

# 1 **Genomic characterization of the *C. tuberculostearicum*** 2 **species complex, a ubiquitous member of the human skin** 3 **microbiome**

4 Nashwa M. Ahmed <sup>a</sup>, Payal Joglekar <sup>a</sup>, Clayton Deming <sup>a</sup>, NISC Comparative Sequencing  
5 Program <sup>b</sup>, Katherine P. Lemon <sup>c,d</sup>, Heidi H. Kong <sup>e</sup>, Julia A. Segre <sup>a</sup>, Sean Conlan <sup>a,†</sup>

6  
7 <sup>a</sup> Microbial Genomics Section, Translational and Functional Genomics Branch, NHGRI, NIH,  
8 Bethesda, Maryland, USA

9 <sup>b</sup> NIH Intramural Sequencing Center, NHGRI, NIH, Rockville, Maryland, USA

10 <sup>c</sup> Alkek Center for Metagenomics & Microbiome Research, Department of Molecular Virology &  
11 Microbiology, Baylor College of Medicine, Houston, Texas, USA

12 <sup>d</sup> Division of Infectious Diseases, Texas Children's Hospital, Department of Pediatrics, Baylor  
13 College of Medicine, Houston, Texas, USA

14 <sup>e</sup> Cutaneous Microbiome and Inflammation Section, NIAMS, NIH, Bethesda, Maryland, USA

15 <sup>†</sup> Corresponding author [conlans@mail.nih.gov](mailto:conlans@mail.nih.gov)

16  
17 **ABSTRACT** *Corynebacterium* is a predominant genus in the skin microbiome, yet its genetic  
18 diversity on skin is incompletely characterized and lacks a comprehensive set of reference  
19 genomes. Our work aims to investigate the distribution of *Corynebacterium* species on the skin,  
20 as well as to expand the existing genome reference catalog to enable more complete  
21 characterization of skin metagenomes. We used V1-V3 16S rRNA gene sequencing data from  
22 14 body sites of 23 healthy volunteers to characterize *Corynebacterium* diversity and distribution  
23 across healthy human skin. *Corynebacterium tuberculostearicum* is the predominant species  
24 found on human skin and we identified two distinct *C. tuberculostearicum* ribotypes (A & B) that  
25 can be distinguished by variation in the 16S rRNA V1-V3 sequence. One is distributed across all  
26 body sites and the other found primarily on the feet. We performed whole genome sequencing  
27 of 40 *C. tuberculostearicum* isolates cultured from the skin of five healthy individuals across  
28 seven skin sites. We generated five closed genomes of diverse *C. tuberculostearicum* which  
29 revealed that *C. tuberculostearicum* isolates are largely syntenic and carry a diversity of  
30 methylation patterns, plasmids and CRISPR/Cas systems. The pangenome of *C.*  
31 *tuberculostearicum* is open with a core genome size of 1806 genes and a pangenome size of  
32 5451 total genes. This expanded pangenome enabled the mapping of 24% more *C.*  
33 *tuberculostearicum* reads from shotgun metagenomic datasets derived from skin body sites.

34 Finally, while the genomes from this study all fall within a *C. tuberculostearicum* species  
35 complex, the ribotype B isolates may constitute a new species.

36

37 **IMPORTANCE** Amplicon sequencing data combined with isolate whole genome sequencing  
38 has expanded our understanding of *Corynebacterium* on the skin. Human skin is characterized  
39 by a diverse collection of *Corynebacterium* species but *C. tuberculostearicum* predominates  
40 many sites. Our work supports the emerging idea that *C. tuberculostearicum* is a species  
41 complex encompassing several distinct species. We produced a collection of genomes that help  
42 define this complex including a potentially new species which we are calling *C. hallux* based on  
43 a preference for sites on the feet, whole-genome average nucleotide identity, pangenomics and  
44 growth in skin-like media. This isolate collection and high-quality genome resource sets the  
45 stage for developing engineered strains for both basic and translational clinical studies.

46

47           Microbiomes are shaped by taxa that are both characteristic to those sites and  
48 functionally important to that community. The genus *Corynebacterium* is one such taxa for the  
49 human skin and nares. Foundational studies using 16S rRNA gene sequencing and shotgun  
50 metagenomics by our lab (1, 2) and others (3) have established *Corynebacterium* as common  
51 members of the skin microbiome. While *Corynebacterium* have been positively correlated with  
52 the resolution of dysbiosis associated with eczema flares (4), the importance of the  
53 *Corynebacterium* spp. is less defined for skin disease severity in primary immune deficient  
54 patients (5, 6). *Corynebacterium* spp. are predominant members of the human aerodigestive  
55 tract microbiome (nares, oral cavity and respiratory tract) (3) and participate in microbe-microbe  
56 interactions with members of nasal microbiome (7, 8). *Corynebacterium* have been shown to  
57 engage with the host immune system, specifically *C. accolens*-promoted IL23-dependent  
58 inflammation in mice on a high-fat diet (9). *C. bovis* and *C. mastiditis* have been shown to  
59 predominate the microbiome of a ADAM10-deficient mouse model (10) as well as an ADAM17-  
60 deficient mouse model of eczema (11). Finally, *C. tuberculostearicum* has been shown to  
61 induce inflammation in human epidermal keratinocyte cell cultures (12). These studies establish  
62 *Corynebacterium* spp. as key members of the skin microbiome capable of both microbe-microbe  
63 and microbe-host interactions.

64           A critical resource for understanding the biology of *Corynebacterium* on the skin is a  
65 robust collection of complete reference genomes, including isolates collected from a variety of  
66 individuals and body sites. Previously published genome collections from skin- or nares-resident  
67 species include *Staphylococcus epidermidis* (13), *Cutibacterium acnes* (14) and the recent  
68 comparative analysis of *Dolosigranulum pigrum* (15). Of note, while emerging bioinformatic  
69 methods and pipelines are now being employed to extract nearly-complete genomes (MAGs)  
70 from metagenomic assemblies of skin samples (16), MAGs are not yet a substitute for genomes  
71 from cultured isolates to understand strain level or pangenomic diversity. In addition to  
72 functional prediction, comparative genomics is increasingly being used to augment conventional  
73 microbiological methods to define or redefine taxonomic boundaries (17, 18), as well as  
74 describe the full extent of diversity within these boundaries (19). A pangenome, which  
75 encompasses the complete set of genes present within a set of genome sequences, enables  
76 the characterization of gene-level heterogeneity within a taxonomic group. The pangenome is  
77 commonly subdivided into the 'core' genome, referring to genes present in all strains, and the  
78 'accessory' or 'dispensable' genome, referring to those present in only one or some isolates.  
79 (The accessory pangenome can be further subdivided to reflect a wider range of gene  
80 uniqueness, e.g. singletons.) Thorough characterization of taxa is limited by the availability of

81 representative and high-quality genome assemblies. Unfortunately, with the exceptions of  
82 clinically relevant *Corynebacterium* spp. (e.g., *C. diphtheriae*, *C. striatum* and *C.*  
83 *pseudotuberculosis*), the genus is inadequately sequenced, with 75% of species having fewer  
84 than six genomes. This includes common skin-associated species like *C. tuberculostearicum*  
85 with just five unique isolate genomes, only two of which are from skin.

86 This work seeks first to characterize the distribution of *Corynebacterium* across 14 skin  
87 sites from 23 healthy volunteers. The second goal of this work focuses on what we identify as  
88 the predominant skin *Corynebacterium* species, *C. tuberculostearicum*. We have sequenced 23  
89 distinct *C. tuberculostearicum* strains (n=40 genomes before dereplication), a five-fold increase  
90 in the number of publicly available, unique genomes (n=5). In addition to short-read assemblies,  
91 we generated five complete genomes which, along with the type strain (DSM44922),  
92 demonstrate that *C. tuberculostearicum* genomes are largely syntenic and carry a number of  
93 methylation systems as well as a CRISPR/Cas system. Genes from the *C. tuberculostearicum*  
94 genomes in our collection fall into 5451 gene clusters comprising the species pangenome. This  
95 expanded pangenome, as compared to existing public references, improved the mapping of *C.*  
96 *tuberculostearicum* metagenomic reads from unrelated healthy volunteers. In addition, we have  
97 identified a distinct *C. tuberculostearicum* clade that is highly enriched on the feet that may  
98 represent a new species, tentatively designated *Corynebacterium hallux*.

99

## 100 **Results**

101 ***Corynebacterium* spp. are predominant members of the healthy skin microbiome** To  
102 explore the tropism of *Corynebacterium*, we surveyed the microbial diversity of healthy human  
103 skin using existing 16S rRNA V1-V3 amplicon sequencing data (5, 20). Clinical samples were  
104 obtained from 23 healthy volunteers across 14 body sites: sebaceous (back, Ba; occiput, Oc;  
105 external auditory canal, Ea; retroauricular crease, Ra; manubrium, Mb; glabella, Gb), moist  
106 (inguinal crease, Ic; antecubital crease, Ac), dry (hypothenar palm, Hp; volar forearm, Vf), foot  
107 (toe nail, Tn; toe web, Tw; plantar heel, Ph) and (N)ares. An average of 10,000 sequences per  
108 sample were generated which yielded a total of 8334 amplicon sequence variants (ASV), or  
109 unique 16S rRNA gene signatures. After rarefying the dataset to an even depth, 5967 ASVs  
110 remained. As expected, the dominant genera identified on the skin, present in 94% of skin  
111 samples, were *Cutibacterium* (41% of reads, ASV1 is *C. acnes*), *Staphylococcus* (9% of reads,  
112 ASV2 is *S. epidermidis*), and *Corynebacterium* (9% of reads, ASV3 is *C. tuberculostearicum*).

113

114 The genus *Corynebacterium* was present in 96% of the skin sites sequenced, averaging 17% of  
115 reads. With a preference for moist over sebaceous skin sites (Fig. S1), *Corynebacterium* thrives  
116 in the humid, temperate environments of the feet and nares. While variation in species  
117 composition was observed between individuals, some sites and habitats displayed species  
118 enrichment at specific locations across multiple individuals (Fig. 1A). We observed that *C.*  
119 *accolens* was enriched in the nares, with a prevalence of 83-87% across nares samples and  
120 constituted an average of 33-41% of *Corynebacterium* reads. *C. afermentans* were enriched  
121 across feet sites, where they were present in 54% of samples and comprised an average of  
122 17% of *Corynebacterium* reads. Most notably, however, we found that *C. tuberculostearicum*  
123 was present in 94% of body sites and was often the most abundant *Corynebacterium*. *C.*  
124 *tuberculostearicum* reads represented 67% of corynebacterial reads in the feet, 47% in dry  
125 sites, 58% in sebaceous sites, and 46% in the nares.

