

Global Genetic Cartography of Urban Metagenomes and Anti-Microbial Resistance

David Danko^{1, 2, †}, Daniela Bezdán^{1, 2, †}, Ebrahim Afshinnekoo^{1, 2, *}, Sofia Ahsanuddin^{3, *}, Chandrima Bhattacharya^{1, 2, *}, Daniel J Butler^{1, 2, *}, Kern Rei Chng^{4, *}, Daisy Donnellan^{1, 2, *}, Jochen Hecht^{5, *}, Katerina Kuchin^{1, 2, *}, Mikhail Karasikov^{6, *}, Abigail Lyons^{1, 2, *}, Lauren Mak^{1, 2, *}, Dmitry Meleshko^{1, 2, *}, Harun Mustafa^{6, *}, Beth Mutai^{8, 9, *}, Russell Y Neches^{7, *}, Amanda Ng^{4, *}, Olga Nikolayeva^{10, *}, Tatyana Nikolayeva^{10, *}, Eileen Png^{4, *}, Krista Ryon^{1, 2, *}, Jorge L Sanchez^{1, 2, *}, Heba Shaaban^{1, 2, *}, Maria A Sierra^{1, 2, *}, Dominique Thomas^{1, 2, *}, Ben Young^{1, 2, *}, Omar O. Abudayyeh^{11, *}, Josue Alicea^{1, 2, *}, Malay Bhattacharyya^{12, 13, *}, Ran Blekhan^{14, *}, Eduardo Castro-Nallar^{15, *}, Ana M Cañas^{1, 2, *}, Aspasia D Chatziefthimiou^{1, 2, *}, Robert W Crawford^{16, *}, Francesca De Filippis^{17, 18, *}, Youping Deng^{19, *}, Christelle Desnues^{20, *}, Emmanuel Dias-Neto^{21, *}, Marius Dybwad^{22, *}, Eran Elhaik^{23, *}, Danilo Ercolini^{17, 18, *}, Alina Frolova^{24, *}, Dennis Gankin^{11, *}, Jonathan S. Gootenberg^{11, *}, Alexandra B Graf^{25, *}, David C Green^{26, *}, Iman Hajirasouliha^{1, 2, *}, Mark Hernandez^{27, *}, Gregorio Iraola^{28, 29, 30, *}, Soojin Jang^{31, *}, Andre Kahles^{6, *}, Frank J Kelly^{26, *}, Kaymisha Knights^{1, 2, *}, Nikos C Kyrpides^{7, *}, Paweł P Łabaj^{59, *}, Patrick K H Lee^{32, *}, Marcus H Y Leung^{32, *}, Per Ljungdahl^{33, *}, Gabriella Mason-Buck^{26, *}, Ken McGrath^{34, *}, Cem Meydan^{1, 2, *}, Emmanuel F Mongodin^{35, *}, Milton Ozorio Moraes^{36, *}, Niranjan Nagarajan^{4, *}, Marina Nieto-Caballero^{27, *}, Houtan Noushmehr^{37, *}, Manuela Oliveira^{38, *}, Stephan Ossowski^{39, 40, *}, Olayinka O Osulale^{41, *}, Orhan Özcan^{45, *}, David Paez-Espino^{7, *}, Nicolas Rascovan^{42, *}, Hugues Richard^{43, *}, Gunnar Rättsch^{6, *}, Lynn M Schriml^{35, *}, Torsten Semmler^{44, *}, Osman U Sezerman^{45, *}, Leming Shi^{46, 47, *}, Tielu Shi^{48, *}, Le Huu Song^{49, *}, Haruo Suzuki^{50, *}, Denise Syndercombe Court^{26, *}, Scott W Tighe^{51, *}, Xinzhao Tong^{32, *}, Klas I Udekwi^{33, *}, Juan A Ugalde^{52, *}, Brandon Valentine^{1, 2, *}, Dimitar I Vassilev^{53, *}, Elena Vayndorf^{54, *}, Thirumalaisamy P Velavan^{55, *}, Jun Wu^{48, *}, María M Zambrano^{56, *}, Jifeng Zhu^{1, 2, *}, Sibozhu^{57, 58, *}, Christopher E Mason^{1, 2, ‡}, and The International MetaSUB Consortium*

†Equal contribution

*Listed alphabetically

‡Corresponding author

*Full list attached

¹Weill Cornell Medicine

²The Bin Talal Bin Abdulaziz Al Saud Institute for Computational Biomedicine

³Icahn School of Medicine at Mount Sinai

⁴Genome Institute of Singapore

⁵Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Barcelona, Spain

⁶ETH Zurich, Department of Computer Science, Biomedical Informatics Group

⁷Department of Energy, Joint Genome Institute, Walnut Creek, California 94598, USA.

