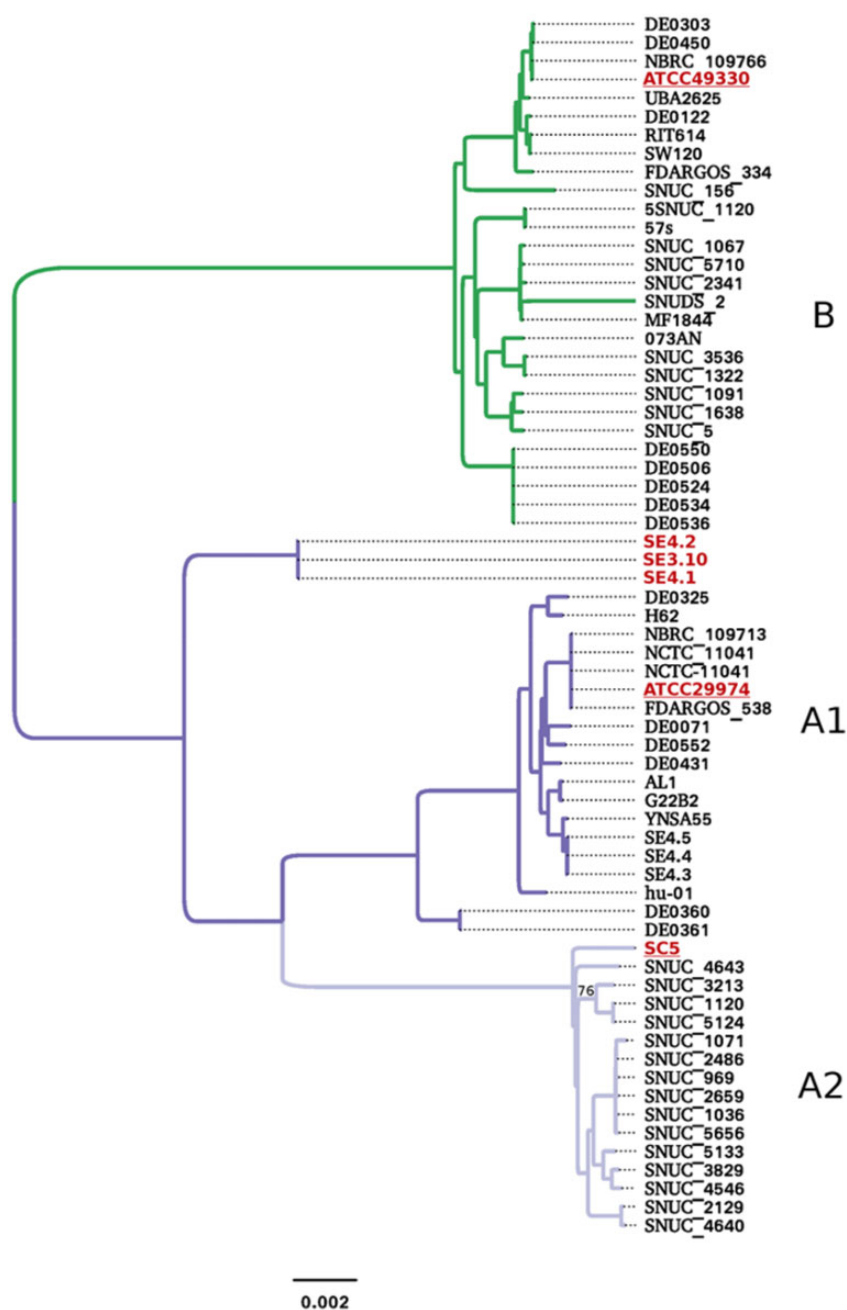
SCU should be considered as two distinct species. Consistent with this consideration, we observed that the ANI<sub>b</sub> value recovered from the comparison of the two type strains SCC ATCC 29974 and SCU ATCC 49330 is 91.5%. Notably, comparisons between the draft genome assembly of SC5 with the SCC ATCC 29974 and SCU ATCC 49330 type strains resulted in ANI<sub>b</sub> values of 95.4% and 91.0%, respectively (table 1).