126

127 ***C. tuberculostearicum* is the most common skin *Corynebacterium*** A variety of marker  
128 gene approaches have been employed to determine the phylogenetic relationships between  
129 *Corynebacterium* species including combinations of 16S rRNA, *rpoB*, *rpoC* and *gyrA* genes (for  
130 review see (21)). In general, it is difficult to accurately classify *Corynebacterium* to the species-  
131 level using amplicon data and standard reference databases. The Human Oral Microbiome  
132 Database (3) is a curated database that includes a training set with a supraspecies taxonomic  
133 level enabling assignment of sequences to multiple species where ambiguity exists. In our case,  
134 >99.5% of sequences classified as *C. tuberculostearicum* using the Refseq classification, were  
135 also classified as *C. tuberculostearicum* (part of the *accolens/macginleyi/tuberculostearicum*  
136 superspecies) by eHOMD, including the two predominant 16S rRNA sequence variants, ASV3  
137 and ASV13 which differed by a SNP and a single-base indel.

138 ASV3 constituted 83% of *C. tuberculostearicum*-classified reads (compared to < 8% for  
139 all other ASVs of this species) and showed a cosmopolitan distribution across body sites (Fig.  
140 1B). Found in 100% of healthy volunteers (N=23) and 87% (254/293) of skin samples, ASV3  
141 was predominant and ubiquitous across human skin. Relative abundance analysis revealed  
142 ASV3 abundance > 85% within all habitats except foot sites, where it made up 66% of *C.*  
143 *tuberculostearicum*-classified reads. As of this writing, the existing *C. tuberculostearicum* NCBI  
144 reference genomes containing complete V1-V3 sequences are all ASV3 as are 100% of 16S  
145 rRNA gene *C. tuberculostearicum* references in the SILVA reference database.

146 In contrast to cosmopolitan ASV3, ASV13 was enriched primarily on feet, constituting  
147 28% of *C. tuberculostearicum*-classified reads from the Ph, Tw, and Tn sites (8%, 9%, and 70%,

148 respectively). In 9 of 23 HVs, ASV13 constituted over 90% of *C. tuberculostearicum*-mapped  
149 reads within a single foot site (Fig S2); notably, much of this predominance was observed in Tn  
150 sites. In addition, we observed that some individuals exhibited within-site predominance by  
151 other less common ASVs, with some individuals colonized by a single non-dominant ASV  
152 across multiple body sites. In HV 12, for example, 52-100% of *C. tuberculostearicum*-mapped  
153 reads in each body site excluding the toenail are classified as ASV39. We also noted that, while  
154 sites on the feet (Ph, Tn, Tw) were often colonized by multiple ASVs, other body sites tended to  
155 be colonized by a single ASV.

156 We searched the SILVA database (v138.1) for perfect matches to the ASV13 sequence  
157 and found 152 matches, all associated with uncultured *Corynebacterium*. The majority of them  
158 were from our own full-length 16S rRNA gene sequencing of skin microbiome samples (1). This  
159 observation, combined with the fact that all the existing *C. tuberculostearicum* reference  
160 genomes had the more common ASV3 sequence variant, led us to hypothesize that the ASV13  
161 sequence, which we hereafter refer to as ribotype B, could be associated with an unrecognized  
162 species or subspecies. For the purposes of the current work, we will use the term *C.*  
163 *tuberculostearicum* species complex (22, 23) to refer to all *C. tuberculostearicum*-like isolates  
164 found on skin. Additionally, we will refer to the predominant ASV3 OTU as ribotype A.

165  
166 **Expanding the *C. tuberculostearicum* complex reference catalog** Prior to this study, only  
167 five *C. tuberculostearicum* species complex isolates had been sequenced and submitted to  
168 NCBI. Only two of those were from human skin and neither was a closed genome. To expand  
169 the *C. tuberculostearicum* complex reference genome catalog, we sequenced isolates from five  
170 different HVs (Supplementary Table S1). To enrich for the previously unsequenced ribotype B  
171 *Corynebacterium*, we screened skin-associated isolates by sequencing their 16S rRNA gene.  
172 This screen identified eight ribotype B isolates for further study. In total, we shotgun sequenced  
173 40 isolates in the *C. tuberculostearicum* complex– 30 from ribotype A, 8 from ribotype B and 2  
174 from other ASVs. Initial genome clustering using mash indicated that some of the isolates we  
175 sequenced were closely related. Therefore, we used dRep (24) to identify groups of highly similar  
176 genomes (ANImf > 99.5%) and chose the best representative genome for each genome set based  
177 on sequence assembly statistics: maximal N50, minimal number of contigs, and maximal  
178 genome size. In cases of comparable assembly quality, genomes were selected to increase  
179 body site representation. This resulted in a final set of 23 dereplicated *C. tuberculostearicum*  
180 complex genomes, with 18 from ribotype A, four from ribotype B, and one from ASV30.

181

182 **Whole-genome features of five complete *C. tuberculostearicum* complex genomes** In  
183 addition to a paucity of *C. tuberculostearicum* reference genomes at the time of this study, the  
184 ones that did exist were not associated with publications describing their general features. To  
185 address this, we selected five of the dereplicated genomes, three ribotype A and two ribotype B,  
186 for long-read sequencing on the PacBio platform. The subsequent finished or complete *C.*  
187 *tuberculostearicum* complex genomes revealed four copies of the 16S rRNA gene in each  
188 genome. For each genome, we performed a multi-sequence alignment containing each V1-V3  
189 region copy along with the predominant ribotype A sequence first identified in our amplicon  
190 sequencing dataset. Within a genome, copies of the V1-V3 region are almost entirely identical  
191 across alignment, with the ribotype A *C. tuberculostearicum* complex genomes carrying four  
192 copies of ASV3. Admittedly, one exception was a single nucleotide variant identified in one copy  
193 of 16S rRNA 5' region of CTNIH12 (Fig. S3). Notably, no variation was found in the ribotype B  
194 genomes, which both carried four identical gene copies marked by the two characteristic  
195 sequence variants as identified in the amplicon sequencing dataset. The within-genome  
196 homogeneity of 16S rRNA genes confirmed its usefulness as a marker.

197 These complete *C. tuberculostearicum* complex genomes also enabled us to directly  
198 compare the type strain (DSM 44922/FDAARGOS\_1117; human bone marrow) to our ribotype  
199 A and B isolates without the ambiguity introduced by unfinished genomes. Supplemental Figure  
200 4 shows that the five PacBio genomes from this study were largely co-linear, with >80% of the  
201 genome in large syntenic blocks, with the type strain DSM 44922. All five of the genomes in this  
202 study had a 440 kb region that was reorganized relative to the type strain. This region,  
203 comprising around 17.6% of the genome encoded 392 genes (387 coding). The breakpoints for  
204 inversions or translocated blocks in the reference were marked by mobile element families (e.g.,  
205 IS3, IS256, IS481, IS6) that suggest a mechanism for these rearrangements.

206 We extracted the methylation profiles from the PacBio reads of our five genomes  
207 (Supplementary Table S2). The most common methylation pattern, found in all five genomes,  
208 was N6-methyladenine modification (m6A) of GATC motifs (~16,000 sites/genome; 99%  
209 methylated), typically associated the Dam methylase. A second motif AAAAC was also found to  
210 be methylated (m6A) in all five genomes (~75% methylated). In addition to these two ubiquitous  
211 methylation patterns, ribotype A isolate CTNIH10 had two additional methylated motifs present  
212 in hundreds of copies (GGCANNNNNATC, GATDNNNNTGCC). CTNIH20 had an additional  
213 three methylated motifs present at 520-1626 sites/genome. Finally, CTNIH23 had evidence of  
214 an additional five methylation motifs across the genome that were all >98% methylated and

215 present at 225-2159 sites/genome. Methylation systems are important for horizontal gene  
216 transfer, phage resistance and potential recombinant engineering of these strains.

217 The presence of CRISPR-Cas as well as other phage defense systems pose additional  
218 barriers to horizontal gene transfer (HGT). We detected an eight-gene Type I-E Cas gene  
219 cluster and two large repeat arrays (24 and 19 spacers) in the CTNIH20 ribotype B genome, but  
220 not in any of the other full-length genomes from this study. Additional CRISPR-Cas systems  
221 were detected in the short-read assemblies of CTNIH9 (ribotype A; Type I-E Cas gene cluster, 8  
222 spacer CRISPR) and CTNIH22 (ribotype B; Type I-E Cas gene cluster, 12 spacer CRISPR).  
223 Prior to this, the only public *C. tuberculostearicum* complex reference genome with a CRISPR-  
224 Cas system was strain SK141 (ACVP01). A variety of other defense systems including  
225 restriction modification systems were also identified using the DefenseFinder tool  
226 (Supplementary Table S2) (25, 26).

227 Plasmids are important for the mobilization of virulence factors, antibiotic resistance  
228 genes and as tools for recombinant engineering. A single 4.2 kb plasmid was deposited in the  
229 public databases associated with the *C. tuberculostearicum* species complex, p1B146  
230 (NC\_014912) (27). Across the five long-read genomes sequenced here, we detected five  
231 plasmids, none of which aligned to p1B146. Two plasmids, pCT3-020e and pCT4-9116 from  
232 CTNIH23 and CTNIH12, had the same backbone as the *C. diphtheriae* plasmid pNG2 (*ORF9-  
233 traA-ORF11-parAB-repA*) but lacked the erythromycin resistance cassette. CTNIH20 carries  
234 three plasmids ranging in size from 21.2 kb to 27.3 kb. All three carried a *traA/recD2* ortholog  
235 encoding a relaxase/helicase but are otherwise unrelated. While most of the proteins on these  
236 three plasmids were annotated as hypothetical, pCT1-afe7 carried an *ebrB* efflux pump and  
237 pCT1-0563 carried an *stp* (spectinomycin/tetracycline) efflux pump predicted to be involved in  
238 resistance to dyes and antibiotics. All five plasmids were characterized by the presence of a  
239 TraA/RecD2 encoding gene, suggesting a common mobility mechanism. Furthermore, we found  
240 fragments of these plasmids in many of the contig-level genomes. For instance, the *C.  
241 tuberculostearicum* CIP 102622 genome (JAEHFL01) carried both the *stp* gene and a nearby  
242 transcription factor (>99.6% identity) on a 9.8 kb contig, showing the value of these plasmid  
243 references for identifying HGT elements.

