# A genome catalog of the early-life human skin microbiome

Zeyang Shen[1], Lukian Robert[1], Milan Stolpman[2], You Che[2], Audrey Walsh[3,4,5], Richard Saffery[3,5], Katrina J. Allen[3,5], Jana Eckert[3], Angela Young[3], Clay Deming[1], Qiong Chen[1], Sean Conlan[1], Karen Laky[6], Jenny Min Li[6], Lindsay Chatman[6], Sara Saheb Kashaf[1], NISC Comparative Sequencing Program, VITALITY team[3], Heidi H. Kong[2], Pamela A. Frischmeyer-Guerrerio[6]*, Kirsten P. Perrett[3,4,5,7]*, Julia A. Segre[1]*†

[1]Microbial Genomics Section, Translational and Functional Genomics Branch, National Human Genome Research Institute, NIH, Bethesda, Maryland, USA
[2]Dermatology Branch, National Institute of Arthritis and Musculoskeletal and Skin Diseases, NIH, Bethesda, Maryland, USA
[3]Murdoch Children's Research Institute, Parkville, Victoria, Australia
[4]Department of Paediatrics, University of Melbourne, Parkville, Victoria, Australia
[5]Centre for Food and Allergy Research, Murdoch Children's Research Institute, Parkville, Victoria, Australia
[6]Laboratory of Allergic Diseases, National Institute of Allergy and Infectious Diseases, Bethesda, Maryland, USA
[7]Department of Allergy & Immunology, Royal Children's Hospital, Parkville, Victoria, Australia

*These authors contributed equally.
†Correspondence: jsegre@nhgri.nih.gov

## Abstract

Metagenome-assembled genomes have greatly expanded the reference genomes for skin microbiome. However, the current reference genomes are largely based on samples from adults in North America and lack representation from infants and individuals from other continents. Here we used ultra-deep shotgun metagenomic sequencing to profile the skin microbiota of 215 infants at age 2-3 months and 12 months who were part of the VITALITY trial in Australia as well as 67 maternally-matched samples. Based on the infant samples, we present the Early-Life Skin Genomes (ELSG) catalog, comprising 9,194 bacterial genomes from 1,029 species, 206 fungal genomes from 13 species, and 39 eukaryotic viral sequences. This genome catalog substantially expands the diversity of species previously known to comprise human skin microbiome and improves the classification rate of sequenced data by 25%. The protein catalog derived from these genomes provides insights into the functional elements such as defense mechanisms that distinguish early-life skin microbiome. We also found evidence for vertical transmission at the microbial community, individual skin bacterial species and strain levels between mothers and infants. Overall, the ELSG catalog uncovers the skin microbiome of a previously underrepresented age group and population and provides a comprehensive view of human skin microbiome diversity, function, and transmission in early life.

## Background

In direct contact with the environment, human skin is both a barrier and a habitat for microbes, including bacteria, fungi, and viruses, which help modulate immune responses and provide colonization resistance from adverse species[1,2]. Skin microbial community composition is shaped both by the ecology of the body site (oily, moist, dry) and skin physiology[1]. For example, during the transition through puberty, the maturation of sebaceous glands creates a lipid-rich environment to facilitate growth of *Cutibacterium*[3]. Compared to adults, early-life skin is characterized by higher water content, lower natural moisturizing factor concentration, and fewer lipids[4,5], which provides a distinct cutaneous environment for microbes and a unique habitat to study the skin microbiome.

Human skin microbiota is initially seeded at birth largely from maternal microbiome in association with the mode of delivery[6–8]. This relationship fades within 4-6 weeks[6,7], but skin microbial communities at the species level were found to be similar between babies and mothers over weeks to years after delivery[6,9,10]. Even though multiple studies have investigated vertical transmission and development of the human gut microbiome[11–14], mother-to-infant transmission of skin microbiome remains underexplored. Specifically, resolution of vertical transmission of strains on the skin has never been directly demonstrated.

One major challenge in studying the early-life skin microbiome is the lack of microbial reference genomes. Previous skin metagenomic studies found approximately 50% of the metagenomic reads do not match genomes in public databases[1,15]. Recent advancement in metagenome-assembled genomes (MAGs) has made it possible to generate large genome collections beyond culture-dependent methods[16]. We have recently published the Skin Microbial Genome Collection (SMGC)[17], which greatly expanded the reference genomes for skin microbiome in adults and substantially improved the classification rate of metagenomic reads. Comprehensive genome collections are also available for human gut microbiome[18–21]. In particular, the recent Early-Life Gut Genomes (ELGG) catalog has indicated great diversity and novelty of early-life gut microbiome compared to later in life[19]. To date, there have been no reports of skin microbial genomes in the first year of life. Comparative research investigating the gut microbiome in different populations also demonstrated great diversity of microbiome in people living in different geographic locations[18,20,21]. However, the current skin microbial genomes are derived from mostly adults residing in North America[17] and lack representation of individuals from other continents.

Here, we sequenced and assembled metagenomes from over 500 skin swabs collected longitudinally at age 2-3 months and 12 months from two body sites of 215 infants born in Australia, providing a catalog of 9,439 nonredundant genomes across multiple kingdoms for early-life skin microbiome. Using these data, we characterized the taxonomic and functional profile of the early-life skin microbiome and investigated the vertical transmission of the skin microbiome between mothers and infants.

## Results

### Deep sequencing of early-life skin metagenomes resulted in 9,439 nonredundant microbial genomes

To obtain comprehensive skin microbiome on early life, we conducted ultra-deep shotgun metagenomic sequencing on 565 skin swabs collected from the cheek and antecubital fossa (inside bend of the elbow) of 215 infants who were part of the VITALITY trial[22] (Fig. 1a, Supplementary table 1, 2). Among these infants, 69 were sampled longitudinally at 2-3 months and 12 months, 3 were sampled at 2-3 months only, and the rest were sampled at 12 months only. The two skin sites were selected as being representative of sebaceous and moist sites, which are usually inhabited by distinct microbiomes[1] and have clinical importance for future eczema studies as these are commonly affected sites[23]. Each sample yielded a median of 28.6 million non-human reads (IQR = 11.7-48.6 million). We applied a previously established bioinformatic pipeline[24] to build MAGs from single samples. To increase MAG quality and the detection of rare species[17], we pooled reads from the two skin sites within the same individual at each time point to generate MAGs from an additional 283 co-assemblies (Fig. 1a). To generate MAGs, single and pooled sample were assembled with MEGAHIT[25] and binned with a combination of MetaBAT 2[26], MaxBin 2[27], and CONCOCT[28]. Prokaryotic MAGs were refined with metaWRAP[29] and checked for chimerism with GUNC[30], while eukaryotic MAGs were checked for quality with EukCC[31]. Eukaryotic viral sequences were detected by aligning the contigs from MEGAHIT to the nucleotide collection database (nt) with BLASTn[32] and checked for quality with CheckV[33]. After removing redundant genomes across the entire dataset, our analyses yielded 9,194 nonredundant prokaryotic MAGs, 206 nonredundant eukaryotic MAGs and 39 eukaryotic viral sequences, comprising the Early-life Skin Genome (ELSG) catalog.

Among the 9,194 nonredundant prokaryotic MAGs (Fig. 1b, Fig. S1a, Supplementary table 3), 1,550 were classified as "high-quality" (completeness >90%, contamination <5%, and the presence of 5S, 16S and 23S rRNA genes and at least 18 of the standard tRNAs); 2,880 as "near-complete" (completeness >90%, contamination <5%, and didn't meet the rRNA or tRNA requirement of high-quality MAGs); and 4,764 as "medium-quality" (completeness >50%, contamination <10%, and quality score defined as completeness-5×contamination[18] > 50) based on the Metagenome-Assembled Genome standard[34]. As a complement to the standard quality metrics, we estimated the level of strain heterogeneity of each MAG using CMSeq[16] and obtained the median at 0.16% for prokaryotic MAGs. We applied similar criteria to 206 nonredundant eukaryotic MAGs, resulting in 5 "high-quality" (completeness >90%, contamination <5%, and the presence of 5S, 18S, 26S rRNA genes as well as at least 18 of the standard tRNAs), 42 "near-complete" (completeness >90%, contamination <5%, and didn't meet the rRNA or tRNA requirement of high-quality MAGs), and 159 "medium-quality" MAGs (completeness >50%, contamination <10%) (Fig. 1b, Fig. S1a, Supplementary table 4). Higher quality of MAGs was usually associated with a lower number of contigs, a larger N50, a lower level of strain heterogeneity, a higher read depth, and the presence of more unique tRNAs (Fig. S1a). Among the 39 eukaryotic viral sequences in the ELSG catalog, 9 were classified as "complete" (completeness=100%), 20 as "high-quality" (completeness >90%), 8 as "medium-quality" (completeness >50%), and only 2 as "low-quality" (completeness <50%) according to

127    CheckV[33] (Fig. 1c, Fig. S1b, Supplementary table 5). Considering the challenge of assembling
128    complete viral sequences from short-read metagenomes[33], we decided to include the two low-
129    quality sequences in the ELSG catalog.
130
131    To investigate transmission, we collected 67 skin swabs from the antecubital fossa of mothers
132    during the 12-month infant visit (Fig. 1a). These samples underwent DNA sequencing and were
133    assembled into individual sample-level MAGs using the aforementioned bioinformatic pipeline.
134    The mother samples yielded a total of 721 bacterial MAGs, 55 fungal MAGs, and 3 eukaryotic
135    viral sequences of medium quality or higher.
136

137    ### Species diversity in the ELSG catalog
138    To characterize the phylogenetic diversity of the ELSG catalog, we used 95% average nucleotide
139    identity (ANI) threshold to further cluster the nonredundant MAGs into 1,029 prokaryotic and
140    13 eukaryotic species-level clusters[35]. We assigned species-level taxonomy to representative
141    prokaryotic MAGs with GDTB-Tk[36] and to eukaryotic MAGs with >95% ANI to GenBank fungal
142    genomes. Rarefaction analysis showed that the number of species in the ELSG was not
143    saturated, when including MAGs recovered from a single sample. Excluding species recovered
144    from only one sample, which may be transient in nature or individual-specific, the number of
145    species came close to saturation, indicating that the ELSG catalog captured most of the
146    common species present on the early-life skin (Fig. 2a).
147