#### Cluster of orthologous genes and Phylogenetic Analyses

Cluster of orthologous genes (COGs) as well as core and accessory genomes were established using an approach based on best reciprocal BLAST hits. A total of 5,456 clusters of putative orthologs with more than one gene and 5,044 singleton genes were identified. Hierarchical clustering of phenetic patterns of gene presence/absence profiles was applied

to identify SC isolates with a similar gene content. Consistent with our previous observations, three distinct groups were observed, identical in size and composition with A1, A2, and B. Notably, although the A1 and A2 clades were delineated very clearly by this analysis (fig. 2), suggesting a similar gene content within isolates of these groups, a somewhat more heterogeneous pattern was observed for group B, suggesting a higher plasticity of the pan-genome, possibly associated with lateral gene transfer.

Analysis of the SC core genome provides additional evidence for the different gene content in the A1, A2, and B. Indeed, when all the 66 available genomes are considered, a core genome of 1,468 genes was recovered, whereas notable differences were observed in the size of the core genome between the three groups (fig. 3A).



**Fig. 4.**—Phylogenetic tree of *Staphylococcus cohnii* isolates based on concatenated alignment of 1,468 core genes. Branch colors indicate the different groups identified in this study, according to the color code defined in figure 1. Bootstrap values below 95 are reported on the corresponding branches. The SE4.1, SE 4.2, and SE3.10 strains are marked in red. SCS isolate and the two type strains SCC ATCC 29974 and SCU ATCC 49330 are highlighted in red and underlined.

Although the estimated core genome size for the A2 cluster was 1,889 genes, the core genomes of A1 and B were larger, with an estimated sizes of 1,993 and 1,953 genes, respectively. As our analyses are based on nearly equivalent numbers of isolates for every group, this difference is unlikely to be the result of a biased sampling but might reflect a tendency for a more compact genome with a reduced number of genes in the A2 clade. Consistent with this hypothesis, we

observe that the average number of predicted protein coding genes is significantly reduced in A2 (average 2,566) with respect to A1 (2,618) and B (2,651), with  $P$ -values of 0.06 and 0.022, respectively, according to a Welch  $t$ -test.

Although our analyses suggest that the accessory genome of SC is relatively open (fig. 3B) and that additional genes are likely to be discovered as new genomic sequences become available, we note once again that the



**Table 2.**  
A1, A2, and B Strain-Specific Genes

Group	Accession	Annotation
A1	WP_103211109.1	ABC-F family ATP-binding cassette domain-containing protein
A1	WP_019468720.1	tRNA (N6-isopentenyl adenosine(37)-C2)-methylthiotransferase MiaB
A1	WP_019468192.1	MFS transporter
A1	WP_040030451.1	Aldehyde dehydrogenase family protein
A1	WP_019468481.1	Trigger factor
A1	WP_103211478.1	M15 family metallopeptidase
A1	WP_040030229.1	Class 1b ribonucleoside-diphosphate reductase
A1	WP_019468907.1	Hypothetical protein
A2	WP_107523479.1	Orotidine-5'-phosphate decarboxylase
A2	WP_107505199.1	BtrH N-terminal domain-containing protein
A2	WP_019468295.1	Uracil phosphoribosyltransferase
A2	WP_181187692.1	Hypothetical protein
A2	WP_107384484.1	Arsenate reductase
A2	WP_107384163.1	LytTR family transcriptional regulator
B	WP_073342256.1	NAD(P)/FAD-dependent oxidoreductase
B	WP_073344697.1	DUF4352 domain-containing protein
B	WP_046206585.1	DUF4064 domain-containing protein
B	WP_073344420.1	ABC transporter substrate-binding protein
B	WP_073345446.1	NAD(P)H-dependent oxidoreductase
B	WP_103161674.1	Anion permease
B	WP_103161408.1	Glucose 1-dehydrogenase
B	WP_073343781.1	Energy-coupling factor transporter
B	WP_073341550.1	Amino acid ABC transporter ATP-binding
B	WP_073342088.1	Thymidylate synthase

NOTE.—The protein ID and the relevant annotation is shown for each gene.

accessory genome in A2 is substantially reduced with respect to A1 and B, again consistent with systematic differences in gene content.