244

#### 245 **Taxonomic structure of the skin-associated *C. tuberculostearicum* species complex**

246 While 16S rRNA amplicon sequencing enabled us to group *C. tuberculostearicum* species  
247 complex isolates into two predominant ASVs, the dereplicated genomes enable further high-  
248 resolution taxonomic analysis of this species complex. We used GET\_HOMOLOGUES to



249 extract core genes and build a phylogenetic tree based on core genome SNPs (Fig. 2). We  
250 noted additional taxonomic structure particularly amongst ribotype A isolates. The five public  
251 reference genomes were in the ribotype A-dominated portion of the tree as expected based on  
252 their 16S rRNA gene sequence.

253 While there was good correlation between 16S rRNA ASVs and the core SNP tree, we  
254 noted a single isolate, CTNIH19, which carried a ribotype A allele but localized with the ribotype  
255 B isolates on the tree. CTNIH19 was isolated from the inguinal crease and is the most basal  
256 member of this clade. Work by Cappelli and colleagues recently defined a number of new  
257 *Corynebacterium* species and CTNIH19 was >99.9% identical to a species they designate *C.*  
258 *curieae* (28). We calculated the average nucleotide identity across isolates using pyANI (Fig.  
259 S5) and determined that ribotype B isolates share ANI >97% with themselves and <94% with  
260 other *C. tuberculostearicum* complex genomes, including *C. curieae*. We submitted our ribotype  
261 B isolate genomes to the DSMZ type strain genome server (TYGS) (29) to obtain a  
262 taxonomic/nomenclature assignment. TYGS predicted that ribotype B genomes belong to a new  
263 species in both the whole genome and 16S rDNA trees. The closest TYGS references were *C.*  
264 *tuberculostearicum* DSM 44922 and *C. kefirresidentii*. *C. kefirresidentii* was first described in  
265 2017 (30) after isolation from kefir grains but has not been accepted as an official species yet.  
266 Three of our isolates (CTNIH2, CTNIH6, CTNIH14) from three different healthy volunteers were  
267 98% identical to the *C. kefirresidentii* reference, calling into question whether kefir is the only  
268 natural host for this bacterium (22, 23).

269

270 **Pangenome of the skin-derived *C. tuberculostearicum* complex** We performed a  
271 pangenomic analysis to describe the coding diversity of the *C. tuberculostearicum* species  
272 complex (Figs. 3). We generated an anvio (31) pangenomic map to illustrate genomic variation  
273 across the combined set (N=28) of NCBI reference genomes and our dereplicated genomes.  
274 (Fig. 3A). Pangenome openness was estimated using the Heap's law model (Fig. 3B) as  
275 proposed by Tettelin et al (32). The model indicated an open pangenome ( $0.30, \pm 0.1, \gamma > 0$ ),  
276 predicting that the *C. tuberculostearicum* pangenome would increase with more genomes  
277 analyzed. With our additional 23 genomes, the total pangenome size increases from 3080  
278 genes to 5451 genes, resulting in an expansion of the non-core, or accessory genome by over  
279 300% (Fig. 3C). We performed a functional characterization of 23 lab-sequenced and 5 NCBI-  
280 derived *C. tuberculostearicum* species complex genomes using the eggNOG-mapper  
281 annotation tool (Fig. S6), which returned annotations for 83.2% of orthologous gene clusters (of  
282 which 21% are annotated COG category "S", *Function Unknown*). Interestingly, among the non-

283 core genes, inorganic ion metabolism and transport-related genes were among the most  
284 abundant. We performed a principal components analysis (PCA) of gene presence/absence  
285 data describing our 23 genomes and 2 skin-derived reference sequences (Fig. 4A). We  
286 observed site-specific clustering of genomes isolated from the feet and moist environments. In  
287 addition, we observed distinct clustering of ribotype B isolates (circles) away from other foot-  
288 derived *C. tuberculostearicum* complex genomes, which agreed with the core phylogenetic  
289 clustering. In addition, we identified 11 genes (A=2, B=9) that were unique to and carried by  
290 every member of a ribotype (Supplementary Table. S3), four of which we were able to assign  
291 functional annotation using the UniprotKB sequence similarity search tool, including a  
292 bacteriocin and ferric uptake protein.

293

294 **Improved metagenomic read mapping using an expanded *C. tuberculostearicum***  
295 **pangenome** We tested whether the expanded genomic reference set could improve the rate  
296 of *C. tuberculostearicum* read mapping in a set of metagenomic datasets from 12 healthy  
297 volunteers at 6 body sites (2). Reads were mapped with bowtie2 (33) against a genome  
298 database consisting of the five unique NCBI references or a database of the NCBI references  
299 plus the dereplicated genomes from this study. Overall, 27% more reads were assigned to *C.*  
300 *tuberculostearicum* using the expanded genome set as compared to NCBI references alone.  
301 While the five HVs with isolate genomes in the expanded database showed slightly better  
302 classification, median improvement of 32%, over those without isolates in the database  
303 (median=24%), the difference was not statistically significant, showing the broad utility of these  
304 genomes. Furthermore, when broken down by body site, toenail (Tn) sites showed the largest  
305 improvement in *C. tuberculostearicum* read assignment (72%) while nares, which only  
306 contributed a single genome to the expanded database, improved by 55% (Fig 4B). To control  
307 for spurious read mapping to repetitive elements or other assembly artifacts, these analyses  
308 were repeated using only the predicted gene catalogs, rather than the whole genomes, and very  
309 similar improvements in *C. tuberculostearicum* read mapping were observed, 25% median  
310 improvement and a similar site-dependence.

311

312 **Growth of skin-derived *C. tuberculostearicum* in sweat media** Members of the *C.*  
313 *tuberculostearicum* species complex are widely distributed across the skin's microenvironments.  
314 Differences in body site physiology and nutrient composition inherent to each niche may provide  
315 selective growth advantages (and disadvantages) to a subset of strains. We performed a pilot  
316 experiment to investigate differential growth phenotypes of *C. tuberculostearicum* species

317 complex ribotypes in skin-like media. (Fig 5) *Corynebacterium* are often cultured on brain-heart  
318 infusion (BHI) media plates supplemented with 1% Tween-80 (BHI + 1% Tween80) so this  
319 media was used as a positive control for growth in liquid medium (Fig. 5B). Isolates were  
320 cultured on two types of medium consisting of a complex mixture of amino acids, lipids, and  
321 other metabolites that mimic human eccrine sweat, with one medium supplemented to include a  
322 sebum-like synthetic lipid mixture. Both simulated sweat medias were supplemented with 0.1%  
323 Tween-80. A collection of eight skin-derived *Corynebacterium* strains consisting of four ribotype  
324 A and four ribotype B strains were grown for 20 hours in triplicate and in two separate  
325 experiments for each strain and medium condition (N=6) (Fig. 5 B-D). In all three growth  
326 conditions, ribotype B isolates demonstrated a lower mean OD<sub>600</sub> over time than ribotype A  
327 isolates (Fig. 5A). Using ANOVA and the Tukey method, we determined that the area-under-  
328 the-curve (AUC) difference between the two ribotypes is statistically significant ( $p < 0.0001$ ) for  
329 all media conditions. This pattern was particularly pronounced in the BHI + 1% Tween80 and  
330 Sweat media + 0.1% Tween80 conditions. We observed that the addition of synthetic lipid  
331 mixture to the eccrine sweat-like medium attenuated, however still maintained the growth  
332 difference between ribotype B and other strains, suggesting lipid-limited growth for members of  
333 ribotype B.

334

## 335 DISCUSSION

336 In this study, we investigated the genomic diversity of the predominant yet under-  
337 sequenced *Corynebacterium* genus. Our survey of microbial diversity across human skin  
338 revealed niche-specific enrichment of *Corynebacterium* species and identified *C.*  
339 *tuberculostearicum* as a predominant and widespread species on human skin. Our amplicon-  
340 based analysis was able to identify a site-specific novel 16S rRNA gene ribotype which led to an  
341 expanded sequencing of the *C. tuberculostearicum* species complex. In total, we sequenced 23  
342 distinct isolates belonging to the *C. tuberculostearicum* species complex including *C.*  
343 *tuberculostearicum* (n=15), *C. kefirresidentii* (n=3), *C. curieae* (n=1) and a novel species we are  
344 calling *C. hallux* (n=4). Discovery of *C. kefirresidentii* on human skin and nares suggests that  
345 humans are a natural host for this species.

346 *C. hallux* is likely a new species of skin-associated *Corynebacterium* and merits further  
347 work to formally name it. It was cultured from three different healthy volunteers, detected by  
348 amplicon sequencing in most HVs, represented in the recently published SMGC (SMGC\_122)  
349 (16) and detected in public 16S rRNA gene databases entries associated with skin. In our  
350 healthy volunteers, it was enriched in sites on the feet, particularly the toenail and toe web.

351 Microbial communities on the feet are highly diverse and relatively unstable (2) subject to  
352 temperature fluctuations and invasion by environmental microorganisms.

353 This study helps to resolve the diversity of *C. tuberculostearicum* species complex  
354 strains and provides an important genetic resource for future study. Our whole-genome  
355 sequencing uncovered insights into the genetic diversity of the complex and improved read-  
356 mapping overall by >24%, which will in turn bolster future sequencing efforts and lead to better  
357 characterization of *Corynebacterium* across human skin. While our bioinformatic analysis  
358 greatly expands the non-core genome, a significant proportion of these genes are putative and  
359 lack definitive annotation. Overall, we did not detect obvious gene-level differences between  
360 ribotype B and other strains that would explain the observed differences in site distribution  
361 pattern and growth on synthetic media. Only 11 genes perfectly segregated the two ribotypes  
362 and limitations of functional annotation tools resulted in only hypothetical functional annotations.