148    Next, we explored the novelty of the species diversity in the ELSG catalog. We compared the
149    ELSG catalog with the Skin Microbial Genome Collection (SMGC)[17], a collection of cultured and
150    uncultured skin microbial genomes primarily based on adult samples in North America, and the
151    Early-Life Gut Genome (ELGG) catalog[19]. Among the 1,029 representative prokaryotic MAGs in
152    the ELSG catalog, 699 clustered independently of any genome from the SMGC and the ELGG,
153    expanding the phylogenetic diversity by 56% (Fig. 2b, Fig. S2a, b, Supplementary table 3).
154    Among these, 313 were not assigned with species-level taxonomy based on GTDB (red in Fig.
155    2b, Fig. S2b). Note that 79 (11%) species-level clusters overlapped with MAGs built from
156    mothers skin samples (blue in Fig. 2b, Fig. S2b), suggesting these species are likely population-
157    specific rather than early-life-specific. ELSG-specific species spanned 16 different phyla greatly
158    expanding the current knowledge of skin microbiome. Top genera of the early-life-specific
159    species were *Streptococcus*, *Corynebacterium*, *Neisseria*, *Bifidobacterium*, and *Prevotella* (Fig.
160    2c). Early-life species-level clusters that were also present in the SMGC-specific species were
161    from the genera *Streptococcus*, *Corynebacterium*, and *Prevotella*. As some of the best studied
162    skin genera, *Staphylococcus* harbored very few ELSG-specific species, and similarly
163    *Cutibacterium* species were almost always found in both the ELSG and the SMGC.
164

165    Among the eukaryotic genera covered by the ELSG catalog, *Malassezia* was unsurprisingly the
166    dominant genus, followed by *Saccharomyces* (Supplementary table 4), which was not present in
167    the SMGC or mother samples. To compare the *Malassezia* species in the ELSG and the SMGC,
168    we clustered 7 species-level representative MAGs from the ELSG classified to be *Malassezia*, 7
169    *Malassezia* MAGs from the SMGC, and representative GenBank reference genomes (Fig. 2d).
170    *Malassezia obtusa* was assembled from the early-life skin but not found in the SMGC, whereas

171    *M. rara*, which is a novel species found on the human skin in the SMGC, were not found in the
172    ELSG catalog. Together these findings suggest fungal specificity of early-life skin.

173

174    Next, we explored the species diversity of 39 eukaryotic viral sequences in the ELSG catalog.
175    The most prevalent viruses found on infant skin were torque teno virus and
176    gammapapillomavirus (Fig. 2e, Supplementary table 5). Interestingly, the majority of these viral
177    sequences were found exclusively in 12-month infants, except for the gammapapillomavirus
178    discovered on the cheeks of three infants at 2-3 months.

179

180    Considering the novel species discovered on early-life skin, we used the ELSG catalog as an
181    additional source of reference genomes to classify shotgun metagenomic reads. By adding the
182    ELSG to a Kraken 2 database[37] created from the default RefSeq genomes and the SMGC, we
183    obtained a median classification rate of 77% (IQR = 69%-82%) for the early-life skin
184    metagenomic datasets, which was a median of 25% improvement over the standard RefSeq
185    database (Fig. 2f, Fig. S2c). Interestingly, the ELSG also substantially improved the classification
186    rate for metagenomic data of mothers (Fig. 2f, Fig. S2c) and slightly improved read mapping for
187    the antecubital metagenomes of the SMGC (Fig. S2d), suggesting the value of the ELSG in
188    capturing age- or population-specific species.

189

190    Comparison of taxonomic profiles between early-life and adult skin microbiome
191    We next explored similarities of the infant skin microbial community at two time points as well
192    as the relatedness of infant skin to mothers. The microbial community of infants demonstrated
193    strong skin-site differentiation with cheek and antecubital samples separated on a principal
194    coordinate analysis as well as age differentiation with 2-3 months and 12 months separated for
195    each skin site (Fig. 3a). Interestingly, the microbial community on the antecubital fossa of
196    mothers was most similar to the antecubital fossa of infants at 12 months (Fig. 3a), suggesting a
197    potential trajectory of maturation in the microbial community from early life to adulthood. We
198    calculated Bray-Curtis dissimilarity between the antecubital fossa of babies and mothers and
199    saw a significantly lower beta diversity (p < 1e-4) between related infant-mother pairs
200    compared to unrelated infant-mother pairs, consistent for both infant sexes (Fig. 3b). We also
201    calculated the beta diversity between the two time points of the same infant as compared to
202    different individuals. For both body sites, we saw a significantly lower beta diversity (p < 0.01)
203    within the same individuals, indicating an individualized trajectory of maturation that starts as
204    early as 2-3 months (Fig. S3a). Together, this suggests that the microbial communities on infant
205    skin may be influenced by individual factors, including the mother's skin microbiome.

206

207    Overall, the skin microbiome of early life contained roughly 97.8% bacteria, 2% fungi, and 0.2%
208    viruses (Fig. 3c), or 92% bacteria, 1% fungi, and 7% viruses after genome size normalization (Fig.
209    S3b). Antecubital fossa of infants generally had a more diverse microbial community than the
210    cheek (Fig. S3c). We also saw an increase in diversity from 2-3 months to 12 months at both
211    body sites (Fig. S3c). At the phylum level, Actinobacteria were more abundant on antecubital
212    fossa, whereas more Firmicutes, particularly *Streptococcus*, was found on cheek (Fig. 3c). Both
213    Bacteroidetes and Proteobacteria gained abundances over time (Fig. 3c). Differential
214    abundance analysis indicated 222 genera significantly (adjusted p < 0.01) gained abundance at

215  antecubital fossa over time and 257 genera increased on the cheek, including *Neisseria* and
216  *Saccharomyces* (Fig. 3d). Another 62 genera and 43 genera lost abundance at 12 months on
217  antecubital fossa and cheek, respectively, including *Staphylococcus* (Fig. 3d), which is consistent
218  with previous studies that also found a decrease in *Staphylococcus* over time[7,38]. The
219  prevalence of abundant species was correlated with the number of genomes in the ELSG (Fig.
220  S3d). For instance, *Cutibacterium acnes* was the most prevalent species found on early-life skin
221  and contributed the largest number of MAGs in the ELSG (Fig. 3e). Consistent with a higher
222  abundance of *Staphylococcus* at 2-3 months, most of the *Staphylococcus* genomes were
223  assembled from infants at 2-3 months even though the sample size at 2-3 months is much
224  smaller than 12 months (Fig. 3e).
225  
226  ## Comparison of the early-life and adult skin microbiome protein catalogs
227  To estimate the functional capacity in the ELSG catalog, we predicted protein-coding sequences
228  for each of the 9,194 nonredundant bacterial MAGs, resulting in a total of ~3.5 million protein
229  clusters at 90% amino acid identity. According to the rarefaction analysis, the protein clusters
230  found in the ELSG catalog were not saturated, but close to saturation when only considering ~2
231  million protein clusters that were identified in at least two MAGs (Fig. S4a), consistent with
232  previous findings in gut microbiome[18,19]. When examining individual species, we discovered
233  that some of the most prominently represented species had either reached a saturation point
234  or were nearing saturation (Fig. 4a). The conspecific gene frequency had a bimodal distribution
235  (Fig. S4b), consistent with observations in the SMGC[17]. We defined those genes shared by at
236  least 90% conspecific genomes of each species as core genes and the rest as accessory genes[18]
237  (Fig. S4c) and then compared the functions encoded in the core and accessory genes based on
238  several annotation databases. Core genes were generally better annotated than accessory
239  genes in all databases (Fig. S4d). According to COG annotations[39], core genes were enriched for
240  functions related to metabolism and translation, whereas accessory genes were enriched for
241  functions related to replication, defense mechanisms, and transcription (Fig. 4b). A similar
242  pattern of functional roles performed by core and accessory genes has previously been
243  reported for gut microbiomes[18].
244  
245  We next compared the pan-genome of early-life skin microbiome with that of SMGC. The pan-
246  genome size was variable between the two genome collections for several species (Fig. S4e).
247  For example, *Micrococcus luteus* had a 14% larger pan-genome in the ELSG catalog, while, in
248  contrast, *Cutibacterium acnes* had a 5% larger pan-genome in the SMGC. Besides the pan-
249  genome size difference, many genes were specific to one collection (Fig. 4c). Interestingly,
250  ELSG- or SMGC-specific genes were enriched in COG categories such as cell motility and defense
251  mechanisms while collection-shared genes were enriched for functions related to metabolism
252  (Fig. 4d).
253  
254  ## Intraspecies genomic diversity between infants and mothers indicates vertical transmission
255  To characterize the genomic diversity across species-level clusters within the ELSG catalog, we
256  calculated the rate of intraspecies single-nucleotide variants (SNVs). *Rothia mucilaginosa*, a
257  prevalent species on early-life skin, contained the highest SNV density, 40 SNVs per kb,