Phylogenetic analyses of the concatenated alignment of 1,468 core genes (fig. 4), recovered a tree with a topology consistent with the clustering of the isolates based on genome identity levels, providing an additional line of evidence for the presence of three distinct clades within the SC species complex. However, we notice that according to our phylogenetic analyses, A1 is not strictly monophyletic and a distinct basal clade formed by three isolates: SE4.1, SE 4.2, and SE3.10 is observed. Interestingly, a similar pattern is replicated also in figure 1, where the same group of isolates (SE4.1, SE 4.2, and SE3.10) form an identical basal clade, suggesting overall reduced levels of genome identity of these three strains with the other strains included in A1. Notably, the same pattern is not recovered when clustering of isolates based on gene content is considered (fig. 2). This indicates that the gene content of SE4.1, SE 4.2, and SE3.10 is highly consistent with that of other strains in the A1 cluster. Taken all together, these observations might suggest widespread lateral gene transfer between SE4.1, SE 4.2, and SE3.10 and other SC isolates or alternatively, that

while having a similar gene content with strains included in the A1 group, overall SE4.1, SE 4.2, and SE3.10 are highly divergent at sequence level. Possibly indicating faster evolutionary rates and/or ongoing diversifying selection. Importantly, we underscore that these three strains were isolated from a similar environment (rice seeds), and geographic location (India), consistent with the hypothesis that the observed reduction in genome/protein identity levels might reflect the regional/environmental diversity of bacterial communities (Lozupone and Knight, 2008).

#### Molecular Discrimination of A1, A2, and B Strains

Interestingly, analyses of core genome composition between the three SC groups proposed here, identify 8, 6, and 10 genes that are universally present in the genomes of A1, A2, and B strains, respectively, and consistently absent from genomes of the other groups. These genes (Table 2) are not adjacent in the genomes suggesting that they do not represent operons. PCR assays based on the presence/absence profiles of these genes might be used to discriminate between members of the A1, A2, and B groups.