363 Our pangenomic analysis did not reveal major metabolic pathways or modules that  
364 differed between ribotype A and B isolates that would explain niche specificity, however there  
365 were two examples of genes with the potential to affect within-niche competition. One of the  
366 genes specific to ribotype B shared sequence similarity with a Lactococcin 972 family  
367 bacteriocin. Bactericidal activity of ribotype B against closely related strains could contribute to  
368 patterns of within-site dominance as observed between ribotypes (Fig. S2). Bactericidal  
369 peptides have recently gained interest as a possible therapeutic intervention for gastrointestinal  
370 disease (34). Furthermore, *Corynebacterium* have been shown to be enriched in a recent study  
371 (35) of post-operative, healing wounds, suggesting an opportunity for biotherapeutic  
372 applications. We also identified a ribotype B-unique copy of a gene encoding ferrous iron  
373 transport protein B, a major regulator of bacterial iron uptake. Iron is an essential nutrient for  
374 survival, requiring the development of highly-efficient sequestering mechanisms by pathogenic  
375 and avirulent bacteria alike (36, 37). Under conditions of limited nutrient bioavailability,  
376 enhanced ferric uptake may prove to be a determining factor of intraspecies competition.

377 On both rich and skin-like media, we observed that ribotype B strains grew less robustly  
378 compared to other strains. Thus, different strains may perform unique roles within their  
379 respective niches. The observed strain-specific distribution pattern may arise from selective  
380 growth advantages including differences in nutritional requirements or nutrient acquisition  
381 mechanisms between strains. Understanding the mechanisms of this variability has important  
382 clinical implications. For example, further characterizing the nutritional limits for sustained  
383 growth may lead to prebiotic therapeutics to augment the growth of beneficial strains within a  
384 given microenvironment, or engineering site-specific, microbe-based drug delivery systems.

385 Understanding the roles and requirements of host-associated microbial communities in  
386 maintaining skin health will provide insight into the emergence of skin disorders in addition to  
387 novel therapeutic interventions to combat them.

388

## 389 **Materials and Methods**

390 **Subject recruitment and sampling** Healthy adult male and female volunteers (HVs)  
391 18–40 years of age were recruited from the Washington, DC metropolitan region. This natural  
392 history study was approved by the Institutional Review Board of the National Human Genome  
393 Research Institute ([clinicaltrials.gov/NCT00605878](https://clinicaltrials.gov/NCT00605878)) and the National Institute of Arthritis and  
394 Musculoskeletal and Skin Diseases (<https://clinicaltrials.gov/ct2/show/NCT02471352>) and all  
395 subjects provided written informed consent prior to participation. Sampling was performed as  
396 described previously (20).

397 **16S rRNA gene sequencing** 16S rRNA gene amplicon sequencing of these samples  
398 has been described previously (5). Briefly, each DNA sample was amplified with universal  
399 primers flanking variable regions V1 (27F, 5'-AGAGTTTGATCCTGGCTCAG) and V3 (534R, 5'-  
400 ATTACCGCGGCTGCTGG). For each sample, the universal primers were tagged with unique  
401 indexes to allow for multiplexing/demultiplexing (38). The following PCR conditions were used: 2  
402  $\mu$ l 10X AccuPrime Buffer II, 0.15  $\mu$ l Accuprime Taq (Invitrogen, Carlsbad, CA), 0.04  $\mu$ l  
403 adapter+V1\_27F (100  $\mu$ M), 2  $\mu$ l primer V3\_354R+barcode (2  $\mu$ M), and 2  $\mu$ l of isolated microbial  
404 genomic DNA. PCR was performed in duplicate for 30 cycles followed by PCR-clean up and  
405 amplicon pooling of ~10 ng DNA. Duplicate amplicons were combined, purified (Agencourt  
406 AMPure XP-PCR Purification Kit (Beckman Coulter, Inc., Brea, CA)), and quantified (QuantIT  
407 dsDNA High-Sensitivity Assay Kit (Invitrogen, Carlsbad, CA)). An average of ~8 ng DNA of 94  
408 amplicons were pooled together, purified (MinElute PCR Purification Kit (Qiagen, Valencia, CA))  
409 and sequenced on a Roche 454 GS20/FLX platform with Titanium chemistry (Roche, Branford,  
410 Connecticut). Flow-grams were processed with the 454 Basecalling pipeline (v2.5.3).

411 **16S rRNA gene amplicon analysis** Sequencing data were processed using DADA2  
412 v1.20.0 (39). Sequences were filtered and trimmed as recommended by the software  
413 developers and truncated to 375 nt: `filterAndTrim(fnFs, filtFs,maxN=0, maxEE=c(2),`  
414 `truncQ=2,truncLen=c(375))`. Sample inference was performed using the `learnErrors`  
415 (`randomize=TRUE`) and the `dada` (`HOMOPOLYMER_GAP_PENALTY=-1, BAND_SIZE=32`)  
416 commands. Chimeras were removed using `removeBimeraDenovo` (`method="consensus",`  
417 `allowOneOff=TRUE`). Taxonomy was assigned using `assignTaxonomy` (`minBoot=70`) command  
418 in DADA2 with the Refseq (<https://zenodo.org/record/3266798>) or eHOMD v15.1 V1V3 (3)

419 training set databases. The resulting amplicon sequence variants (ASVs), taxonomy and  
420 sample metadata were used to build a phyloseq (40) object that was used for further analysis.

421 **Bacterial culturing** *Corynebacterium* isolates were cultured from healthy volunteers as  
422 previously described (16). Briefly, skin samples were collected with eSWabs (COPAN e480C) in  
423 liquid Amies. Samples were diluted and plated on brain heart infusion agar with 1% Tween 80.  
424 Potential *Corynebacterium* isolates were taxonomically classified by amplifying and Sanger  
425 sequencing the full length 16S rRNA gene with primers (8F, 5'-AGAGTTTGATCCTGGCTCAG)  
426 and (1391R, 5'-GACGGGCGGTGWGTRCA).

427 **Bacterial whole genome sequencing** Genomic DNA was purified for each isolate,  
428 from which Nextera XT (Illumina) libraries were generated. Each isolate was sequenced using a  
429 2x151 paired-end dual index run on an Illumina NovaSeq 6000. The reads were subsampled to  
430 achieve 80-100x coverage using seqtk (version 1.2), assembled with SPAdes (version 3.14.1)  
431 (41) and polished using bowtie2 (version 2.2.6) and Pilon (version 1.23) (42). To achieve full  
432 reference genomes for select isolates, genomic DNA was sequenced on the PacBio Sequel II  
433 platform (version 8M SMRTCells, Sequel II version 2.0 sequencing reagents, 15 hr movie  
434 collection). The subreads were assembled using Canu v2.1 and polished using the  
435 pb\_resequencing workflow within PacBio SMRTLink v.9.0.0.92188. Genome annotation was  
436 performed using National Center for Biotechnology Information (NCBI) Prokaryotic Genome  
437 Annotation Pipeline (PGAP: [https://www.ncbi.nlm.nih.gov/genome/annotation\\_prok/](https://www.ncbi.nlm.nih.gov/genome/annotation_prok/)).  
438 Methylation patterns for the assembled genomes were determined using the pb\_basemods  
439 workflow in SMRTLink v.9.0.0.92188. Whole genome and plasmid alignments were generated  
440 in mummer (v3.9.4alpha) and visualized in R.

441 Full-length 16S rRNA gene copies were extracted from each PacBio complete genome.  
442 Briefly, reference *Corynebacterium* 16S rRNA sequences were downloaded from the RDP  
443 database (Good quality, >1200 nt) and used as a BLAST database to identify the coordinates of  
444 the four copies in the genome. To detect intragenomic variation in the 16S rRNA gene, all  
445 copies within each genome were compared against each other using the EMBL-EBI Multiple  
446 Sequence Alignment Tool (MUSCLE). Whole genome alignments were generated in Mauve v  
447 2.4.0.

448 **Phylogenetic analysis** Publicly available genomes were downloaded from NCBI  
449 including *C. tuberculostearicum* (CP068156, CP06979, CP065972, ACVP01, JAEHFL01), *C.*  
450 *kefirresidentii* (CP067012, JAHXPF01), *C. curieae* (JAKMUU01) and *C. accolens* (ACGD01).  
451 GET\_HOMOLOGUES (v09212021) was used to cluster protein sequences from 29 genomes  
452 (28 *C. tuberculostearicum*, 1 *C. accolens*) into orthologous groups and generate a core gene

453 alignment. Prokka GBK files were used as input for clustering. The OrthoMCL (v1.4) option was  
454 used to group sequences utilizing the Markov Clustering Algorithm with a minimum coverage  
455 value of 90% in blast pairwise protein alignments. A strict core consensus genome was  
456 generated by calculating the intersection of single copy genes present in all 29 genomes. The  
457 accompanying GET\_PHYLOMARKERS (v. 2.2.9.1) pipeline was used to identify markers for  
458 phylogenetic inference. IQTREE (v 2.1.2) was used to generate a maximum-likelihood  
459 phylogenetic tree from marker gene cluster alignments with 1000 bootstrap replicates. and a  
460 mean branch support value cutoff of 0.7. The top-scoring tree was visualized and annotated  
461 using the web-based program interactive Tree of Life (iTOL v6). The average nucleotide identity  
462 (ANIb) matrix for all sequences was plotted and annotated using the package heatmap.2/R.

463 **Gene calling and Annotation** The Prokka (v1.14.6) pipeline was used for gene calling  
464 and annotation. GFF3- and GBK- format annotations were generated for 28 *Corynebacterium*  
465 *tuberculostearicum* sequences derived from 23 lab isolates and 5 NCBI references (CIP  
466 102622, FDAARGOS 993, FDAARGOS 1117, FDAARGOS 1198, SK141), in addition to the  
467 *Corynebacterium accolens* representative genome (ATCC 49725) sequence.

468 **Pangenome calculation** Three *C. tuberculostearicum* pangenomes (for all reference  
469 sequences, skin-derived reference sequences and lab sequences, and all sequences) were  
470 calculated from Prokka-derived GFF3 files using Panaroo on sensitive mode with a sequence  
471 identity threshold of 90% and otherwise default parameters. The resultant gene  
472 presence/absence tables were used for downstream analysis.

473 **Pangenome visualization** A pangenomic map was created using anvio (v. 7.1) with  
474 imported Prokka gene calling information and annotations (GBK format) for 23 lab-sequenced  
475 and 5 NCBI reference *C. tuberculostearicum* genomes. Strains were annotated with sample  
476 metadata including skin site, general skin habitat, healthy volunteer ID, as well as phylogenetic  
477 grouping from GET\_HOMOLOGUES analysis. Average nucleotide identity (ANIb) of aligned  
478 regions was calculated within anvio using pyANI. In addition, eggNOG (v. 2.1.7) gene  
479 annotations were used for gene cluster annotation with the NCBI COG Database (2020) and  
480 visualized using ggplot2/R. Core, accessory, and singleton gene counts were derived from gene  
481 presence/absence tables for (1) all reference sequences and (2) all sequences. Counts were  
482 visualized as pie charts.