⁸Kenya Medical Research Institute / US Army medical Research Directorate - Kenya

⁹Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Barcelona, Spain.

¹⁰ETH Zurich, Functional Genomics Center Zurich

¹¹Massachusetts Institute of Technology, McGovern Institute for Brain Research

¹²Machine Intelligence Unit, Indian Statistical Institute, Kolkata

¹³Centre for Artificial Intelligence and Machine Learning, Indian Statistical Institute, Kolkata

¹⁴University of Minnesota

¹⁵Universidad Andrés Bello, Center for Bioinformatics and Integrative Biology, Facultad de Ciencias de la Vida

¹⁶California State University Sacramento

- ¹⁷Department of Agricultural Sciences, Division of Microbiology, University of Naples Federico II
¹⁸Task Force on Microbiome Studies, University of Naples Federico II
¹⁹University of Hawaii John A. Burns School of Medicine
²⁰Aix-Marseille Université, Mediterranean Institute of Oceanology, Université de Toulon, CNRS, IRD, UM 110
²¹A.C. Camargo Cancer Center
²²Norwegian Defence Research Establishment FFI, Kjeller, Norway
²³Department of Animal Plant Sciences, University of Sheffield
²⁴Institute of Molecular Biology and Genetics of National Academy of Science of Ukraine
²⁵University of Applied Sciences Vienna
²⁶Department of Analytical, Environmental and Forensic Sciences
²⁷University of Colorado at Boulder
²⁸Microbial Genomics Laboratory, Institut Pasteur de Montevideo, Uruguay
²⁹Center for Integrative Biology, Universidad Mayor, Santiago de Chile, Chile
³⁰Wellcome Sanger Institute, Hinxton, United Kingdom
³¹Institut Pasteur Korea
³²School of Energy and Environment, City University of Hong Kong, Hong Kong SAR, China
³³Stockholm University
³⁴Microba
³⁵University of Maryland School of Medicine, Institute for Genome Sciences
³⁶Fundação Oswaldo Cruz
³⁷University of São Paulo, Ribeirão Preto Medical School
³⁸Instituto de Patologia e Imunologia Molecular da Universidade do Porto
³⁹Institute of Medical Genetics and Applied Genomics, University of Tübingen, Tübingen, Germany
⁴⁰Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Barcelona, Spain.
3) Universitat Pompeu Fabra, Barcelona, Spain.
⁴¹Applied Environmental Metagenomics and Infectious Diseases Research (AEMIDR), Department of Biological Sciences, Elizade University
⁴²Aix-Marseille Université, IRD, AP-HM, IHU Méditerranée Infection
⁴³Sorbonne University, Faculty of science, Institute of Biology Paris-Seine, Laboratory of Computational and Quantitative Biology
⁴⁴Robert Koch Institute Berlin
⁴⁵Acibadem Mehmet Ali Aydınlar University
⁴⁶Center for Pharmacogenomics, School of Life Sciences and Shanghai Cancer Center, Fudan University
⁴⁷State Key Laboratory of Genetic Engineering (SKLGE) and MOE Key Laboratory of Contemporary Anthropology, School of Life Sciences
⁴⁸The Center for Bioinformatics and Computational Biology, Shanghai Key Laboratory of Regulatory Biology, the Institute of Biomedical Sciences and School of Life Sciences, East China Normal University
⁴⁹Institute of Tropical Medicine, Vietnamese-German Center of Excellence
⁵⁰Keio University
⁵¹University of Vermont
⁵²Millennium Initiative for Collaborative Research on Bacterial Resistance
⁵³Faculty of Mathematics and Informatics, Sofia University "St. Kliment Ohridski"
⁵⁴Institute of Arctic Biology, University of Alaska Fairbanks
⁵⁵Institute of Tropical Medicine, Univeristätsklinikum Tübingen, Tübingen
⁵⁶Corporación Corpogen
⁵⁷State Key Laboratory of Genetic Engineering (SKLGE) and MOE Key Laboratory of Contemporary Anthropology, School of Life Sciences, Fudan University
⁵⁸. Department of Epidemiology, School of Public Health, Fudan University
⁵⁹Małopolska Centre of Biotechnology, Jagiellonian University