**Table 3**  
*Staphylococcus cohnii* Genome Metadata<sup>a,b</sup>

Strain	Isolation Source	Assembly	BioSample	BioProject	Geographic Location	Current Rank	Proposed Rank
SC5	Biological Safety Cabinet	—	SAMN14142771	PRJNA607668	Bari, Italy	SCC	SCB
SNUC 2659	<i>Bos taurus</i> bovine mastitis	GCA_003035875.1	SAMN06172961	PRJNA342349	Quebec, Canada	SC	SCB
SNUC 5656	<i>Bos taurus</i> bovine mastitis	GCA_003035785.1	SAMN06172970	PRJNA342349	Quebec, Canada	SC	SCB
SNUC 5133	<i>Bos taurus</i> bovine mastitis	GCA_003577875.1	SAMN06172969	PRJNA342349	Quebec, Canada	SC	SCB
SNUC 1036	<i>Bos taurus</i> bovine mastitis	GCA_003035975.1	SAMN06172951	PRJNA342349	Ontario, Canada	SC	SCB
SNUC 1120	<i>Bos taurus</i> bovine mastitis	GCA_003035965.1	SAMN06172955	PRJNA342349	Quebec, Canada	SC	SCB
SNUC 969	<i>Bos taurus</i> bovine mastitis	GCA_003039995.1	SAMN06172950	PRJNA342349	Quebec, Canada	SC	SCB
SNUC 4643	<i>Bos taurus</i> bovine mastitis	GCA_003577905.1	SAMN06172967	PRJNA342349	Ontario, Canada	SC	SCB
SNUC 1071	<i>Bos taurus</i> bovine mastitis	GCA_003578125.1	SAMN06172953	PRJNA342349	Ontario, Canada	SC	SCB
SNUC 2129	<i>Bos taurus</i> bovine mastitis	GCA_003039915.1	SAMN06172958	PRJNA342349	Quebec, Canada	SC	SCB
SNUC 5124	<i>Bos taurus</i> bovine mastitis	GCA_003035485.1	SAMN06172968	PRJNA342349	Quebec, Canada	SC	SCB
SNUC 3213	<i>Bos taurus</i> bovine mastitis	GCA_003577955.1	SAMN06172962	PRJNA342349	Ontario, Canada	SC	SCB
SNUC 2486	<i>Bos taurus</i> bovine mastitis	GCA_003035905.1	SAMN06172960	PRJNA342349	Atlantic, Canada	SC	SCB
SNUC 3829	<i>Bos taurus</i> bovine mastitis	GCA_003035865.1	SAMN06172964	PRJNA342349	Ontario, Canada	SC	SCB
SNUC 4546	<i>Bos taurus</i> bovine mastitis	GCA_003577915.1	SAMN06172965	PRJNA342349	Ontario, Canada	SC	SCB
SNUC 4640	<i>Bos taurus</i> bovine mastitis	GCA_003035505.1	SAMN06172966	PRJNA342349	Ontario, Canada	SC	SCB
FDARGOS_538	Human clinical isolate	GCA_003956025.1	SAMN10163250	PRJNA231221	Missing	SC	SCC
SE4.4	Rice seed	GCA_001876785.1	SAMN03097241	PRJNA263233	India	SC	SCC
SE4.3	Rice seed	GCA_001876755.1	SAMN03097240	PRJNA263231	India	SC	SCC
ATCC 29974	Human isolate	GCA_900240165.1	SAMEA104410613	PRJEB22856	Liverpool, UK	SCC	SCC
SE4.5	Rice seed	GCA_001876805.1	SAMN03097242	PRJNA263234	India	SC	SCC
NCTC 11041	Human skin	GCA_002902365.1	SAMN06177162	PRJNA339206	London, UK	SCC	SCC
NCTC 11041	Human skin	GCA_900458255.1	SAMEA3871778	PRJEB6403	USA	SC	SCC
G22B2	Human gall bladder	GCA_000981215.1	SAMN03352186	PRJNA275680	Chandigarh, India	SCC	SCC
hu-01	Human skin swab sample	GCA_000513495.2	SAMN02388844	PRJNA225658	Hangzhou, China	SCC	SCC
AL1	Soy sauce (plant food)	GCA_000292305.1	SAMN02471867	PRJNA171726	Malaysia	SC	SCC
DE0361	Environmental	GCA_007673385.1	SAMN11792521	PRJNA543692	Durham, North Carolina, USA	SC	SCC
DE0431	Environmental	GCA_007668065.1	SAMN11792591	PRJNA543692	Durham, North Carolina, USA	SC	SCC
DE0071	Environmental	GCA_007679825.1	SAMN11792231	PRJNA543692	Durham, North Carolina, USA	SC	SCC
DE0552	Environmental	GCA_007666185.1	SAMN11792712	PRJNA543692	Durham, North Carolina, USA	SC	SCC
DE0325	Environmental	GCA_008764065.1	SAMN11792485	PRJNA543692	Durham, North Carolina, USA	SC	SCC
NBRC 109713	Human skin	GCA_007992675.1	SAMD00172682	PRJDB1638	Missing	SCC	SCC
DE0360	Environmental	GCA_007673395.1	SAMN11792520	PRJNA543692	Durham, North Carolina, USA	SC	SCC
YN5A55	Human	GCA_005861955.1	SAMN11775280	PRJNA543691	Yunnan, China	SC	SCC
H62	Environmental (air)	GCA_001650645.1	SAMN04591361	PRJNA316869	California, USA	SC	SCC
SE4.1	Rice seed	GCA_001876705.1	SAMN03097238	PRJNA263229	India	SC	SCC
SE4.2	Rice seed	GCA_001876735.1	SAMN03097239	PRJNA263230	India	SC	SCC
SE3.10	Rice seed	GCA_001876725.1	SAMN03097237	PRJNA263228	India	SC	SCC
FDARGOS_334	Human peripheral blood	GCA_002984565.1	SAMN06173347	PRJNA231221	Maryland, USA	SC	SU
532	Human catheter	GCA_000972575.1	SAMN03449104	PRJNA279286	Nuevo Leon, Monterrey, Mexico	SCC	SU



**Note.—The protein ID and the relevant annotation is shown for each gene.**

## Discussion

Difficulties in the taxonomic assignment of a novel *Staphylococcus* isolate based on well-established genome identity metrics (ANIb and DDH), prompted us to perform an extensive phylogenomic analysis of the *Staphylococcus cohnii* (SC) species complex, based on all (65) currently available draft genome assemblies. Analyses based on genomic identity levels, core and accessory genome size, gene content, and phylogenetic approaches consistently and strongly suggest that the SC species complex, as currently defined is composed of at least three distinct groups, and advocate a revision of the current phylogenetic classification of SC.