483 A gene rarefaction curve for the *C. tuberculostearicum* pangenome was found by  
484 applying the Vegan/R specaccum() function to the gene presence/absence table, with a random  
485 order of additions of genomes permuted 1000 times. A Heap's law power law model was fitted

486 to the curve using the nls function in stats/R to calculate constants K and  $\alpha$ . The curve was  
487 visualized using ggplot2/R.

488 A principal components analysis was performed on the gene presence/absence table  
489 using the prcomp function in stats(v 3.6.2)/R. The resultant object was visualized using  
490 ggplot2/R.

491 **Unique genes** Scoary (v. 1.6.16) was used to identify genes unique to ribotype A and B  
492 for 23 lab-sequenced *C. tuberculostearicum* complex isolates. Gene sequences were queried  
493 using the UniProtKB online sequence similarity BLAST tool (<https://www.uniprot.org/blast>).

494 **Metagenomic read mapping** Metagenomic reads from 12 healthy volunteers at 6 body  
495 sites, adapter trimmed and host subtracted as described in (2), were aligned to a bowtie2 (v 2-  
496 2.4.5) database built from five NCBI *C. tuberculostearicum* genomes (CIP\_102622, DSM 44922,  
497 FDAARGOS\_1198, FDAARGOS\_993, SK141) with or with supplementation with 23 non-  
498 redundant genomes from this study; default bowtie2 parameters were used: --end-to-end --  
499 sensitive.

500 **Growth curve starter cultures** Isolates for differential growth analysis were selected  
501 on the basis of i) reliable growth in Brain-Heart-Infusion+Tween80 and ii) coverage of the  
502 phylogenetic tree. *C. tuberculostearicum* isolates were grown in overnight liquid culture  
503 consisting of BHI broth (Sigma-Aldrich), augmented with 1% RPI Tween80, and 40ug/ml  
504 Fosfomycin (BHI-T-F) at 37°C with shaking at 220 rpm. To make the “Sweat media + 0.1%  
505 Tween80” media, we filter-sterilized Pickering Artificial Eccrine Perspiration Cat. No. 1700-  
506 0023 (pH 6.5) and RPI Tween80 (1%). This medium was then vortex-combined with 1% volume  
507 of synthetic apocrine sweat (Pickering Cat. No. 1700-070X) to produce the “Sweat media +  
508 0.1% Tween80 + synthetic lipid mixture” medium.

509 **Differential growth experiments** *C. tuberculostearicum* liquid cultures were pelleted,  
510 washed and diluted 10-fold in diH2O to an OD<sub>600</sub> of ~0.1. Differential media were inoculated with  
511 diluted culture at a concentration of 100:1 and plated in triplicate across a 96-well microplate.  
512 Bacterial growth was recorded using the Epoch 2 Microplate. OD<sub>600</sub> readings were taken at 30  
513 minute intervals throughout a 24-hour time span. The experiment was performed in duplicate.  
514 OD<sub>600</sub> measurements were exported, corrected via blank subtraction, and plotted using  
515 ggplot2/R. The package growthcurver/R was used to calculate empirical area under the curve  
516 (AUC) for all isolate:media combinations. Statistical significance testing for ribotype:media  
517 interactions were performed using ANOVA and a post-hoc Tukey test.

518 **Data Availability** Genome data are deposited under the NCBI BioProjects  
519 PRJNA854648, PRJNA694925 and PRJNA854648 (see Table S1). Some amplicon data were



520 published previously (n=145; PRJNA46333)(5) and the remainder are new to this study (n=168;  
521 PRJNA46333)

522

## 523 **ACKNOWLEDGEMENTS**

524 We thank Tommy Hiller Tran for valuable bioinformatics discussions. We thank Dr. Matthew  
525 Kelly for discussions of *Corynebacterium* biology. The computational resources of the NIH High-  
526 Performance Computation Biowulf Cluster (<http://hpc.nih.gov>) were used for this study. This  
527 work was supported by the Intramural Research Programs of the National Human Genome  
528 Reseach Institute and the National Institute of Arthritis and Musculoskeletal and Skin Diseases.

529 N.A. performed bioinformatics and growth curve analyses, prepared figures and wrote  
530 the manuscript; P.J. provided technical support for growth curve experiments; C.D. cultured  
531 bacteria and prepared DNA sequencing libraries; NISC sequenced amplicon and whole genome  
532 libraries; K.L. and H.H. discussed results and provided subject-matter expertise; J.S. and S.C.  
533 conceived the overall study and were responsible for the final version of the manuscript. All  
534 authors read and approved the final manuscript.

535

## 536 **Figure Legends**

537 **FIG 1** *Corynebacterium* species relative abundance in normal human skin microbiome. (A)  
538 Relative abundance of the 15 major *Corynebacterium* species across 14 skin sites: sebaceous  
539 (back, Ba; occiput, Oc; external auditory canal, Ea; retroauricular crease, Ra; manubrium, Mb;  
540 glabella, Gb), moist (inguinal crease, Ic; antecubital crease, Ac), dry (hypothenar palm, Hp;  
541 volar forearm, Vf), foot (toe nail, Tn; toe web, Tw; plantar heel, Ph) and (N)ares. Relative  
542 abundances determined by sequencing of the V1-V3 region of the 16S rRNA gene and  
543 subsetting to *Corynebacterium* reads. (B) Percent of total bacterial reads attributed to  
544 *Corynebacterium* and *C. tuberculostearicum* in each skin habitat. Of the six ASVs assigned to  
545 *C. tuberculostearicum*, mean relative abundance across skin habitats.

546

547 **FIG 2** A maximum-likelihood phylogenetic tree of *C. tuberculostearicum* species complex  
548 genomes from this study and publicly available, calculated from 1315 core gene cluster  
549 alignments. Bootstrap values (located along internal nodes) were calculated from 1000  
550 replicates. Clustering was generated using GET\_HOMOLOGUES OrthoMCL v1.4 option with  
551 minimum coverage 90% in BLAST pairwise alignments. The tree was rooted on outgroup *C.*  
552 *accolens* ATCC 49725. On the right of tree, boxes depict site (body site locations defined in

553 Figure 1) and individual (HV) from which each isolate was cultured. Sites are colored by niche  
554 type, with moist in shades of green; feet in shades of orange; dry in pink; sebaceous in  
555 lavender; and nares in blue. Individuals are randomly but consistently colored.

556

557 **FIG 3** The *Corynebacterium tuberculostearicum* pangenome. (A) Anvi'o pangenomic map for 28  
558 *C. tuberculostearicum* genomes (including 5 NCBI reference genomes). Genomic rings are  
559 annotated by skin site and HV (healthy volunteer) metadata and ordered by pyANI average  
560 nucleotide identity (ANIb). Genome margins are manually adjusted for clarity. (B) Heap's Law  
561 estimate of pangenome openness for 28 genomes. A rarefaction curve showing the total  
562 number of genes accumulated with the addition of new genome sequences in random order  
563 with 1000 permutations. Shaded regions represent the 95% confidence interval. A Heap's law  
564 model was fit to the resultant curve to calculate  $k$  and  $\gamma$  values ( $1977 \pm 38.0$  and  $0.30 \pm 0.01$ ,  
565 respectively). (C) Number of core (belonging to all genomes), accessory (belonging to two or  
566 more genomes), and singleton (belonging to only one genome) genes. The expanded  
567 pangenome contains 5451 genes using 90% sequence identity as a cutoff parameter.

568

569 **FIG 4** *C. tuberculostearicum* complex pangenome clustering and improved metagenomic read  
570 mapping. (A) Principal components analysis of orthologous gene clustering. The gene  
571 presence/absence data for 25 genomes (including two NCBI references, shown in gray) was  
572 analyzed using principal components analysis. Ribotype B genomes are shown as circles; other  
573 genomes are triangles. (B) Improvement in shotgun metagenomic read mapping with a 28  
574 member *C. tuberculostearicum* database as compared to the 5 member NCBI database.  
575 Percent increase in mapped *C. tuberculostearicum* reads by body site. Each point is a healthy  
576 volunteer. Triangles mark healthy volunteers that contributed one or more isolates to the  
577 expanded mapping database.

578

579 **FIG 5** Growth phenotypes of select *C. tuberculostearicum* complex strains in synthetic sweat  
580 media. (A) Empirical area under curve comparison of *C. tuberculostearicum* species complex  
581 strains from ribotype A and ribotype B, with biological replicates grouped by color. Strains were  
582 grown in Brain Heart Infusion (BHI) + 1% Tween; Sweat media + 0.1% Tween80; Sweat media  
583 + 0.1% Tween80 + synthetic lipid mixture. Medium composition is described in further detail in  
584 Methods. (B-D) Selected growth curves from a representative experiment plotted with standard  
585 error. Ribotype B isolates are shown in shades of blue; Ribotype A isolates are shown in shades  
586 of red.

587

588 **FIG S1** Relative abundances of 4 major bacterial phyla across 14 skin sites from normal human  
589 volunteers. Relative abundances determined by sequencing the V1-V3 region of 16S rRNA  
590 followed by classification with the DADA2 and the eHOMD v15.1 database.

591

592 **FIG S2** Relative abundance of the major *C. tuberculostearicum* ASVs across 14 skin sites in  
593 normal human volunteers.

594

595 **FIG S3** A schematic representation of 16S rRNA intra-genome variation. Gray rectangles  
596 represent the complete V1-V3 region of each of the four 16S rRNA gene copies found in each  
597 genome. Variant positions are marked with an X and colored blue for variants characteristic of  
598 ribotype B, or orange for additional variants or variants that are found outside the trimmed ASV  
599 (red dotted line).

600

601 **FIG S4** Alignment of the *C. tuberculostearicum* reference genome with five PacBio genomes  
602 from this study. Aligned regions are shown as bands colored by the percent identity.

603

604 **FIG S5** Average Nucleotide Identity (ANI) of *C. tuberculostearicum* species complex genomes.  
605 The distance matrix was calculated using fastANI and bidirectional percent identities were  
606 averaged. Distances were hierarchically clustered and visualized using heatmap.2/R. Species  
607 are abbreviated as *C. tuberculostearicum*, Ctub; *C. curieae*, Ccur; *C. kefirresidentii*, Ckef.  
608 Pairwise comparisons at >95% identity are marked with a (\*).