Abstract

1
2 Although studies have shown that urban environments and mass-transit systems have distinct
3 genetic profiles, there are no systematic worldwide studies of these dense, human microbial ecosys-
4 tems. To address this gap in knowledge, we created a global metagenomic and antimicrobial resis-
5 tance (AMR) atlas of urban mass transit systems from 60 cities, spanning 4,728 samples and 4,424
6 taxonomically-defined microorganisms collected for three years. This atlas provides an annotated,
7 geospatial profile of microbial strains, functional characteristics, antimicrobial resistance markers,
8 and novel genetic elements, including 10,928 novel predicted viral species, 1302 novel bacteria, and
9 2 novel archaea. Urban microbiomes often resemble human commensal microbiomes from the skin
10 and airways, but also contain a consistent “core” of 31 species which are predominantly not human
11 commensal species. Samples show distinct microbial signatures which may be used to accurately

12 predict properties of their city of origin including population, proximity to the coast, and taxonomic
13 profile. These data also show that AMR density across cities varies by several orders of magnitude,
14 including many AMRs present on plasmids with cosmopolitan distributions. Together, these results
15 constitute a high-resolution, global metagenomic atlas, which enables the discovery of new genetic
16 components of the built human environment, highlights potential forensic applications, and provides
17 an essential first draft of the global AMR burden of the world's cities.

18 **Keywords:** Built Environment, metagenome, global health, antimicrobial resistance

1 Introduction

The high-density urban environment has historically been home to only a fraction of all people, with the majority living in rural areas or small villages. In the last two decades, the situation has reversed; 55% of the world's population now lives in urban areas (Ritchie and Roser, 2020; United Nations, 2018). Since the introduction of germ theory and John Snow's work on cholera, it has been clear that people in cities interact with microbes in ways that can be markedly different than in rural areas (Neiderud, 2015). Microbes in the built environment have been implicated as a possible source of contagion (Cooley et al., 1998) and certain syndromes, like allergies, are associated with increasing urbanization (Nicolaou et al., 2005). It is now apparent that cities in general have an impact on human health though the mechanisms of this impact are broadly variable and often little understood. Indeed, our understanding of microbial dynamics in the urban environment outside of pandemics has only begun (Gilbert and Stephens, 2018).

Technological advances in next-generation sequencing (NGS) and metagenomics have created an unprecedented opportunity for rapid, global studies of microorganisms and their hosts, providing researchers, clinicians, and policymakers with a more comprehensive view of the functional dynamics of microorganisms in a city. NGS facilitates culture-independent sampling of the microorganisms in an area with the potential for both taxonomic and functional annotation; this is particularly important for surveillance of microorganisms as they acquire antimicrobial resistance (AMR) (Fresia et al., 2019). Metagenomic methods enable nearly real-time monitoring of organisms, AMR genes, and pathogens as they emerge within a given geographical location, and have the potential to reveal hidden microbial reservoirs and detect microbial transmission routes as they spread around the world (Zhu et al., 2017). There are several different drivers and sources for AMR; including agriculture, farming, and livestock in rural and suburban areas, household and industrial sewage, usage of antimicrobials, hard metals, and biocides, as well as human and animal waste, all these factors contribute to the complexity of AMR transmission (Allen et al., 2009; Martínez, 2008; Singer et al., 2016; Thanner et al., 2016; Venter et al., 2017). A molecular map of urban environments will enable significant new research on the impact of urban microbiomes on human health.

The United Nations projects that by 2050, over two-thirds of the world's population will live in urban areas (Ritchie and Roser, 2020). Consequently, urban transit systems - including subways and buses - are a daily contact interface for billions of people who live in cities. Notably, urban travelers bring their commensal microorganisms with them as they travel and come into contact with organisms and mobile elements present in the environment, including AMR markers. The study of the urban microbiome and the microbiome of the built environment spans several different projects and initiatives including work focused on transit systems (Afshinnekoo et al., 2015; Hsu et al., 2016; Kang et al., 2018; Leung et al., 2014; MetaSUB International Consortium. Mason et al., 2016), hospitals (Brooks et al., 2017; Lax et al., 2017), soil (Hoch et al., 2019; Joyner et al., 2019), and sewage (Fresia et al., 2019; Maritz et al., 2019), among others. However, these efforts for the most part have only been profiled with comprehensive metagenomic methods in a few selected cities on a limited number of occasions. This leaves a gap in scientific knowledge about a microbial ecosystem, with which the global human population readily interacts. Human commensal microbiomes have been found to vary widely based on culture, and thus the geography and geographically constrained studies may to miss key differences (Brito et al., 2016). Moreover, data on urban microbes and AMR genes are urgently needed in developing nations, where antimicrobial drug consumption is expected to rise by 67% by 2030 (United Nations, 2016; Van Boeckel et al., 2015), both from changes in consumer demand for livestock products and an expanding use of antimicrobials - both of which can alter AMR profiles of these cities.