258  suggesting a great potential of functional variability (Fig. 5a). By contrast, *Cutibacterium acnes*,
259  which was even more prevalent, had a much lower density of only about 5 SNVs per kb.
260  Similarly, *Staphylococcus epidermidis*, another common species found on skin, had about 5
261  SNVs per kb.
262
263  Next, we compared paired microbial genomes from infants and mothers. For all seven species
264  for which we had MAGs from at least four related infant-mother pairs, there were significantly
265  fewer SNVs genome-wide (p < 0.01) between related infant-mother pairs as compared to
266  unrelated infants and mothers, consistent with vertical transmission of skin microbes between
267  mother and infant (Fig. 5b). By looking at SNVs at protein-coding regions, four of the seven
268  species including *Cutibacterium acnes* had 62% or less SNVs shared by infants and mothers,
269  whereas the other three species including *Rothia mucilaginosa* had over 78% of SNVs shared by
270  infants and mothers (Fig. 5c). The small proportion of age-group-specific SNVs within these
271  three species was also consistent with the strikingly large differences between related and
272  unrelated babies and mothers (Fig. 5b).
273
274  Besides the genome sharing between infants and mothers, we also investigated the genome
275  sharing at different ages of infants. For the five species with at least four infants that yielded
276  longitudinal pairs of MAGs, the number of SNVs was generally lower within individuals than
277  across individuals (Fig. S5), suggesting temporally persistent microbial genomes on the host.
278  Due to a limited number of samples, further research is needed to examine the applicability of
279  such observation to a broad spectrum of species.
280
281  To further validate the mother-infant transmission of skin microbiome, we cultured
282  *Cutibacterium acnes* from the nasal swabs collected from six pairs of infants and mothers when
283  infants were 12-month-old. Depending on the variable viability of bacteria, we were able to
284  obtain and sequence 4-12 *C. acnes* independent colonies from each individual (Supplementary
285  table 6). Genomes from the related infants and mothers were often closely placed on a
286  phylogenetic tree (Fig. 5d). Consistent with that, we performed multi-locus sequence typing to
287  these genomes and found that four out of six mother-infant pairs shared at least one sequence
288  type (Fig. S5b), which is statistically significant (p = 0.012) based on a permutation test (Fig.
289  S5c). Together, this indicates the mother-infant transmission of skin microbiome is at the strain
290  level.
291

## Discussion

293  We present the first genome collection for early-life skin microbiome and the largest skin
294  microbial genome collection to date containing over a thousand species-level clusters of
295  bacterial and fungal genomes and an additional set of eukaryotic viral sequences. To our
296  knowledge, the ELSG catalog is also the first skin microbial genome collection based on samples
297  from Australia. It is an effective resource of genomes to improve the classification of
298  metagenomic reads for early-life samples and geographically distinct studies. The slightly
299  improved classification of North American samples by including the ELSG catalog could be due
300  to the ultra-deep sequencing and the large sample basis of this study, which recovered ultra-

301    rare and low-abundant species present on human skin across continents. Augmented read
302    mapping would be consistent with species that are more highly abundant in infants and at
303    lower abundance in adults. The ELSG catalog includes hundreds of species previously not
304    characterized for skin, many of which are novel species. Considering that skin is still an
305    understudied organ source of microbiome, this study has demonstrated the importance of
306    profiling different age groups and populations to capture a complete catalog of human skin
307    microbiome. Since the ELSG catalog was based on infant samples at age 12 months or less, this
308    resource will be of particular use in studies of childhood cutaneous disorders, such as atopic
309    dermatitis, which commonly begins in infancy.
310
311    Our study on the vertical transmission of the skin microbiome was empowered by a substantial
312    number of paired samples collected from infants and mothers. Evidence of vertical transmission
313    was found at the microbial community, individual species, and strain levels. Specifically, infants
314    and their mothers had closely related microbial profiles, relatively similar conspecific MAGs for
315    7 species, and shared strains of *Cutibacterium acnes*. Likewise, based on the longitudinal
316    samples of infants at 2-3 months and 12 months, we saw evidence of temporal persistence of
317    the microbiome on infant skin in both microbial profiles and genomes. These findings indicate
318    an important role of mothers in shaping the skin microbiome of early life and suggest
319    microbiome at later time point could be affected by what was preceded. Notably, two out of
320    the six mother-infant pairs where we cultured *C. acnes* isolates shared none of their *C. acnes*
321    strains. It suggests that mothers may not be the only source of skin microbiome for infants, and
322    the skin microbial transmission from other sources such as fathers requires further exploration.
323    Further study is also needed to extend our findings to other species that were not investigated
324    in this study.
325
326    Based on the ELSG catalog, we analyzed the largest published protein catalog for skin
327    microbiome to estimate the functional capacity. By looking at the conspecific pan-genomes, we
328    summarized the functional categories that distinguish core and accessory genes, which
329    replicated the findings in gut microbiome. Interestingly, genes found only in one of the two
330    current skin genome collections were consistently represented by functions related to defense
331    mechanism and replication, recombination, and repair. These categories are potentially the
332    drivers of functional specificity in early-life skin microbiome. Further experiments are needed to
333    validate the function and importance of individual genes in maintaining homeostasis on early-
334    life skin.
335

## Conclusions
337    In summary, our investigation involved profiling the skin metagenomes of infants who had been
338    previously under-represented. This pioneering effort led to the development of the ELSG
339    catalog, which significantly expands the repertoire of skin microbial genomes in infants. The
340    ELSG catalog presents a comprehensive and versatile resource for future studies focused on
341    various aspects of the infant skin microbiome such as microbial transmission and development,
342    and the intricate interplay between disease and the early-life skin microbiota.
343

## Methods

### Participant recruitment, skin sampling and metagenomic sequencing

New mothers along with their infants were recruited as part of the VITALITY trial[22]. Written informed consent was obtained for all participants in this study. Skin samples were collected from the antecubital fossa and cheek of 72 infants at ages 2-3 months. Sixty-nine of these infants together with 140 additional infants were sampled at the same sites at age 12 months. In addition, 67 of these infants' mothers were sampled at the antecubital fossa during the same visit when the 12-month samples were taken. To maximize microbial recovery, no bathing was permitted within 24 h of sample collection. Skin was sampled with an established protocol using pre-moistened Puritan foam swabs collected and stored in 100 µL Yeast Cell Lysis Buffer (Lucigen) buffer at -80° and shipped on dry ice. Concomitant with skin sample collection, air swabs were collected as negative controls to account for any potential environmental or reagent contaminants.

Samples were converted to genomic DNA with an established protocol[40,41]. Briefly, DNA libraries for Illumina sequencing were prepared using the Nextera XT DNA Library Preparation Kit (Illumina) per manufacturer's instructions with the exception of increasing the AMPure XP Bead clean-up volume from 30 µL to 50 µL. 1 ng of extracted DNA was used as input into the fragmentation step. DNA is simultaneously fragmented and tagged with sequencing adapters in a single-tube enzymatic reaction. Libraries were then sequenced with the Illumina NovaSeq 6000 sequencing platform at the NIH Intramural Sequencing Center for 2 × 150 bp, 50 million paired-end reads per sample.

Most of the negative controls yielded <1% of the reads derived from skin samples except for one. We excluded the skin samples collected at the same time of that air swab together with one infant's antecubital fossa sample which yielded less than 10,000 reads. Our final set of samples for analysis includes 565 from infants (424 at 12 months (212 infants x 2 skin sites) + 144 at 2-3 months (72 infants x 2 skin sites) - 3 samples failed) and 67 from mothers.

### Bacteria culturing and sequence typing

Nasal culture samples were obtained from infants and mothers during the same visit when infants were 12-month-old using the COPAN eSwab system in 1 mL AMIES and frozen at -80°C. Broths were diluted and plated on Brain Heart Infusion Agar (BHI + 10 µg/mL Fosfomycin) and incubated in an anaerobic chamber for 7 days at 37°C. Colonies were screened with PCR using C. acnes-specific primers PA-1 5'-GGGTTGTAAACCGCTTTCGCTG-3 and PA-2 5'-GGCACACCCATCTCTGAGCAC-3, then streaked for purity on Blood Agar plates (TSA with 5% Sheep Blood – Remel R01201). gDNA was prepared from isolates and sequenced with an established protocol[17]. C. acnes genomes were assembled from sequenced reads using SPAdes[42] and checked for quality using the 'lineage_wf' workflow of CheckM v1.1.3[43]. The sequence type of each C. acnes genome was identified by multi-locus sequence typing scheme from PubMLST[44]. C. acnes genomes of the same individual were first dereplicated at 99.9% ANI with dRep v3.2.2[45] and then used to build the phylogenetic tree with GToTree v1.6.37[46] based on the single-copy gene set of Actinobacteria.

387

### Pre-processing, metagenomic assembly, and contig binning

Metagenomic reads were trimmed for adapters with Cutadapt v3.4 using the parameters "--nextseq-trim 20 -e 0.15 -m 50"[47] and checked for quality with PRINSEQ-lite v0.20.4 using the parameters "-lc_method entropy -lc_threshold 70 -min_len 50 -min_qual_mean 20 -ns_max_n 5 -min_gc 10 -max_gc 90"[48]. Reads with less than 50 bp length after trimming were removed. The reads were then aligned to the GRCh38 human reference genome with Bowtie2 v2.4.5 using the parameters "--very-sensitive"[49]. The human reads were removed before assembly.

Metagenomic assembly was performed with MEGAHIT v1.2.9 using the default parameters[50]. Pool individual runs were conducted after concatenating the reads from the two skin sites of the same infant at each time point. We performed 283 co-assemblies including 211 from 12 months and 72 from 2-3 months. Contigs were then binned with a combination of MetaBAT 2 v2.15[26], MaxBin 2 v2.2.7[27], and CONCOCT v1.1.0[28] by running the binning module of metaWRAP v1.3.2[29] with the parameter '-l 1500' indicating the minimal contig length 1500 bp.

### Genome quality assessment

To obtain prokaryotic MAGs, the bins produced by each binning tool were refined with the Bin_refinement module of metaWRAP v1.3.2[29] using parameters "-c 50 -x 10" enforcing >50% completeness and <10% contamination. The completeness and contamination of refined bins were evaluated with the 'lineage_wf' workflow of CheckM v1.1.3[43]. The quality score was calculated as: completeness − 5 × contamination. Ribosomal RNAs in each genome were detected with the 'cmsearch' function of INFERNAL v1.1.4 using parameters "--anytrunc --noali"[51] against the Rfam covariance models for the 5S (5S_rRNA), 16S (SSU_rRNA _bacteria) and 23S rRNAs (LSU_rRNA _bacteria)[52]. Transfer RNAs of the standard 20 amino acids were identified with tRNAScan-SE v2.0.11 using the parameter '-B' for bacterial species[53]. Each genome was assessed for chimerism with GUNC v1.0.5[54]. The MAGs with contamination greater than 0.05, clade separation greater than 0.45 and a reference representation score greater than 0.5 were excluded. Based on the Metagenome-Assembled Genome standard[34], MAGs with >90% completeness, <5% contamination, the presence of 5S, 16S and 23S rRNA genes, and at least 18 tRNAs were reported as high-quality draft genomes. MAGs with >90% completeness and <5% contamination but missing the rRNAs or tRNAs were reported as near-complete genomes. MAGs with >50% completeness and <10% contamination were reported as medium quality.