Based on analyses presented in this study, and on the current guidelines and best practices for the classification of bacterial species (Varghese et al. 2015), we propose that the two SC subspecies, as described by Schleifer and Kloos (1975), should be instead regarded as two distinct species: SC, corresponding to groups A1 and A2 in this work (fig. 2) and *Staphylococcus urealyticus* (SU), corresponding to group B. Additionally, the revised classification of SC should include two subspecies: *Staphylococcus cohnii* subsp. *cohnii* (SCC), corresponding to group A1 and *Staphylococcus cohnii* subsp. *barensis* (SCB), corresponding to group A2 and including our novel isolate. Importantly, although the three different clades were not associated with any evident morphological or phenotypic difference, possible criteria for an effective discrimination of SU, SCC, and SCB are proposed in the current study. These include the application of standard approaches based either on DDH and ANIb, two genome-wide similarity metrics, which are considered a reference standard for delineation of bacterial species. Moreover, according to our analyses, a limited, but consistent number of lineage-specific genes is observed in each of SU, SCC, and SCB. In principle, simple tests based on targeted PCR resequencing of these genes could be used to develop a highly effective molecular assays for the discrimination of both species and subspecies proposed here.

Although increased future environmental sampling may help overcome any ascertainment bias inherent in the available cohort, and to resolve the uncertainty of the phylogenetic placement of the three isolates from rice seeds, our study demonstrates a wider host range for SCC than previously hypothesized (Kloos and Wolfshohl 1983, 1991) and includes isolates from plants (rice seed) and vegetable liquid food (soy sauce) (table 3). Conversely, isolates of SCB strains displayed a narrower hosts range, as 15 out of 16 strains were isolated from *Bos taurus* intramammary infections (bovine mastitis) and one (the novel strain described in this study) isolated from a disused biological cabinet. Intriguingly, we observe that isolates of SCB also show a more compact

genome and a significant reduction in gene content, compared to SCC a consideration that might at least in part explain its more reduced/specialized host range. The apparent host range of SU, based on the source of isolation as reported in the metadata associated with GenBank genome submissions, is consistent with previous reports with members of this candidate species isolated both from humans (e.g., skin, blood, catheter) and other animals, including ducks, dogs, cows, and goats (table 3).

Consistent with the problems encountered in the taxonomic classification of our novel SCB isolate, we observed that taxonomic classifications of several SC isolates deposited in GenBank are incomplete or discordant with our analyses. For example, strain 532 (Mendoza-Olazarán et al. 2017) is assigned to SCC in the NCBI biosample but is likely a member of SCU. Furthermore, strain 073AN (NCBI draft genome A.N.: FMPF00000000.1), isolated in Tanzania (Africa) from a goat, is labeled as SCC in the associated Biosample and Bioproject submissions (SAMEA3109313 and ERS576551) (INSDC, last update 2016-09-28, <https://www.ncbi.nlm.nih.gov/biosample/SAMEA3109313/>) although both the original publication (Seni et al. 2019) and our own analyses indicate that it is SCU. These misclassifications are likely consequences of problems in the initial classification of the isolates, arising from ambiguities in the delineation of the SC species complex. Consistent with this hypothesis, we observe that in many cases, the classification of SC strains is limited to species level.

Misclassified sequences, errors in initial taxonomic classification, and annotation can occur and be propagated resulting in the unintentional misclassification of bacterial strains (Federhen 2015). Accordingly, we believe that the approach and criteria presented here may be of general interest and could be applied on a larger scale for the resolution of complex/conflicting taxonomic assignments.

## Data Availability

The SC5 draft genome is available in NCBI under the accession number JAALCY000000000, BioSample accession number SAMN14142771, and BioProject number ID PRJNA607668.

## Supplementary Material

Supplementary tables S1–S2 and figures S1–S2 are available at *Genome Biology and Evolution* online.

## Acknowledgments

This work was supported by INMARE (H2020-BG-2014-2, GA 634486), EMBRIC (H2020-INFRADEV-1-2014-1, GA 654008), EXCELERATE (H2020-INFRADEV-1-2015-1, GA676559), the Molecular Biodiversity Laboratory (Life Watch, Italy), and the Italian Node of ELIXIR. The funder had no role in the study

design, data collection, and interpretation or the decision to submit the work for publication.