609

610 **Fig S6** Functional classifications of orthologous gene clusters. Dotted line demarcates expected  
611 proportion of core, accessory, and singleton genes assigned to each category given the total  
612 (N=3703 including 243 duplicate and triplicate assignments) and core (N=1738, or 47%) number  
613 of category assignments. Total gene counts per category *N* are labeled.

614

615 **TABLE S1** Table of dereplicated whole genomes. N50 is the N50 contig length, n/a for finished  
616 genomes. The number of plasmids are listed for finished genomes. Except for CTNIH9 (\*), all  
617 genomes are from strains associated with HVs in the microbiome analysis. ASV, Amplicon  
618 Sequence Variant; HV, Healthy Volunteer

619

620 **TABLE S2** Phage defense systems and methylation patterns for complete *Corynebacterium*  
621 *tuberculostearicum* species complex genomes. Defense finder reports restrictions modification  
622 systems (RM) and other phage defense systems. Numbers in parenthesis indicate the number  
623 of systems present, if more than one.

624  
625 **TABLE S3** Ribotype-specific genes. Genes uniquely present among all ribotype A or ribotype B  
626 genomes. All ribotype:gene associations p-values (FDR-adjusted) < 0.05.

627

## 628 REFERENCES

629

- 630 1. Grice E, Kong H, Conlan S, Deming C, Davis J, Young A, NISC Comparative Sequencing  
631 Program, Bouffard G, Blakesley R, Murray P, Green E, Turner M, Segre J. 2009.  
632 Topographical and Temporal Diversity of the Human Skin Microbiome. *Science* 324:1192,  
633 1190.
- 634 2. Oh J, Byrd AL, Park M, NISC Comparative Sequencing Program, Kong HH, Segre JA.  
635 2016. Temporal Stability of the Human Skin Microbiome. *Cell* 165:854–866.
- 636 3. Escapa IF, Chen T, Huang Y, Gajare P, Dewhirst FE, Lemon KP. 2018. New Insights into  
637 Human Nostril Microbiome from the Expanded Human Oral Microbiome Database  
638 (eHOMD): a Resource for the Microbiome of the Human Aerodigestive Tract. *mSystems* 3.
- 639 4. Kong HH, Oh J, Deming C, Conlan S, Grice EA, Beatson MA, Nomicos E, Polley EC,  
640 Komarow HD, NISC Comparative Sequence Program, Murray PR, Turner ML, Segre JA.  
641 2012. Temporal shifts in the skin microbiome associated with disease flares and treatment  
642 in children with atopic dermatitis. *Genome Res* 22:850–859.
- 643 5. Oh J, Freeman AF, NISC Comparative Sequencing Program, Park M, Sokolic R, Candotti F,  
644 Holland SM, Segre JA, Kong HH. 2013. The altered landscape of the human skin  
645 microbiome in patients with primary immunodeficiencies. *Genome Res* 23:2103–2114.

- 646 6. Smeekens SP, Huttenhower C, Riza A, van de Veerdonk FL, Zeeuwen PLJM, Schalkwijk J,  
647 van der Meer JWM, Xavier RJ, Netea MG, Gevers D. 2014. Skin microbiome imbalance in  
648 patients with STAT1/STAT3 defects impairs innate host defense responses. *J Innate Immun*  
649 6:253–262.
- 650 7. Brugger SD, Eslami SM, Pettigrew MM, Escapa IF, Henke MT, Kong Y, Lemon KP. 2020.  
651 *Dolosigranulum pigrum* Cooperation and Competition in Human Nasal Microbiota. *mSphere*  
652 5:e00852-20.
- 653 8. Ramsey MM, Freire MO, Gabriliska RA, Rumbaugh KP, Lemon KP. 2016. *Staphylococcus*  
654 *aureus* Shifts toward Commensalism in Response to *Corynebacterium* Species. *Front*  
655 *Microbiol* 7:1230.
- 656 9. Ridaura VK, Bouladoux N, Claesen J, Chen YE, Byrd AL, Constantinides MG, Merrill ED,  
657 Tamoutounour S, Fischbach MA, Belkaid Y. 2018. Contextual control of skin immunity and  
658 inflammation by *Corynebacterium*. *J Exp Med* 215:785–799.
- 659 10. Sakamoto K, Jin S-P, Goel S, Jo J-H, Voisin B, Kim D, Nadella V, Liang H, Kobayashi T,  
660 Huang X, Deming C, Horiuchi K, Segre JA, Kong HH, Nagao K. 2021. Disruption of the  
661 endopeptidase ADAM10-Notch signaling axis leads to skin dysbiosis and innate lymphoid  
662 cell-mediated hair follicle destruction. *Immunity* 54:2321-2337.e10.
- 663 11. Kobayashi T, Glatz M, Horiuchi K, Kawasaki H, Akiyama H, Kaplan DH, Kong HH, Amagai  
664 M, Nagao K. 2015. Dysbiosis and *Staphylococcus aureus* Colonization Drives Inflammation  
665 in Atopic Dermatitis. *Immunity* 42:756–766.
- 666 12. Altonsy MO, Kurwa HA, Lauzon GJ, Amrein M, Gerber AN, Almishri W, Mydlarski PR. 2020.  
667 *Corynebacterium tuberculostearicum*, a human skin colonizer, induces the canonical

- 668 nuclear factor- $\kappa$ B inflammatory signaling pathway in human skin cells. *Immun Inflamm Dis*  
669 8:62–79.
- 670 13. Conlan S, Mijares LA, NISC Comparative Sequencing Program, Becker J, Blakesley RW,  
671 Bouffard GG, Brooks S, Coleman H, Gupta J, Gurson N, Park M, Schmidt B, Thomas PJ,  
672 Otto M, Kong HH, Murray PR, Segre JA. 2012. *Staphylococcus epidermidis* pan-genome  
673 sequence analysis reveals diversity of skin commensal and hospital infection-associated  
674 isolates. *Genome Biol* 13:R64.
- 675 14. Tomida S, Nguyen L, Chiu B-H, Liu J, Sodergren E, Weinstock GM, Li H. 2013. Pan-  
676 genome and comparative genome analyses of *propionibacterium acnes* reveal its genomic  
677 diversity in the healthy and diseased human skin microbiome. *mBio* 4:e00003-00013.
- 678 15. Flores Ramos S, Brugger SD, Escapa IF, Skeete CA, Cotton SL, Eslami SM, Gao W,  
679 Bomar L, Tran TH, Jones DS, Minot S, Roberts RJ, Johnston CD, Lemon KP. 2021.  
680 Genomic Stability and Genetic Defense Systems in *Dolosigranulum pigrum*, a Candidate  
681 Beneficial Bacterium from the Human Microbiome. *mSystems* 6:e0042521.
- 682 16. Saheb Kashaf S, Proctor DM, Deming C, Saary P, Hölzer M, NISC Comparative  
683 Sequencing Program, Taylor ME, Kong HH, Segre JA, Almeida A, Finn RD. 2022.  
684 Integrating cultivation and metagenomics for a multi-kingdom view of skin microbiome  
685 diversity and functions. *Nat Microbiol* 7:169–179.
- 686 17. Caputo A, Fournier P-E, Raoult D. 2019. Genome and pan-genome analysis to classify  
687 emerging bacteria. *Biol Direct* 14:5.
- 688 18. Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. 2018. High throughput  
689 ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun*  
690 9:5114.

- 691 19. Tettelin H, Massignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, Angiuoli SV,  
692 Crabtree J, Jones AL, Durkin AS, Deboy RT, Davidsen TM, Mora M, Scarselli M, Margarit y  
693 Ros I, Peterson JD, Hauser CR, Sundaram JP, Nelson WC, Madupu R, Brinkac LM,  
694 Dodson RJ, Rosovitz MJ, Sullivan SA, Daugherty SC, Haft DH, Selengut J, Gwinn ML, Zhou  
695 L, Zafar N, Khouri H, Radune D, Dimitrov G, Watkins K, O'Connor KJB, Smith S, Utterback  
696 TR, White O, Rubens CE, Grandi G, Madoff LC, Kasper DL, Telford JL, Wessels MR,  
697 Rappuoli R, Fraser CM. 2005. Genome analysis of multiple pathogenic isolates of  
698 *Streptococcus agalactiae*: implications for the microbial “pan-genome.” *Proc Natl Acad Sci*  
699 *USA* 102:13950–13955.
- 700 20. Findley K, Oh J, Yang J, Conlan S, Deming C, Meyer JA, Schoenfeld D, Nomicos E, Park  
701 M, NIH Intramural Sequencing Center Comparative Sequencing Program, Kong HH, Segre  
702 JA. 2013. Topographic diversity of fungal and bacterial communities in human skin. *Nature*  
703 498:367–370.
- 704 21. Oliveira A, Oliveira LC, Aburjaile F, Benevides L, Tiwari S, Jamal SB, Silva A, Figueiredo  
705 HCP, Ghosh P, Portela RW, De Carvalho Azevedo VA, Wattam AR. 2017. Insight of Genus  
706 *Corynebacterium*: Ascertainning the Role of Pathogenic and Non-pathogenic Species. *Front*  
707 *Microbiol* 8:1937.
- 708 22. Salamzade R, Cheong JZA, Sandstrom S, Swaney MH, Stubbendieck RM, Starr NL, Currie  
709 CR, Singh AM, Kalan LR. 2023. Evolutionary investigations of the biosynthetic diversity in  
710 the skin microbiome using *IsaBGC*. *Microb Genom* 9:mgen000988.
- 711 23. Salamzade R, Swaney MH, Kalan LR. 2022. Comparative Genomic and Metagenomic  
712 Investigations of the *Corynebacterium tuberculostearicum* Species Complex Reveals  
713 Potential Mechanisms Underlying Associations To Skin Health and Disease. *Microbiol*  
714 *Spectr* e0357822.