To assess eukaryotic MAGs, the bins from the three binning tools were estimated for completeness and contamination with EukCC v2.1.0[31]. rRNAs and tRNAs were identified using the same approach above except that the Rfam[52] covariance models 5_S_rRNA, SSU_rRNA_eukarya and LSU_rRNA_eukarya were used to find 5S, 18S and 26S, respectively. Bins with >90% completeness, <5% contamination, the presence of 5S, 18S and 26S rRNA genes, and at least 18 tRNAs were reported as high-quality draft genomes. Those with >90% completeness and <5% contamination but not satisfying the rRNAs and tRNAs requirements

429 were defined as near-complete. The rest of the bins with >50% completeness and <10%
430 contamination were reported as medium-quality genomes.
431
432 We further mapped each contig of MAGs to the nt database with BLASTn v2.8.0[32] to assess viral
433 contamination. Contigs with the top hit of a eukaryotic viral genome with >95% nucleotide
434 identity, >1000 bp aligned sequence, and >70% total contig aligned were removed. The contig
435 number and N50 of MAGs were calculated using in-house scripts. Read depth was calculated by
436 first mapping the raw reads back to MAGs Bowtie2 v2.4.5[49] using the default parameters and
437 then calculating mean depth with SAMtools v1.16.1[55]. The strain heterogeneity was estimated
438 by the "polymut.py" script of CMSeq v1.0.4 with parameters "--mincov 10 --minqual 30 --
439 dominant_frq_thrsh 0.8"[16].
440
441 Eukaryotic viral sequences were detected by aligning the contigs from MEGAHIT to the nt
442 database with BLASTn v2.8.0[32] by requiring >90% nucleotide identity, >1000 bp aligned
443 sequence, and >70% total contig aligned. The quality of viral sequences was assessed with
444 CheckV v1.0.0 based on database v1.5[33].
445

446 Redundancy removal and species clustering
447 To remove redundant genomes that were recovered by both single and pooled sample runs, we
448 dereplicated MAGs at a 99.9% ANI threshold with dRep v3.2.2 using parameters '-pa 0.999 --
449 SkipSecondary -cm larger --S_algorithm fastANI -comp 50 -con 10'[45]. fastANI v1.33[35] was used
450 to accelerate the process. Dereplication was performed on prokaryotic MAGs and eukaryotic
451 MAGs separately.
452
453 The MAGs were clustered at the species level by dereplicating at a 95% ANI threshold with
454 dRep v3.2.2 using parameters '-pa 0.90 -sa 0.95 -nc 0.30 -cm larger --S_algorithm fastANI -comp
455 50 -con 10 --run_tertiary_clustering --clusterAlg single'[45]. Representative genome of each
456 species-level cluster was selected based on the dRep scores derived from genome
457 completeness, contamination, strain heterogeneity, and contig N50.
458

459 Taxonomic assignment and phylogenetic analysis
460 Taxonomic annotation of prokaryotic MAGs was assigned with the "classify_wf" workflow of
461 GTDB-Tk v2.1.0 using default parameters and GTDB database release 207[36,56]. The phylogenetic
462 tree of bacterial representative genomes of species-level clusters was built with IQ-TREE
463 v1.6.12 using the parameter "-m MFP"[57] based on the protein sequence alignments generated
464 by GTDB-Tk.
465

466 The eukaryotic MAGs were compared with all of the GenBank fungal genomes first using
467 Mash[58] and then assigned species-level taxonomy with at least 95% ANI calculated by fastANI
468 v1.33[35]. The phylogenetic tree was built with the script BUSCO_phylogenomics.py
469 (https://github.com/jamiemcg/BUSCO_phylogenomics) based on single-copy marker genes
470 identified by BUSCO v4.1.3 using the parameter "-m geno -f --auto-lineage-euk"[59]. The

471    phylogenetic trees were visualized with iTOL[60]. The taxonomic classifications of viral sequences
472    were assigned by the top alignment hit from BLASTn[32].
473

### Metagenomic read classification and microbial abundance estimation

475    Metagenomic reads were mapped with Kraken v2.1.2 using parameters "--confidence 0.1 --
476    paired"[37] against the standard RefSeq database (release 99) and two custom database with
477    additional representative genomes from the SMGC and ELSG catalogs. To integrate the genome
478    catalogs with the RefSeq genomes, we first converted GTDB taxonomy to NCBI taxonomy using
479    the "gtdb_to_ncbi_majority_vote.py" script available in the GTDB-Tk repository[36] and then
480    obtained NCBI taxonomy IDs corresponding to the species- and genus-level taxonomy of each
481    genome with taxonkit v0.12.0[61]. We excluded 22 and 106 representative MAGs from the SMGC
482    and the ELSG, respectively, which did not have a match ID at the genus level. For MAGs with a
483    match ID at the genus level but not at the species level, we created a new taxonomy ID
484    associated with each MAG when building the Kraken databases. Classification improvement
485    was calculated on a per-sample basis as (proportion of reads classified with custom
486    database – proportion of reads classified with RefSeq database)/proportion of reads classified
487    with RefSeq database × 100. Species-level microbial abundances were computed with Bracken
488    v2.5 using parameters "-r 100 -l S"[62].
489

### Alpha and beta diversity calculation

491    Skin metagenomic data with less than 800,000 classified reads were excluded (4% of samples).
492    The remaining samples were first rarefied and then calculated for the number of species with
493    ≥5 reads (richness) and Shannon index with the "diversity" function of vegan package in R v4.1.
494    To calculate the beta diversity, we first removed taxa present in ≤20% samples and then
495    performed log transformation on species abundances after adding pseudocount 1. Bray-Curtis
496    dissimilarity was calculated with the "distance" function of phyloseq v1.38.0[63] in R v4.1.
497    Principal coordinate analysis was conducted based on Bray-Curtis dissimilarity with the
498    "ordination" function of phyloseq package.
499

### Differential abundance analysis

501    Differential abundance was calculated with DESeq2 v1.34.0[64] using the parameters
502    'test="Wald", sfType="poscounts", fitType="local"' based on the rarefied raw reads as used for
503    diversity calculation. Low-prevalence taxa present in less than 10% of samples were removed.
504    Comparisons were conducted for each of the two skin sites, comparing 2-3 months and 12
505    months; and for antecubital fossa, comparing infants at 12 months and mothers. Significantly
506    differential taxa were identified by <0.01 adjusted p-value and >2-fold change.
507

### Pan-genome analysis and functional annotation

509    Protein-coding sequences (CDS) of each genome were predicted and annotated with Prokka
510    v1.14.6 using parameter "--kingdom Bacteria"[65]. Protein clustering across all species of
511    genomes was conducted with the 'easy-linclust' function of MMseqs2 v13-45111 using
512    parameters '--cov-mode 1 -c 0.8 --kmer-per-seq 80 --min-seq-id 0.9' to generate protein
513    clusters at 90% amino acid identity, respectively[66].

514
515 The pan-genome analysis was performed only on near-complete and high-quality genomes.
516 Species with at least ten near-complete or high-quality nonredundant genomes were analyzed
517 with Panaroo v1.3.0 using the parameters '--clean-mode strict --merge_paralogs -c 0.90 --
518 core_threshold 0.90 -f 0.5' for ≥90% amino acid identity and a family threshold of 50%[67].
519 Functional annotation of all protein sequences was performed with eggNOG-mapper v2.1.6[68] to
520 obtain COG[39], KEGG[69], Pfam[70] and GO[71] annotations.
521

522 SNV analysis
523 To assess SNV density of species, we first mapped conspecific genomes to the representative
524 genome using the 'nucmer' program of MUMmer v3.1[72], filtered alignments with the 'delta-
525 filter' program using parameters '-q -r', and then identified SNVs with the 'show-snps' program.
526 SNV density of each genome was computed by dividing the number of SNVs by the size of the
527 representative genome. Only SNVs which occurred in at least two conspecific genomes were
528 included in the analysis. The final SNV density of each species was the mean of SNV densities of
529 all conspecific genomes. The same programs and parameters were used for mother-infant
530 genome comparisons.
531

532 Statistical analysis
533 Statistical analyses were performed using ggpubr package in R v4.1 or scipy package in Python
534 v3.9.9. Two-sided Wilcoxon rank sum tests and t-tests were used to evaluate differences
535 between groups. Pearson correlation coefficient was used to assess correlation. Functional
536 enrichment analysis was performed using two-sided Fisher's exact test, with p-values adjusted
537 by the Bonferroni method. The permutation test (n = 1,000) was applied to assess the
538 significance of sequence type sharing between mothers and infants.
539

540 # Availability of data and materials
541 The raw metagenomic sequencing data are available in the NCBI BioProject database under
542 project number PRJNA971252. The MAGs of the ELSG catalog can be found at
543 https://research.nhgri.nih.gov/projects/ELSG/. Additionally, on the same website, users have
544 access to download nonredundant genomes, species-level representative genomes,
545 phylogenetic tree files, protein catalog, pan-genome annotations, and a custom Kraken 2
546 database based on the ELSG catalog. All the code utilized in this study is available on GitHub at
547 https://github.com/skinmicrobiome/ELSG.
548

549 Other publicly available data used in this project: SMGC is available at
550 http://ftp.ebi.ac.uk/pub/databases/metagenomics/genome_sets/skin_microbiome. Shotgun
551 metagenomic sequencing data used in the SMGC is accessed from the NCBI Sequence Read
552 Archive under accession number SRP002480. ELGG catalog is available at
553 https://doi.org/10.5281/zenodo.6969520.
554

## Acknowledgements

## References

1. Oh, J. *et al.* Biogeography and individuality shape function in the human skin metagenome. *Nature* **514**, 59–64 (2014).