## Literature Cited

- Adeyemi AI, Sulaiman AA, Solomon BB, Chinedu OA, Victor IA. 2010. Bacterial bloodstream infections in HIV-infected adults attending a Lagos Teaching Hospital. *J Health Popul Nutr.* 28(4):318–326.
- Afgan E, et al. 2018. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2018 update. *Nucleic Acids Res.* 46(W1):W537–W544.
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol.* 215(3):403–410.
- Argemi X, Hansmann Y, Prola K, Prévost G. 2019. Coagulase-negative Staphylococci pathogenomics. *Int J Mol Sci.* 20(5):1215.
- Auch AF, Klenk HP, Göker M. 2010. Standard operating procedure for calculating genome-to-genome distances based on high-scoring segment pairs. *Stand Genomic Sci.* 2(1):142–148.
- Bankevich A, et al. 2012. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol.* 19(5):455–477.
- Basaglia G, Moras L, Bearz A, Scalone S, Paoli PD. 2003. *Staphylococcus cohnii* septicaemia in a patient with colon cancer. *J Med Microbiol.* 52(1):101–102.
- Bean DC, Wigmore SM, Wareham DW. 2017. Draft genome sequence of *Staphylococcus cohnii* subsp. *urealyticus* isolated from a healthy dog. *Genome Announc.* 5(7):e01628-16.
- Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*30(15):2114–2120.
- Camacho C, et al. 2009. BLAST+: architecture and applications. *BMC Bioinformatics*10(1):421.
- Castresana J. 2000. Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Mol Biol Evol.* 17(4):540–552.
- Chiara M, et al. 2018. A-GAME: improving the assembly of pooled functional metagenomics sequence data. *BMC Genomics*19(1):44.
- Colston S M, et al. 2014. Bioinformatic genome comparisons for taxonomic and phylogenetic assignments using *Aeromonas* as a test case. *mBio.* 5(6):e02136 10.1128/mBio.02136-14PMC: 25406383
- Crossley KB, Jefferson KK, Archer GL, Fowler VG, editors. 2009. *Staphylococci in human disease*. 2nd ed. Chichester (West Sussex): Wiley-Blackwell.
- Darriba D, Taboada GL, Doallo R, Posada D. 2011. ProtTest 3: fast selection of best-fit models of protein evolution. *Bioinformatics*27(8):1164–1165.
- David MD, Elliott T. 2015. Coagulase-negative Staphylococci. *Br J Hosp Med (Lond).* 76(8):C126–C128.
- Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* 32(5):1792–1797.
- Federhen S. 2015. Type material in the NCBI Taxonomy Database. *Nucleic Acids Res.* 43(D1):D1086–1098.
- Garrido-Sanz D, et al. 2016. Correction: Genomic and Genetic Diversity within the *Pseudomonas fluorescens* Complex. *PLoS ONE.* 11(4):e0153733 10.1371/journal.pone.0153733
- Guindon S, Delsuc F, Dufayard J-F, Gascuel O. 2009. Estimating maximum likelihood phylogenies with PhyML. In: Posada D, editor. *Bioinformatics for DNA sequence analysis*. Totowa (NJ): Methods Molecular Biology Humana Press. pp. 113–137.
- Kloos WE, Wolfshohl JF. 1983. Deoxyribonucleotide sequence divergence between *Staphylococcus cohnii* subspecies populations living on primate skin. *Curr Microbiol.* 8(2):115–121.
- Kloos WE, Wolfshohl JF. 1991. *Staphylococcus cohnii* subspecies: *Staphylococcus cohnii* subsp. *cohnii* subsp. *nov.* and *Staphylococcus cohnii* subsp. *urealyticus* subsp. *nov.* *Int J Syst Bacteriol.* 41(2):284–289.
- Larsen J, Andersen PS, Winstel V, Peschel A. 2017. *Staphylococcus aureus* CC395 harbours a novel composite staphylococcal cassette chromosome mec element. *J Antimicrob Chemother.* 72(4):1002–1005.
- Lavecchia A, et al. 2018. Draft genome sequences of three novel *Staphylococcus arlettae* strains isolated from a disused biological safety cabinet. *Microbiol Resour Announc.* 7(13). doi:10.1128/MRA.01012-18.
- Lozupone CA, Knight R. 2008. Species divergence and the measurement of microbial diversity. *FEMS Microbiol Rev.* 32(4):557–578.
- Mastroianni A, Coronado O, Nanetti A, Manfredi R, Chiodo F. 1996. *Staphylococcus cohnii*: an unusual cause of primary septic arthritis in a patient with AIDS. *Clin Infect Dis.* 23(6):1312–1313.