- 715 24. Olm MR, Brown CT, Brooks B, Banfield JF. 2017. dRep: a tool for fast and accurate  
716 genomic comparisons that enables improved genome recovery from metagenomes through  
717 de-replication. *ISME J* 11:2864–2868.
- 718 25. Tesson F, Hervé A, Touchon M, d’Humières C, Cury J, Bernheim A. 2021. Systematic and  
719 quantitative view of the antiviral arsenal of prokaryotes. *bioRxiv*  
720 <https://doi.org/10.1101/2021.09.02.458658>.
- 721 26. Abby SS, Néron B, Ménager H, Touchon M, Rocha EPC. 2014. MacSyFinder: a program to  
722 mine genomes for molecular systems with an application to CRISPR-Cas systems. *PLoS*  
723 *One* 9:e110726.
- 724 27. Wieteska Ł, Szewczyk EM, Szemraj J. 2011. Characterization of novel plasmid p1B146  
725 from *Corynebacterium tuberculostearicum*. *J Microbiol Biotechnol* 21:796–801.
- 726 28. Cappelli EA, Ksiezarek M, Wolf J, Neumann-Schaal M, Ribeiro TG, Peixe L. 2023.  
727 Expanding the Bacterial Diversity of the Female Urinary Microbiome: Description of Eight  
728 New *Corynebacterium* Species. *Microorganisms* 11:388.
- 729 29. Meier-Kolthoff JP, Carbasse JS, Peinado-Olarte RL, Göker M. 2022. TYGS and LPSN: a  
730 database tandem for fast and reliable genome-based classification and nomenclature of  
731 prokaryotes. *Nucleic Acids Res* 50:D801–D807.
- 732 30. Blasche S, Kim Y, Patil KR. 2017. Draft Genome Sequence of *Corynebacterium*  
733 *kefirresidentii* SB, Isolated from Kefir. *Genome Announc* 5:e00877-17.
- 734 31. Eren AM, Kiefl E, Shaiber A, Veseli I, Miller SE, Schechter MS, Fink I, Pan JN, Yousef M,  
735 Fogarty EC, Trigodet F, Watson AR, Esen ÖC, Moore RM, Clayssen Q, Lee MD, Kivenson  
736 V, Graham ED, Merrill BD, Karkman A, Blankenberg D, Eppley JM, Sjödin A, Scott JJ,

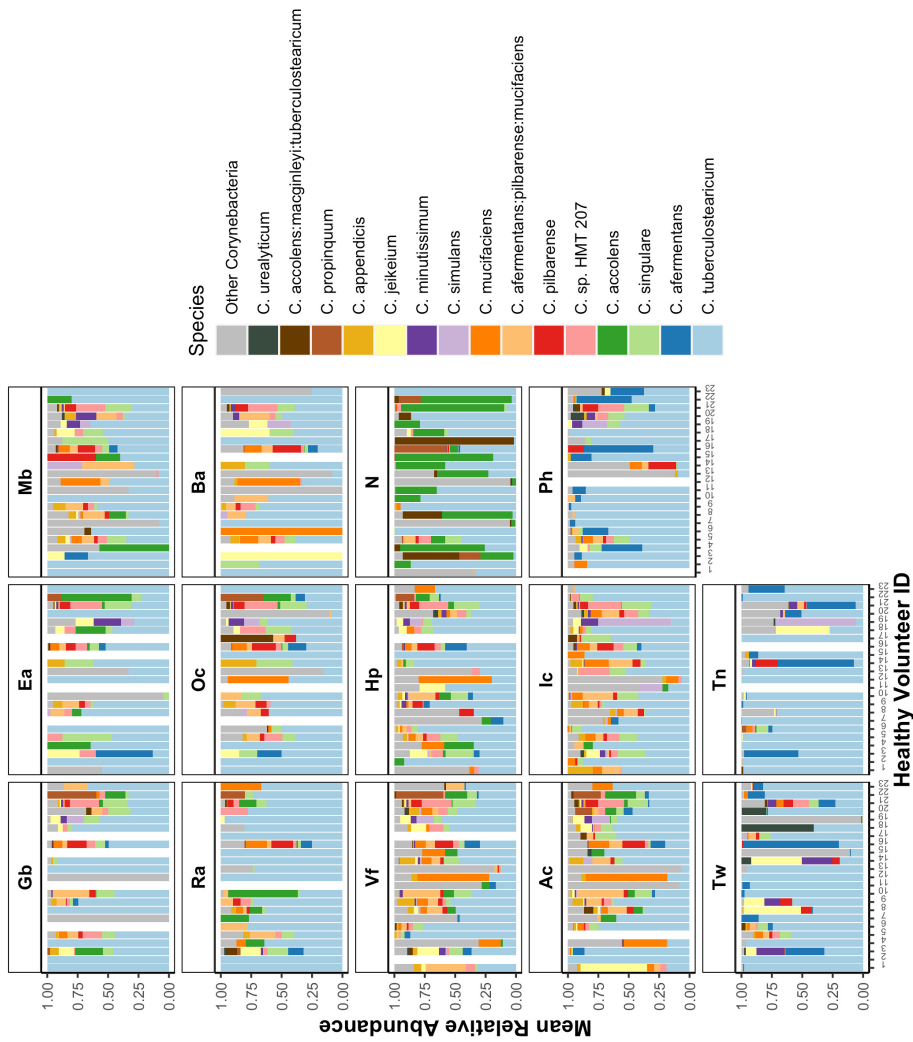


- 737 Vázquez-Campos X, McKay LJ, McDaniel EA, Stevens SLR, Anderson RE, Fuessel J,  
738 Fernandez-Guerra A, Maignien L, Delmont TO, Willis AD. 2021. Community-led, integrated,  
739 reproducible multi-omics with anvio. *Nat Microbiol* 6:3–6.
- 740 32. Tettelin H, Riley D, Cattuto C, Medini D. 2008. Comparative genomics: the bacterial pan-  
741 genome. *Curr Opin Microbiol* 11:472–477.
- 742 33. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods*  
743 9:357–359.
- 744 34. Lopetuso LR, Giorgio ME, Saviano A, Scaldaferri F, Gasbarrini A, Cammarota G. 2019.  
745 Bacteriocins and Bacteriophages: Therapeutic Weapons for Gastrointestinal Diseases? *Int J*  
746 *Mol Sci* 20:183.
- 747 35. Gupta S, Poret AJ, Hashemi D, Eseonu A, Yu SH, D’Gama J, Neel VA, Lieberman TD.  
748 2022. Cutaneous Surgical Wounds Have Distinct Microbiomes from Intact Skin. *Microbiol*  
749 *Spectr* e0330022.
- 750 36. Nairz M, Schroll A, Sonnweber T, Weiss G. 2010. The struggle for iron - a metal at the host-  
751 pathogen interface. *Cell Microbiol* 12:1691–1702.
- 752 37. Lau CKY, Krewulak KD, Vogel HJ. 2016. Bacterial ferrous iron transport: the Feo system.  
753 *FEMS Microbiol Rev* 40:273–298.
- 754 38. Lennon NJ, Lintner RE, Anderson S, Alvarez P, Barry A, Brockman W, Daza R, Erlich RL,  
755 Giannoukos G, Green L, Hollinger A, Hoover CA, Jaffe DB, Juhn F, McCarthy D, Perrin D,  
756 Ponchner K, Powers TL, Rizzolo K, Robbins D, Ryan E, Russ C, Sparrow T, Stalker J,  
757 Steelman S, Weiland M, Zimmer A, Henn MR, Nusbaum C, Nicol R. 2010. A scalable, fully

- 758 automated process for construction of sequence-ready barcoded libraries for 454. *Genome*  
759 *Biol* 11:R15.
- 760 39. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. 2016. DADA2:  
761 High-resolution sample inference from Illumina amplicon data. *Nat Methods* 13:581–583.
- 762 40. McMurdie PJ, Holmes S. 2013. phyloseq: an R package for reproducible interactive analysis  
763 and graphics of microbiome census data. *PLoS One* 8:e61217.
- 764 41. Nurk S, Bankevich A, Antipov D, Gurevich AA, Korobeynikov A, Lapidus A, Prjibelski AD,  
765 Pyshkin A, Sirotkin A, Sirotkin Y, Stepanauskas R, Clingenpeel SR, Woyke T, McLean JS,  
766 Lasken R, Tesler G, Alekseyev MA, Pevzner PA. 2013. Assembling single-cell genomes  
767 and mini-metagenomes from chimeric MDA products. *J Comput Biol* 20:714–737.
- 768 42. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,  
769 Wortman J, Young SK, Earl AM. 2014. Pilon: an integrated tool for comprehensive microbial  
770 variant detection and genome assembly improvement. *PLoS ONE* 9:e112963.
- 771  
772