2. Harris-Tryon, T. A. & Grice, E. A. Microbiota and maintenance of skin barrier function. *Science (1979)* **376**, 940–945 (2022).

3. Park, J. *et al.* Shifts in the Skin Bacterial and Fungal Communities of Healthy Children Transitioning through Puberty. *Journal of Investigative Dermatology* **142**, 212–219 (2022).

4. Casterline, B. W. & Paller, A. S. Early development of the skin microbiome: therapeutic opportunities. *Pediatr Res* **90**, 731–737 (2021).

5. Stamatas, G. N., Nikolovski, J., Mack, M. C. & Kollias, N. Infant skin physiology and development during the first years of life: a review of recent findings based on in vivo studies. *Int J Cosmet Sci* **33**, 17–24 (2011).

6. Chu, D. M. *et al.* Maturation of the infant microbiome community structure and function across multiple body sites and in relation to mode of delivery. *Nat Med* **23**, 314–326 (2017).

7. Capone, K. A., Dowd, S. E., Stamatas, G. N. & Nikolovski, J. Diversity of the human skin microbiome early in life. *Journal of Investigative Dermatology* **131**, 2026–2032 (2011).

8. Dominguez-Bello, M. G. *et al.* Delivery mode shapes the acquisition and structure of the initial microbiota across multiple body habitats in newborns. *Proceedings of the National Academy of Sciences* **107**, 11971–11975 (2010).

9. Bogaert, D. *et al.* Mother-to-infant microbiota transmission and infant microbiota development across multiple body sites. *Cell Host Microbe* **31**, 447-460.e6 (2023).

10. Zhu, T. *et al.* Age and Mothers: Potent Influences of Children's Skin Microbiota. *Journal of Investigative Dermatology* **139**, 2497-2505.e6 (2019).

11. Ferretti, P. *et al.* Mother-to-Infant Microbial Transmission from Different Body Sites Shapes the Developing Infant Gut Microbiome. *Cell Host Microbe* **24**, 133-145.e5 (2018).

12. Yassour, M. *et al.* Strain-Level Analysis of Mother-to-Child Bacterial Transmission during the First Few Months of Life. *Cell Host Microbe* **24**, 146-154.e4 (2018).

13. Valles-Colomer, M. *et al.* The person-to-person transmission landscape of the gut and oral microbiomes. *Nature* (2023) doi:10.1038/s41586-022-05620-1.

14. Valles-Colomer, M. *et al.* Variation and transmission of the human gut microbiota across multiple familial generations. *Nat Microbiol* **7**, 87–96 (2022).

15. Oh, J., Byrd, A. L., Park, M., Kong, H. H. & Segre, J. A. Temporal Stability of the Human Skin Microbiome. *Cell* **165**, 854–866 (2016).

16. Pasolli, E. *et al.* Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* **176**, 649-662.e20 (2019).

17. Saheb Kashaf, S. *et al.* Integrating cultivation and metagenomics for a multi-kingdom view of skin microbiome diversity and functions. *Nat Microbiol* **7**, 169–179 (2022).

18. Almeida, A. *et al.* A unified catalog of 204,938 reference genomes from the human gut microbiome. *Nat Biotechnol* **39**, 105–114 (2021).

19. Zeng, S. *et al.* A compendium of 32,277 metagenome-assembled genomes and over 80 million genes from the early-life human gut microbiome. *Nat Commun* **13**, (2022).

609   20.   Jin, H. *et al.* A high-quality genome compendium of the human gut microbiome of Inner
610          Mongolians. *Nat Microbiol* **8**, 150–161 (2023).
611   21.   Kim, C. Y. *et al.* Human reference gut microbiome catalog including newly assembled
612          genomes from under-represented Asian metagenomes. *Genome Med* **13**, 134 (2021).
613   22.   Allen, K. J. *et al.* VITALITY trial: protocol for a randomised controlled trial to establish the
614          role of postnatal vitamin D supplementation in infant immune health. *Open* **5**, 9377
615          (2015).
616   23.   Kennedy, E. A. *et al.* Skin microbiome before development of atopic dermatitis: Early
617          colonization with commensal staphylococci at 2 months is associated with a lower risk of
618          atopic dermatitis at 1 year. *Journal of Allergy and Clinical Immunology* **139**, 166–172
619          (2017).
620   24.   Saheb Kashaf, S., Almeida, A., Segre, J. A. & Finn, R. D. Recovering prokaryotic genomes
621          from host-associated, short-read shotgun metagenomic sequencing data. *Nature
622          Protocols* vol. 16 2520–2541 Preprint at https://doi.org/10.1038/s41596-021-00508-2
623          (2021).
624   25.   Li, D., Liu, C.-M., Luo, R., Sadakane, K. & Lam, T.-W. MEGAHIT: an ultra-fast single-node
625          solution for large and complex metagenomics assembly via succinct de Bruijn graph.
626          *Bioinformatics* **31**, 1674–1676 (2015).
627   26.   Kang, D. D. *et al.* MetaBAT 2: an adaptive binning algorithm for robust and efficient
628          genome reconstruction from metagenome assemblies. *PeerJ* **7**, e7359 (2019).
629   27.   Wu, Y.-W., Simmons, B. A. & Singer, S. W. MaxBin 2.0: an automated binning algorithm
630          to recover genomes from multiple metagenomic datasets. *Bioinformatics* **32**, 605–607
631          (2016).
632   28.   Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat
633          Methods* **11**, 1144–1146 (2014).
634   29.   Uritskiy, G. V, DiRuggiero, J. & Taylor, J. MetaWRAP—a flexible pipeline for genome-
635          resolved metagenomic data analysis. *Microbiome* **6**, 158 (2018).
636   30.   Orakov, A. *et al.* GUNC: detection of chimerism and contamination in prokaryotic
637          genomes. *Genome Biol* **22**, 178 (2021).
638   31.   Saary, P., Mitchell, A. L. & Finn, R. D. Estimating the quality of eukaryotic genomes
639          recovered from metagenomic analysis with EukCC. *Genome Biol* **21**, 244 (2020).
640   32.   Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment
641          search tool. *J Mol Biol* **215**, 403–410 (1990).
642   33.   Nayfach, S. *et al.* CheckV assesses the quality and completeness of metagenome-
643          assembled viral genomes. *Nat Biotechnol* **39**, 578–585 (2021).
644   34.   Bowers, R. M. *et al.* Minimum information about a single amplified genome (MISAG) and
645          a metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nat Biotechnol* **35**,
646          725–731 (2017).
647   35.   Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High
648          throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.
649          *Nat Commun* **9**, 5114 (2018).
650   36.   Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk v2: memory friendly
651          classification with the genome taxonomy database. *Bioinformatics* **38**, 5315–5316 (2022).

652   37.   Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with Kraken 2.
653         *Genome Biol* **20**, 257 (2019).
654   38.   Rapin, A. *et al.* The skin microbiome in the first year of life and its association with atopic
655         dermatitis. *Allergy* **n/a**, (2023).
656   39.   Galperin, M. Y., Makarova, K. S., Wolf, Y. I. & Koonin, E. V. Expanded microbial genome
657         coverage and improved protein family annotation in the COG database. *Nucleic Acids Res*
658         **43**, D261–D269 (2015).
659   40.   Tirosh, O. *et al.* Expanded skin virome in DOCK8-deficient patients. *Nat Med* **24**, 1815–
660         1821 (2018).
661   41.   Byrd, A. L. *et al.* Staphylococcus aureus and Staphylococcus epidermidis strain diversity
662         underlying pediatric atopic dermatitis. *Sci Transl Med* **9**, eaal4651 (2017).
663   42.   Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its Applications to
664         Single-Cell Sequencing. *Journal of Computational Biology* **19**, 455–477 (2012).
665   43.   Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
666         assessing the quality of microbial genomes recovered from isolates, single cells, and
667         metagenomes. *Genome Res* **25**, 1043–1055 (2015).
668   44.   Jolley, K. A., Bray, J. E. & Maiden, M. C. J. Open-access bacterial population genomics:
669         BIGSdb software, the PubMLST.org website and their applications [version 1; referees: 2
670         approved]. *Wellcome Open Res* **3**, (2018).
671   45.   Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate
672         genomic comparisons that enables improved genome recovery from metagenomes
673         through de-replication. *ISME J* **11**, 2864–2868 (2017).
674   46.   Lee, M. D. GToTree: a user-friendly workflow for phylogenomics. *Bioinformatics* **35**,
675         4162–4164 (2019).
676   47.   Martin, M. Cutadapt removes sequences from high-throughput sequencing reads.
677         *EMBnet Journal* **17**, 1 (2013).
678   48.   Schmieder, R. & Edwards, R. Quality control and preprocessing of metagenomic datasets.
679         *Bioinformatics* **27**, 863–864 (2011).
680   49.   Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**,
681         357–359 (2012).
682   50.   Li, D., Liu, C.-M., Luo, R., Sadakane, K. & Lam, T.-W. MEGAHIT: an ultra-fast single-node
683         solution for large and complex metagenomics assembly via succinct de Bruijn graph.
684         *Bioinformatics* **31**, 1674–1676 (2015).
685   51.   Nawrocki, E. P. & Eddy, S. R. Infernal 1.1: 100-fold faster RNA homology searches.
686         *Bioinformatics* **29**, 2933–2935 (2013).
687   52.   Kalvari, I. *et al.* Rfam 14: expanded coverage of metagenomic, viral and microRNA
688         families. *Nucleic Acids Res* **49**, D192–D200 (2021).
689   53.   Chan, P. P., Lin, B. Y., Mak, A. J. & Lowe, T. M. tRNAscan-SE 2.0: improved detection and
690         functional classification of transfer RNA genes. *Nucleic Acids Res* **49**, 9077–9096 (2021).
691   54.   Orakov, A. *et al.* GUNC: detection of chimerism and contamination in prokaryotic
692         genomes. *Genome Biol* **22**, 178 (2021).
693   55.   Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, giab008
694         (2021).