- Meier-Kolthoff JP, Auch AF, Klenk H-P, Göker M. 2013. Genome sequence-based species delimitation with confidence intervals and improved distance functions. *BMC Bioinformatics*14(1):60.
- Mendoza-Olazarán S, et al. 2017. Draft genome sequences of two opportunistic pathogenic strains of *Staphylococcus cohnii* isolated from human patients. *Stand Genomic Sci.* 12:49.
- Moavwad AA, et al. 2019. Evolution of antibiotic resistance of coagulase-negative Staphylococci isolated from healthy Turkeys in Egypt: first report of linezolid resistance. *Microorganisms*7(10):476.
- Okudera H, Kobayashi S, Hongo K, Mizuno M. 1991. Fatal meningitis due to *Staphylococcus cohnii*. Case report. *Neurosurg Rev.* 14(3):235–236.
- Otto M. 2008. Staphylococcal biofilms. *Curr Top Microbiol Immunol.* 322:207–228.
- Otto M. 2013. Coagulase-negative Staphylococci as reservoirs of genes facilitating MRSA infection. *Bioessays*35(1):4–11.
- Richter M, Rosselló-Móra R. 2009. Shifting the genomic gold standard for the prokaryotic species definition. *Proc Natl Acad Sci USA.* 106(45):19126–19131.
- Schleifer KH, Kloos WE. 1975. Isolation and characterization of Staphylococci from human skin I. Amended descriptions of *Staphylococcus epidermidis* and *Staphylococcus saprophyticus* and descriptions of three new species: *Staphylococcus cohnii*, *Staphylococcus haemolyticus*, and *Staphylococcus xylosus*. *Int Syst Evol Micr.* 25(1):50–61.
- Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*30(14):2068–2069.
- Seni J, et al. 2019. Draft genome sequence of a multidrug-resistant caprine isolate of *Staphylococcus cohnii* subsp. *urealyticus* from Tanzania encoding *ermB*, *tet(K)*, *dfcG*, *fusF* and *fosD*. *J Glob Antimicrob Resist.* 18:163–165.
- Sneath PH. 1992. Correction of orthography of epithets in *Pasteurella* and some problems with recommendations on latinization. *Int J Syst Bacteriol.* 42(4):658–659.
- Song Y, et al. 2017. cfr-mediated linezolid-resistant clinical isolates of methicillin-resistant coagulase-negative Staphylococci from China. *J Glob Antimicrob Resist.* 8:1–5.
- Szczuka E, Jabłońska L, Kaznowski A. 2016. Coagulase-negative Staphylococci: pathogenesis, occurrence of antibiotic resistance genes and in vitro effects of antimicrobial agents on biofilm-growing bacteria. *J Med Microbiol.* 65(12):1405–1413.
- Szewczyk EM, Piotrowski A, Różalska M. 2000. Predominant Staphylococci in the intensive care unit of a paediatric hospital. *J Hosp Infect.* 45(2):145–154.
- Szewczyk EM, Różalska M, Cieślowski T, Nowak T. 2004. Plasmids of *Staphylococcus cohnii* isolated from the intensive-care unit. *Folia Microbiol.* 49(2):123–131.
- Varghese NJ, et al. 2015. Microbial species delineation using whole genome sequences. *Nucleic Acids Res.* 43(14):6761–6771.
- Waldon E, Sobiś-Glinkowska M, Szewczyk EM. 2002. Evaluation of selected features of *Staphylococcus cohnii* enabling colonization of humans. *Folia Microbiol.* 47(5):565–571.

- Whelan S, Goldman N. 2001. A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach. *Mol Biol Evol.* 18(5):691–699.
- Winstel V, et al. 2013. Wall teichoic acid structure governs horizontal gene transfer between major bacterial pathogens. *Nat Commun.* 4(1):2345.
- Yamashita S, Yonemura K, Sugimoto R, Tokunaga M, Uchino M. 2005. *Staphylococcus cohnii* as a cause of multiple brain abscesses in Weber-Christian disease. *J Neurol Sci.* 238(1–2):97–100.
- Yong YY, Dykes GA, Choo WS. 2019. Biofilm formation by staphylococci in health-related environments and recent reports on their control using natural compounds. *Crit Rev Microbiol.* 45(2):201–222.
- Zhang H, Ma Y, Liu P, Li X. 2016. Multidrug resistance operon *emrAB* contributes for chromate and ampicillin co-resistance in a *Staphylococcus* strain isolated from refinery polluted river bank. *Springerplus*5(1):1648.
- Zhang J, Kobert K, Flouri T, Stamatakis A. 2014. PEAR: a fast and accurate Illumina Paired-End reAd mergeR. *Bioinformatics*30(5):614–620.

**Associate editor:** Maria Costantini