Figure 1

A



B

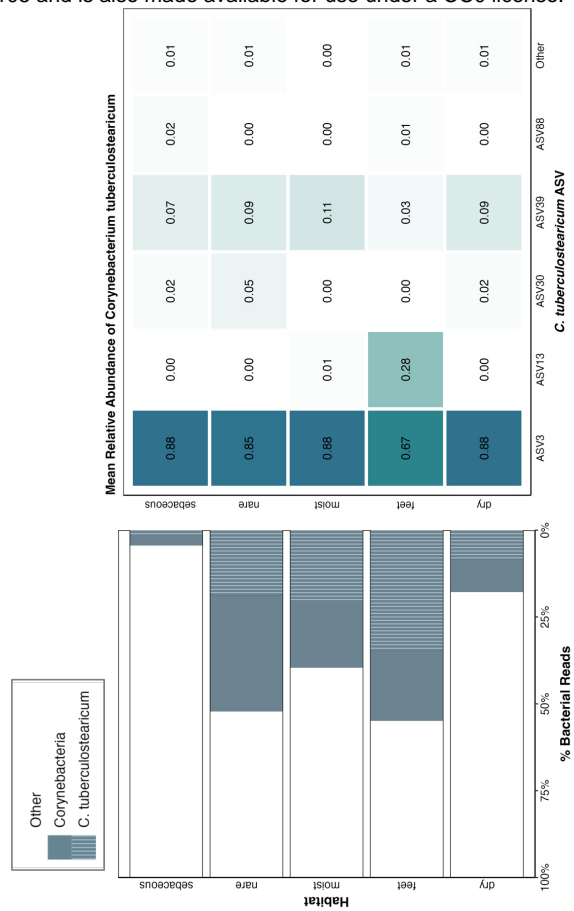


Figure 2

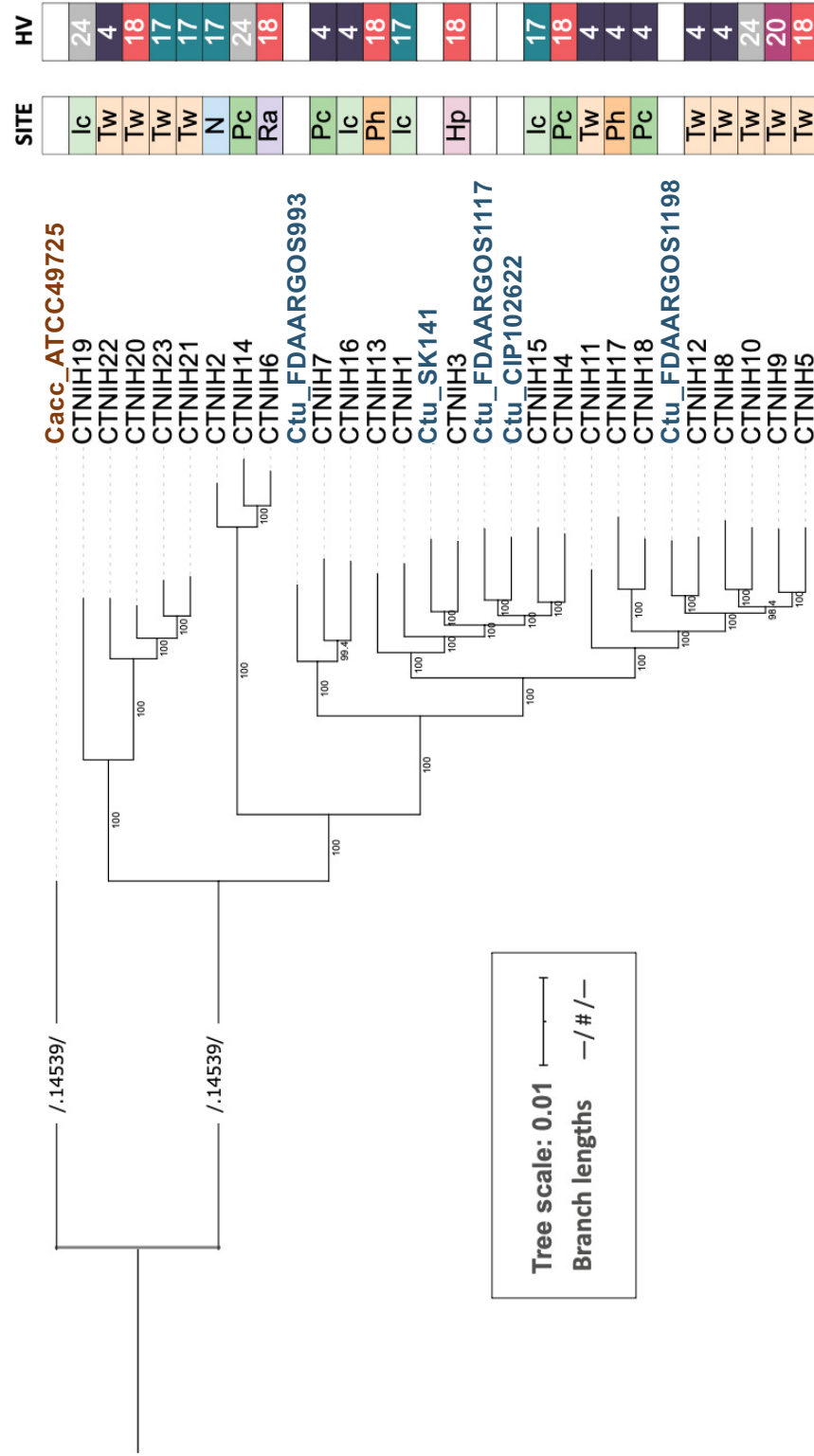
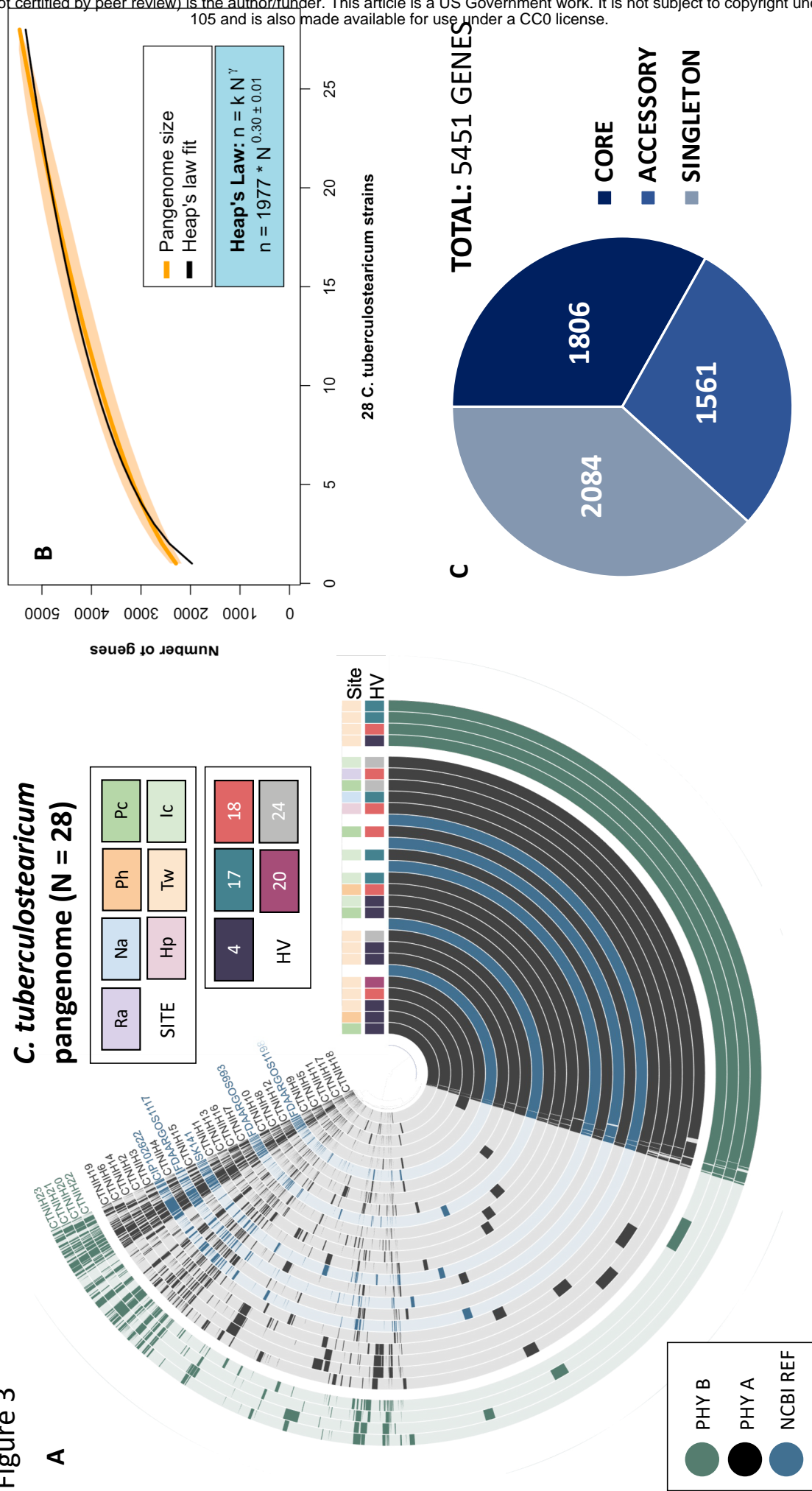


Figure 3



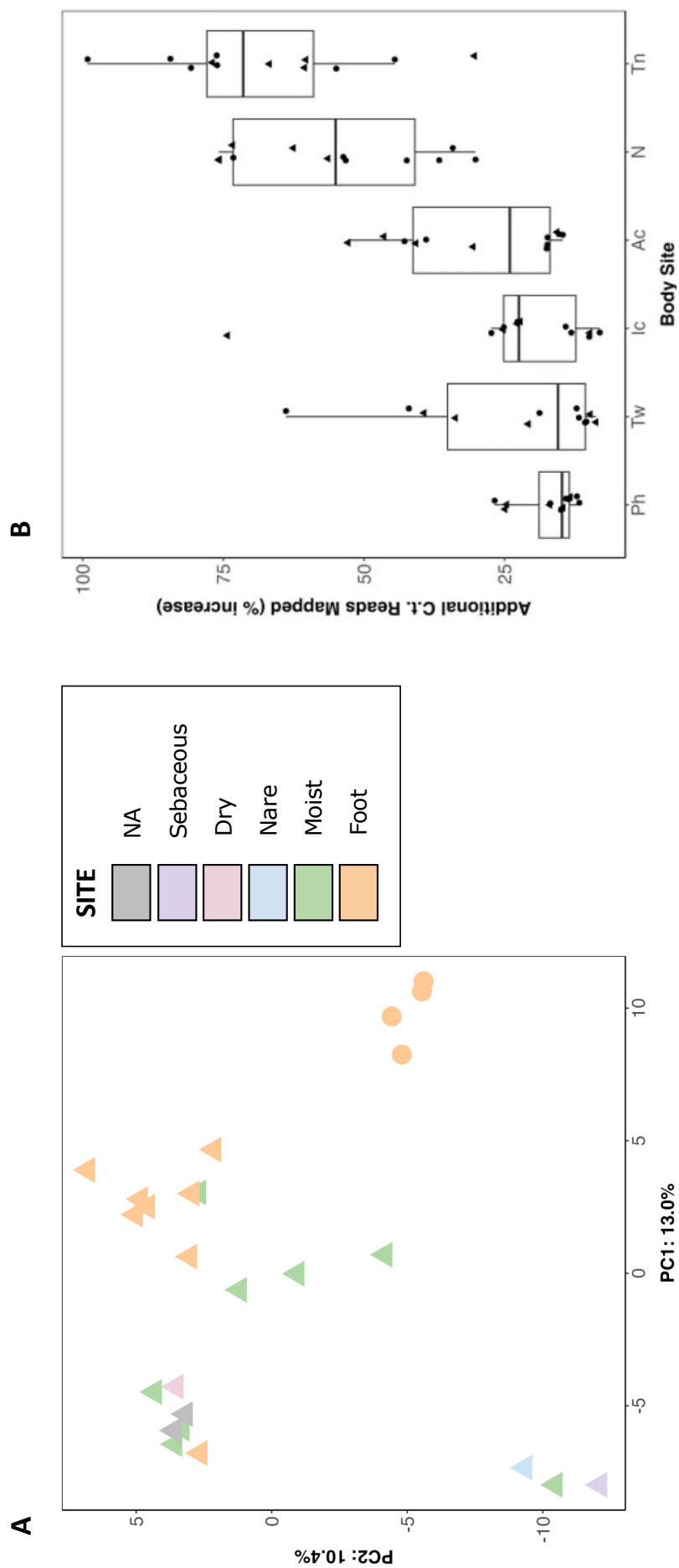


Figure 4

Figure 5