695   56.   Parks, D. H. *et al.* A complete domain-to-species taxonomy for Bacteria and Archaea. *Nat*
696         *Biotechnol* **38**, 1079–1086 (2020).
697   57.   Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and Effective
698         Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol Biol Evol* **32**,
699         268–274 (2015).
700   58.   Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using
701         MinHash. *Genome Biol* **17**, 132 (2016).
702   59.   Manni, M., Berkeley, M. R., Seppey, M., Simão, F. A. & Zdobnov, E. M. BUSCO Update:
703         Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage
704         for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes. *Mol Biol Evol* **38**, 4647–4654
705         (2021).
706   60.   Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic
707         tree display and annotation. *Nucleic Acids Res* **49**, W293–W296 (2021).
708   61.   Shen, W. & Ren, H. TaxonKit: A practical and efficient NCBI taxonomy toolkit. *Journal of*
709         *Genetics and Genomics* **48**, 844–850 (2021).
710   62.   Lu, J., Breitwieser, F. P., Thielen, P. & Salzberg, S. L. Bracken: estimating species
711         abundance in metagenomics data. *PeerJ Comput Sci* **3**, e104 (2017).
712   63.   McMurdie, P. J. & Holmes, S. phyloseq: An R Package for Reproducible Interactive
713         Analysis and Graphics of Microbiome Census Data. *PLoS One* **8**, e61217- (2013).
714   64.   Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion
715         for RNA-seq data with DESeq2. *Genome Biol* **15**, 550 (2014).
716   65.   Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–
717         2069 (2014).
718   66.   Steinegger, M. & Söding, J. Clustering huge protein sequence sets in linear time. *Nat*
719         *Commun* **9**, 2542 (2018).
720   67.   Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo
721         pipeline. *Genome Biol* **21**, 180 (2020).
722   68.   Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J. eggNOG-
723         mapper v2: Functional Annotation, Orthology Assignments, and Domain Prediction at the
724         Metagenomic Scale. *Mol Biol Evol* **38**, 5825–5829 (2021).
725   69.   Kanehisa, M. *et al.* Data, information, knowledge and principle: back to metabolism in
726         KEGG. *Nucleic Acids Res* **42**, D199–D205 (2014).
727   70.   Finn, R. D. *et al.* Pfam: the protein families database. *Nucleic Acids Res* **42**, D222–D230
728         (2014).
729   71.   Consortium, T. G. O. Gene Ontology Consortium: going forward. *Nucleic Acids Res* **43**,
730         D1049–D1056 (2015).
731   72.   Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol* **5**,
732         R12 (2004).

733

734

735   Figure legends
736   **Figure 1. The genome catalog assembled from the early-life skin samples.**

737    a. Schematic of study design from sampling to analysis. MAGs were constructed from single
738       samples and pooled samples based on the two body sites of the same infant at each time
739       point. MAGs from infant samples comprise the ELSG catalog. MAGs from mother samples
740       were used for comparative analysis.
741    b. Completeness and contamination for each of nonredundant prokaryotic and eukaryotic
742       MAGs included in the ELSG catalog, colored by the quality level.
743    c. Quality and completeness distribution for eukaryotic viral sequences included in the ELSG
744       catalog.
745
746    **Figure 2. Expansion of species diversity in skin microbiome.**
747    a. Rarefaction analysis of the number of species as a function of the number of nonredundant
748       genomes. Curves are depicted both for all the ELSG species and after excluding singleton
749       species (represented by only one genome).
750    b. Phylogenetic tree of the 1,029 representative bacterial MAGs in the ELSG catalog. Clades
751       are colored by GTDB phylum annotation (outer ring) and whether these are novel species
752       (inner shades). Bar graphs in the outermost layer indicate the number of nonredundant
753       genomes within each species-level cluster.
754    c. Comparison of species diversity between the ELSG catalog and the SMGC. Species-level
755       clusters were binned into the genus level in the bar graphs, ordered by a decreasing
756       number of ELSG-specific species.
757    d. Phylogenetic tree of the *Malassezia* genomes from the ELSG and the SMGC together with
758       GenBank reference genomes with *Saccharomyces cerevisiae* as the outgroup.
759    e. Number of infant samples harboring eukaryotic viruses included in the ELSG catalog.
760    f. Proportion of metagenomic reads from skin samples classified with Kraken 2 databases
761       based upon RefSeq, augmented by the SMGC and the ELSG. The boxes represent the
762       interquartile range, and the whiskers indicate the lowest and highest values within 1.5
763       times the interquartile range.
764
765    **Figure 3. Early-life skin microbial community structure.**
766    a. Principal coordinate analysis (PCoA) on Bray-Curtis dissimilarity between the microbial
767       profiles. Each point represents a single sample and is colored by body site and age group.
768       Ellipses represent a 95% confidence interval around the centroid of each sample group.
769    b. The Bray-Curtis dissimilarity of mother-infant pairs comparing related versus unrelated
770       dyads. Median value of each infant and all unrelated mothers was used. Statistical
771       difference was tested by two-sided Wilcoxon rank sum test.
772    c. Relative abundance of skin microbiome averaged for each sample group. Two of the most
773       abundant genera within each bacterial phylum were shown.
774    d. Differential taxa at the genus level between infants of different ages and between infants at
775       12 months and mothers. The size of the dots represents the log-transformed adjusted p-
776       value from DESeq2, and the color indicates fold changes. The top differentially abundant
777       genera for each comparison were shown.
778    e. Number of species-level MAGs recovered from infants at 2-3 months and 12 months, sorted
779       by the total number of MAGs.
780

**Figure 4. Proteins and functions of early-life skin microbiome.**
a. Rarefaction curves of the number of protein clusters obtained as a function of the number of species-level genomes. Each curve represents one species. The curves for species with more than 60 genomes are truncated for visualization purpose.
b. Comparison of the functional categories assigned to the core and accessory genes for species with at least 10 near-complete or high-quality genomes (>90% completeness, <5% contamination). Each dot represents one species. Odds ratio was calculated from the contingency table with core and accessory genes on one axis and the tested and the other functional categories on the other axis. Only significantly enriched functional categories are shown. Significance was calculated with a two-tailed t-test on log-transformed odds ratios and further adjusted for multiple comparisons using the Bonferroni correction.
c. Comparison of the protein clusters between the ELSG and the SMGC for species with at least 5 near-complete or high-quality genomes in each catalog.
d. Functional categories enriched in ELSG-specific and SMGC-specific genes compared to shared genes. Each dot represents a species. Only statistically significant categories are shown.

**Figure 5. Single-nucleotide variation indicates vertical transmission of skin microbiome.**
a. Top species with the largest intraspecies SNV density. The size of dots indicates the number of MAGs corresponding to each species.
b. Number of SNVs in pairwise comparisons between mother-infant pairs and between infants and unrelated mothers. Only species with genomes from at least 4 mother-infant pairs were considered for analysis. Statistical significance was tested by two-tailed Wilcoxon rank sum test. **P<0.01, ***P<0.001, ****P<0.0001.
c. Proportion of SNVs that were found in genomes from infants only or mothers only or both. SNVs were called based on the species-level representative MAG as the reference genome.
d. Phylogenetic tree of representative *C. acnes* cultured isolates with *C. modestum* as the outgroup. Source of individual is indicated in the label name and label color. Sequence type is displayed in parentheses.

**Figure S1. Quality metrics of the nonredundant genomes in the ELSG catalog.**
a. Genomes were stratified by quality level with colors matching those in Figure 1. Box lengths represent the interquartile range, and whiskers indicate the lowest and highest values within 1.5 times the interquartile range.
b. Eukaryotic viral sequences were stratified by quality level with colors matching those in Figure 1.

**Figure S2. Expansion of species diversity in the ELSG catalog.**
a. Comparison of species-level representative MAGs from three genome catalogs: ELSG, SMGC, and ELGG. The numbers indicate MAGs from each catalog that did or did not cluster with other catalogs at the species level.
b. Phylogenetic diversity accounted by known and novel species. Colors match Figure 2.
c. The improvement rate of read classification over the standard Kraken 2 RefSeq database.
d. Kraken 2 classification rate of reads from published skin metagenomic data.

825
826 **Figure S3. Early-life skin microbial community structure.**
827 a. The Bray-Curtis dissimilarity of microbial community between two time points of the same
828 infant compared with that of different infants. Median value of each infant and all other
829 infants was used to plot. Statistical difference was tested by two-sided Wilcoxon rank sum
830 test.
831 b. Relative abundance of skin microbiome averaged for each sample group after genome size
832 normalization, which emphasized the viral community.
833 c. Alpha diversity (richness and Shannon index) of skin samples divided by age group and skin
834 site.
835 d. Relationship between species prevalence and number of MAGs. Each dot represents a
836 species, the color of which indicates maximum relative abundance among all samples. The
837 prevalence of a species was calculated as the number of samples with >0.1% relative
838 abundance of that species. Pearson's correlation coefficient and p-value indicate a
839 significant correlation.
840
841 **Figure S4. Proteins and functions of early-life skin microbiome.**
842 a. Rarefaction curves of the number of protein clusters as a function of the number of
843 genomes for all species combined. Separate curves are depicted for proteins clustered at
844 90% amino acid identity for all protein clusters and after excluding singleton protein clusters
845 (containing only one protein sequence).
846 b. The number of genes in relation to the fraction of conspecific genomes where genes were
847 found. Only near-complete and high-quality genomes were considered in the analysis.
848 c. Number of genes shared by conspecific genomes in relation to the cutoff on the fraction of
849 conspecific genomes. Vertical dashed line represents the 90% threshold used in this study
850 to define core genes.
851 d. Proportion of core and accessory genes annotated by various databases including Clusters
852 of Orthologous Genes (COG), Kyoto Encyclopedia of Genes and Genomes (KEGG), Pfam, and
853 Gene Ontology (GO). Each dot represents a species. Only species with at least 10 near-
854 complete or high-quality genomes are included. Statistical significance was tested by two-
855 tailed Wilcoxon rank sum test.
856 e. Comparison of the rate of gene gain between the ELSG and the SMGC. Only near-complete
857 and high-quality genomes are included. Dashed line connects the end point of the collection
858 with fewer number of genomes of the same species.
859
860 **Figure S5. Intraspecies single-nucleotide variation and vertical transmission.**
861 a. Number of SNVs in pairwise comparisons between genomes of the same infant at 2-3
862 months and 12 months (same infant) and between any two genomes assembled from
863 different times (different infant). Only species with genomes from at least 3 infants at
864 different times were considered for analysis. Statistical significance was tested by two-tailed
865 Wilcoxon rank sum test. *P<0.05, **P<0.01. Non-significance ("ns") indicates P>0.05.
866 b. Number of sequence types of *C. acnes* cultured isolates from each individual.

867 c. Average number of shared sequence types of *C. acnes* cultured isolates between related
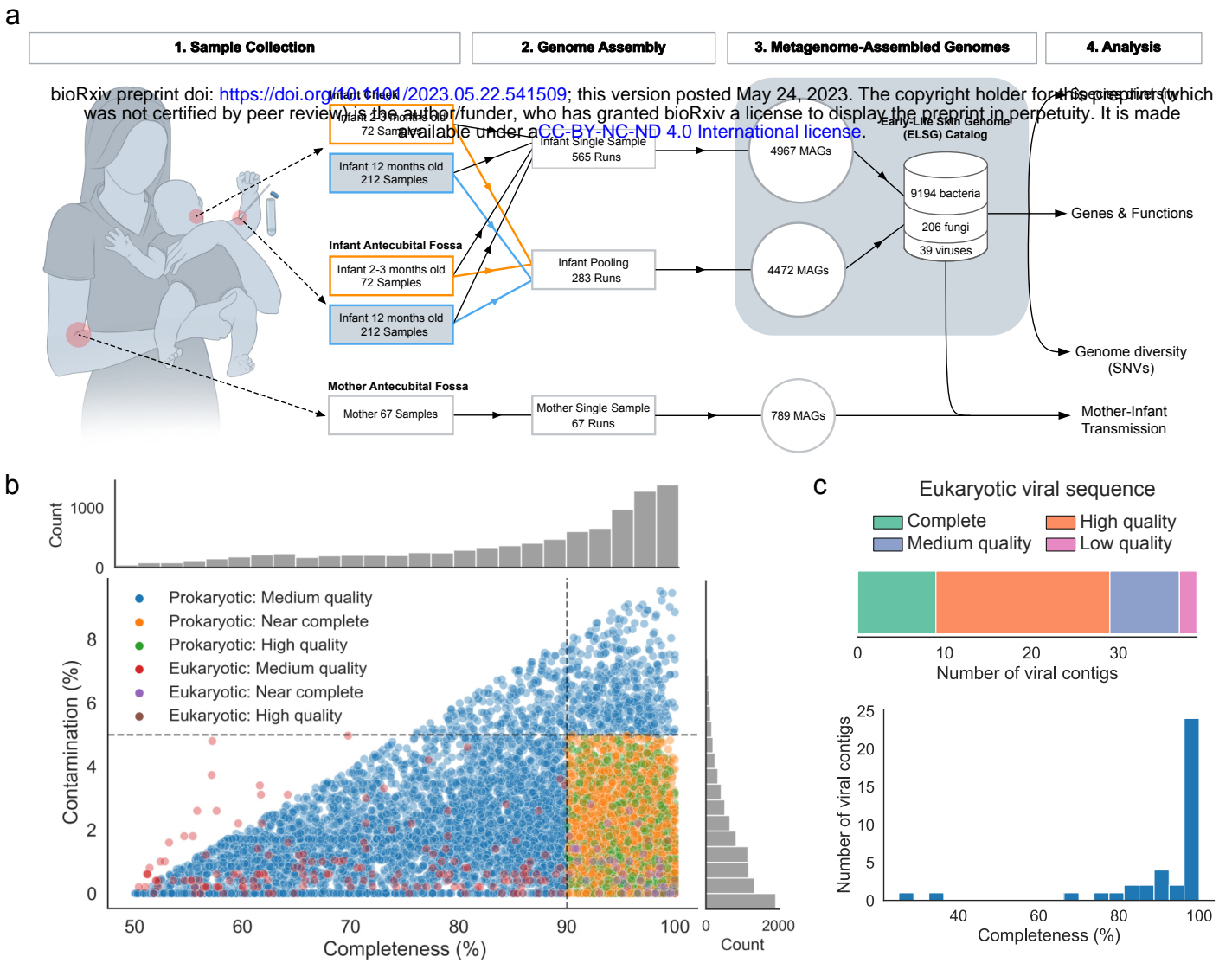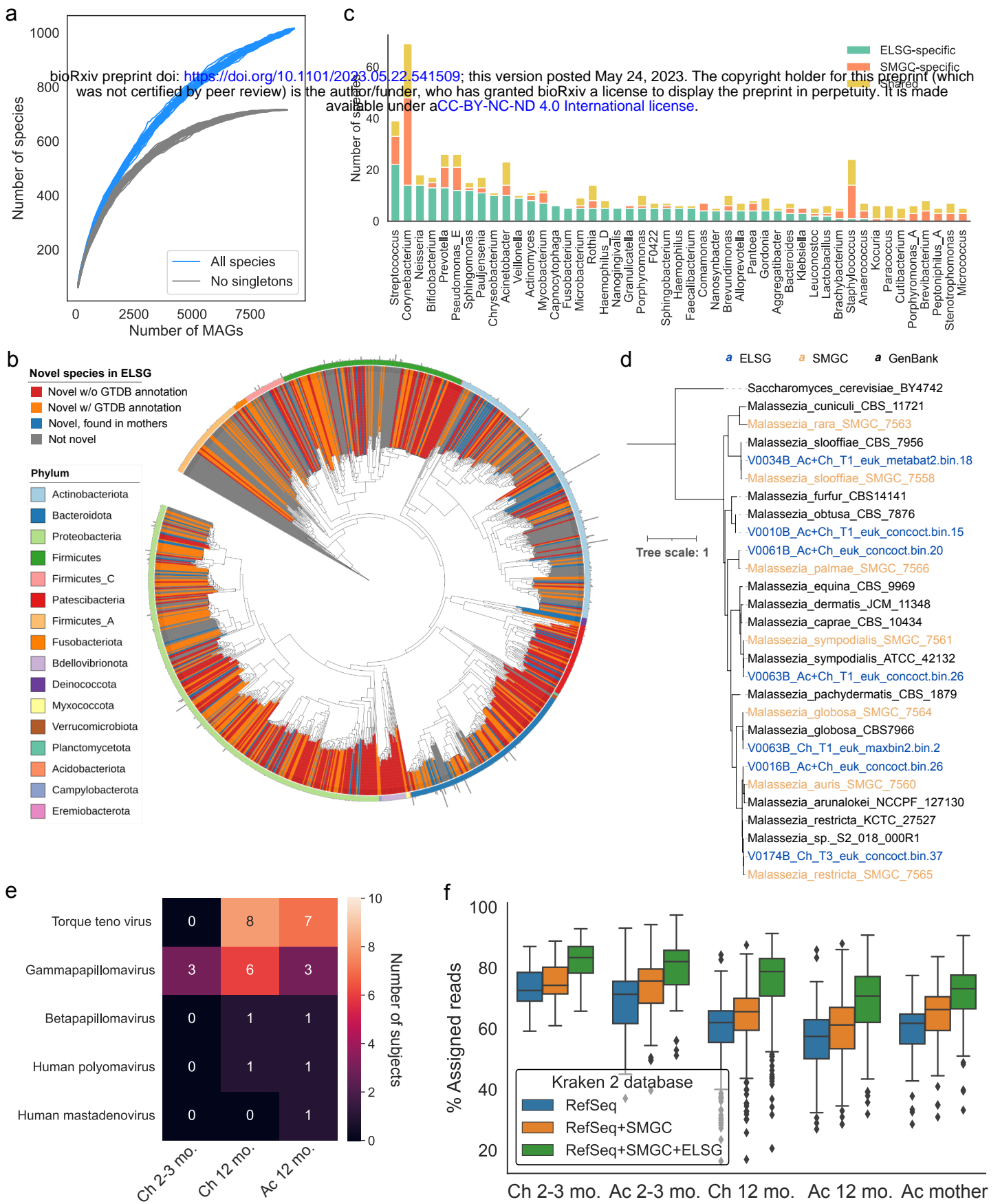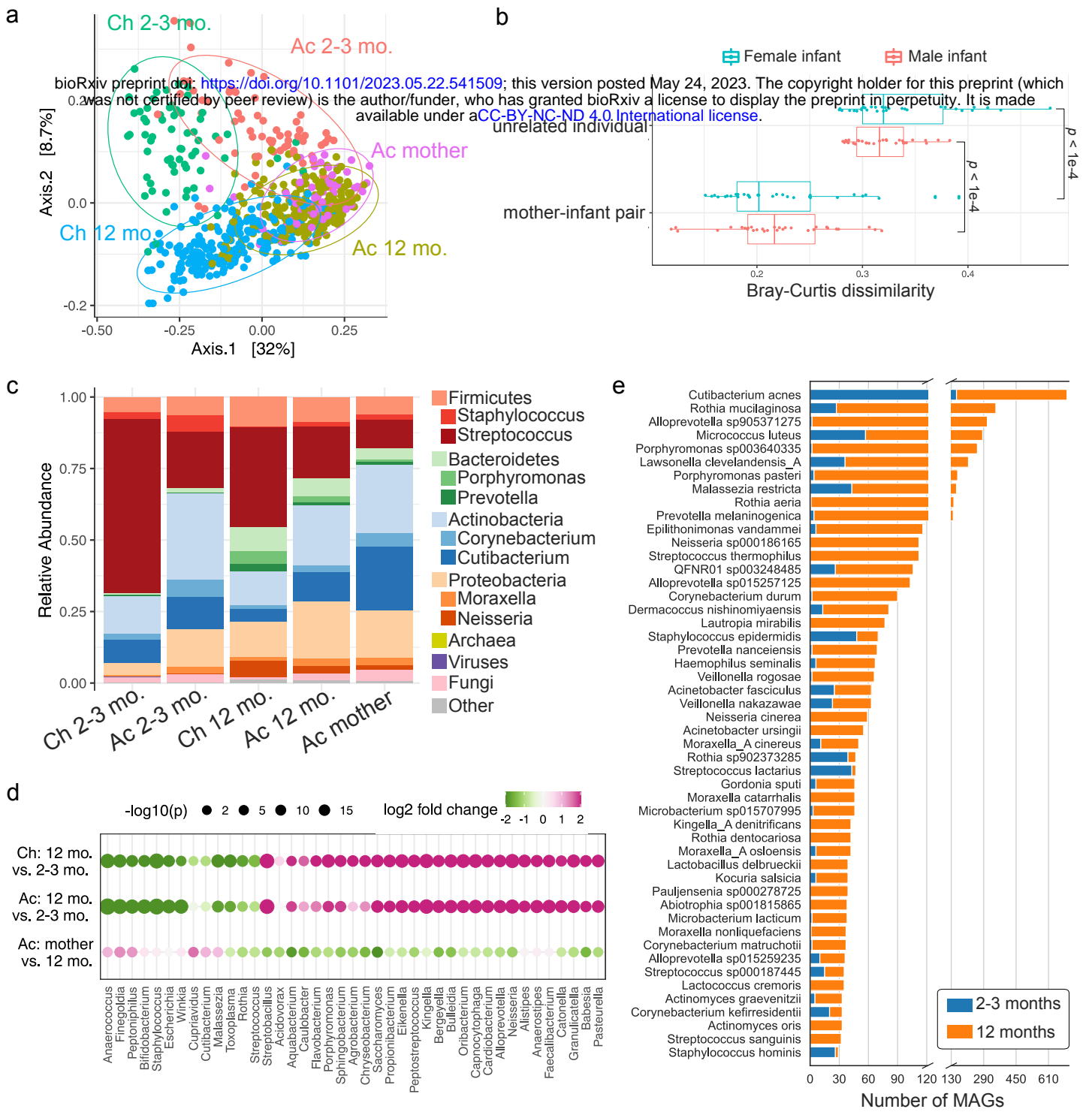868  infants and mothers (orange dashed line) and between any two individuals after 1,000
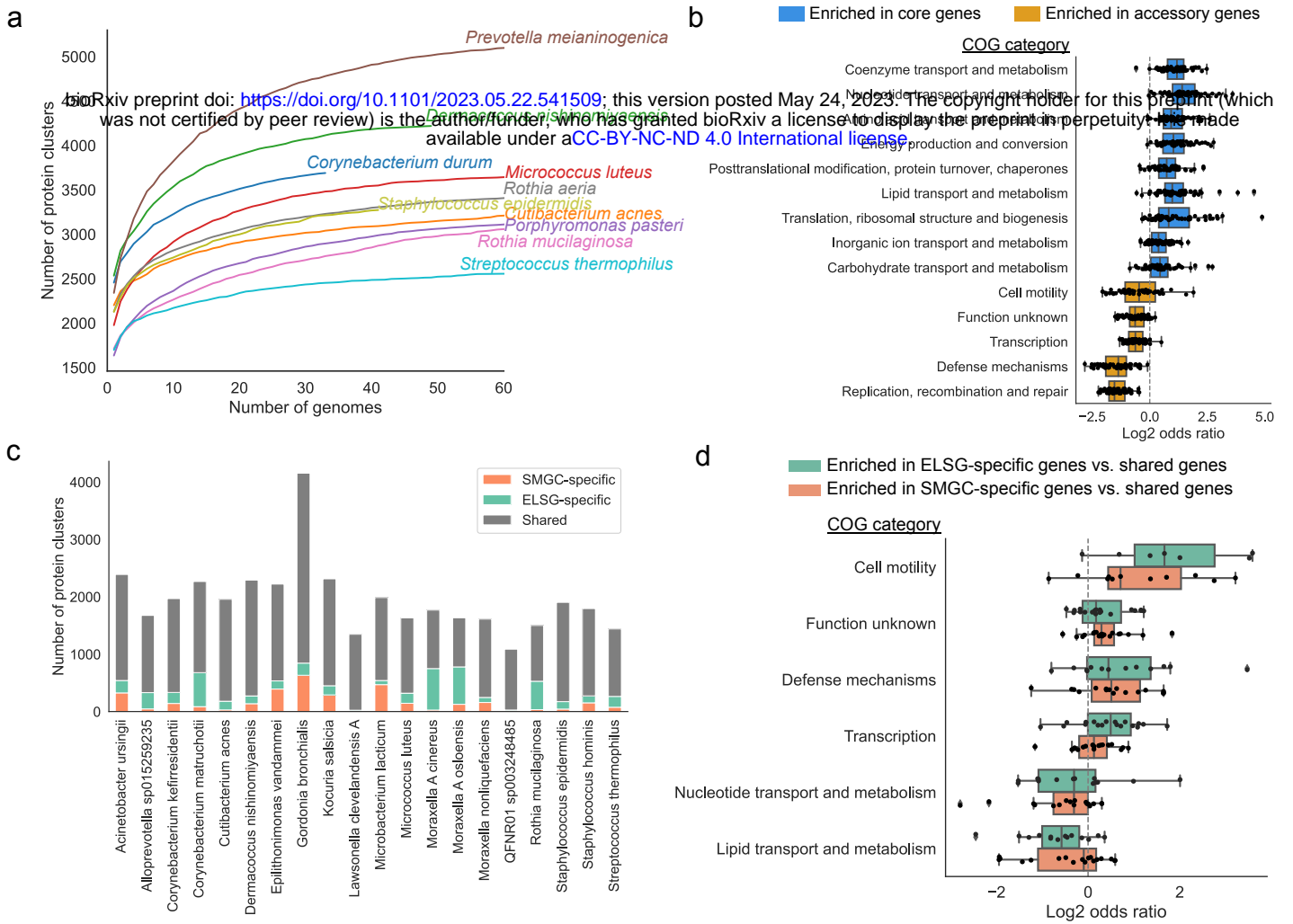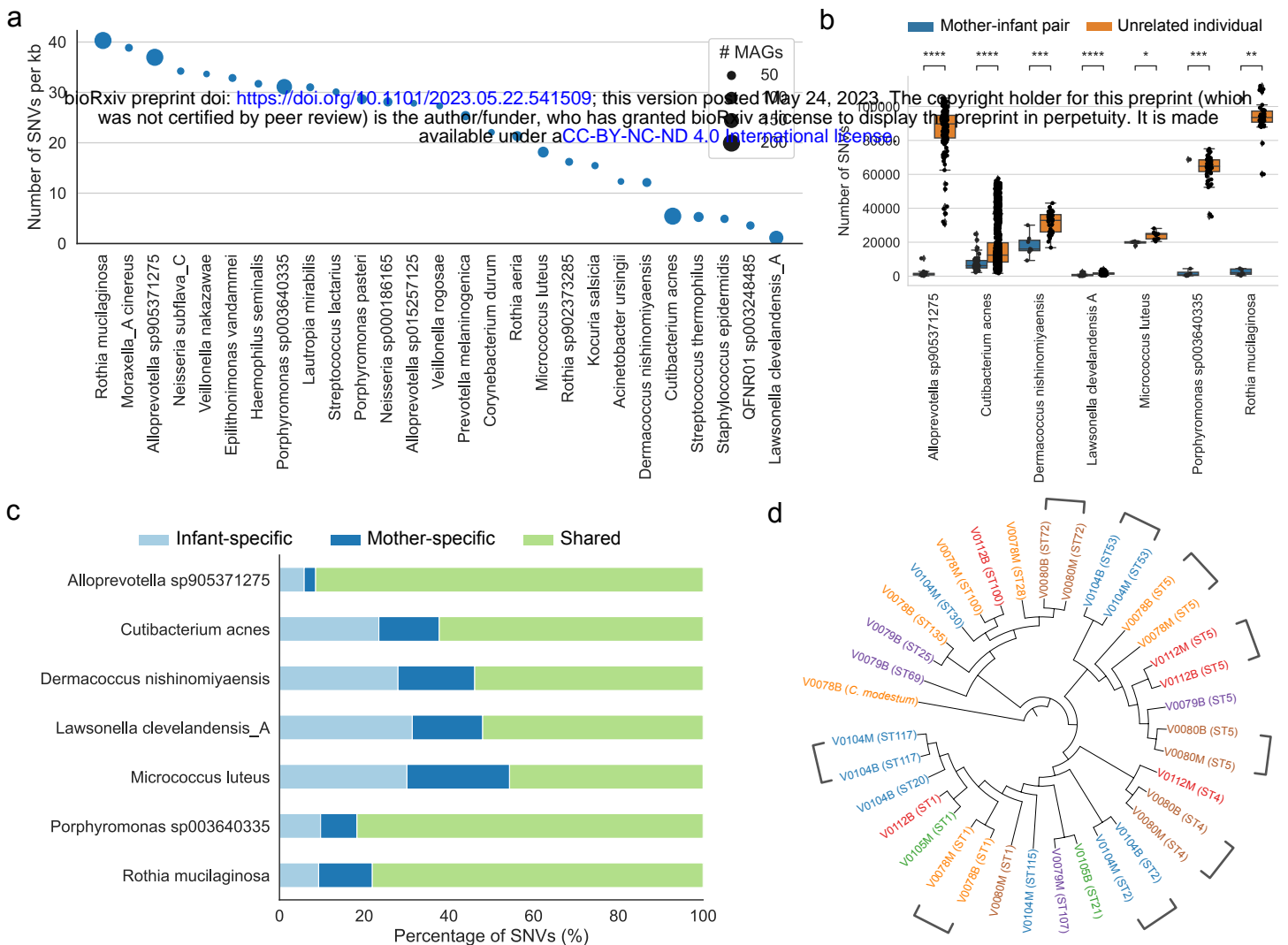869  permutations (histogram).

Figure 1

Figure 2

Figure 3

Figure 4

Figure 5