The International Metagenomics and Metadesign of Subways and Urban Biomes (MetaSUB) Consortium was launched in 2015 to address this gap in knowledge on the density, types, and dynamics of urban metagenomes and AMR profiles. Since then, we have developed standardized collection and sequencing protocols to process 4,728 samples across 60 cities worldwide (Table S1). Sampling took place at three major time points: a pilot study in 2015-16 and two global city sampling days (gCSD, June 21st) in 2016 and 2017. Each sample was sequenced with 5-7M 125bp paired-end reads using Illumina NGS sequencers (see Methods). To deal with the challenging analysis of our large dataset, we generated an open-source analysis pipeline (MetaSUB Core Analysis Pipeline, CAP), which includes a comprehensive set of state-of-the-art, peer-reviewed, metagenomic tools for taxonomic identification, *k*-mer analysis, AMR gene prediction, functional profiling, de novo assembly, annotation of particular microbial species, and geospatial mapping.

To our knowledge this study represents the first and largest global metagenomic study of urban microbiomes - with a focus on transit systems - that reveals a consistent "core" urban microbiome across

76 all cities, as well as distinct geographic variation that may reflect epidemiological variation and that
77 enables a new forensic, source-tracking capabilities. More importantly, our data demonstrate that a
78 significant fraction of the urban microbiome remains to be characterized. Though 1,000 samples are
79 sufficient to discover roughly 80% of the observed taxa and AMR markers, we continued to observe
80 taxa and genes at an ongoing discovery rate of approximately one new species (previously non-observed)
81 and one new AMR marker for every 10 samples. Notably, this genetic variation is affected by various
82 environmental factors (e.g., climate, surface type, latitude, etc.) and samples show greater diversity near
83 the equator. Moreover, sequences associated with AMR markers are widespread, though not necessarily
84 abundant, and show geographic specificity. Here, we present the results of our global analyses and a
85 set of tools developed to access and analyze this extensive atlas, including: two interactive map-based
86 visualizations for samples (metasub.org/map) and AMRs (resistanceopen.org), an indexed search tool
87 over raw sequence data (dnaloc.ethz.ch/), a Git repository for all analytical pipelines and figures, and
88 application programming interfaces (APIs) for computationally accessing results ([github.com/metasub/
89 metasub_utils](https://github.com/metasub/metasub_utils)).

90 2 Results

91 We have collected 4,728 samples from from the mass transit systems of 60 cities around the world
92 (Table 1, Supplementary table S1). These samples were collected from various common surfaces in the
93 mass transit systems such as railings, benches, and ticket kiosks and were subjected to metagenomic
94 sequencing. We use the microbiome of mass transit systems as a proxy for the urban microbiome as a
95 whole and present our key findings here.

96 A Core Urban Microbiome Centers Global Diversity

97 We first investigated the distribution of microbial species across the global urban environment. Specifi-
98 cally, we asked whether the urban environment represents a singular type of microbial ecosystem or a set
99 of related, but distinct, communities, especially in terms of biodiversity. We observed a bi-modal distri-
100 bution of taxa prevalence across our dataset, which we used to define two separate sets of taxa based on
101 the inflection points of the distribution: the putative “sub-core” set of urban microbial species that are
102 consistently observed (>70% of samples) and the less common “peripheral” (<25% of samples) species.
103 We also defined a set of true “core” taxa which occur in essentially all samples (>97% of samples). Apply-
104 ing these thresholds, we identified 1,145 microbial species (Figure 2C) that make up the sub-core urban
105 microbiome with 31 species in the true core microbiome (Figure 2A). Core and sub-core taxa classifica-
106 tions were further evaluated for sequence complexity and genome coverage on a subset of samples. Of
107 the 1,206 taxa with prevalence greater than 70%, 69 were flagged as being low quality classifications (see
108 methods). The sub-core microbiome was principally bacterial, with just one eukaryotic taxon identified
109 and not flagged: *Saccharomyces cerevisiae*. Notably, no archaea or viruses were identified in the group of
110 sub-core microorganisms (note that this analysis did not include viruses newly discovered in this study).
111 For viruses in particular, this may be affected by the sampling or DNA extraction methods used, by
112 limitations in sequencing depth, or by missing annotations in the reference databases used for taxonomic

Table 1: Sample Counts, The number of samples collected from each region.

Region	Pilot	CSD16	CSD17	Other	Total
North America	28	284	371	276	959
East Asia	34	26	1297	0	1357
Europe	177	310	939	1	1427
Sub Saharan Africa	0	116	192	0	308
South America	20	44	199	68	331
Middle East	0	100	15	0	115
Oceania	0	94	32	0	126
Background Control	0	0	40	0	40
Lab Control	0	0	20	6	26
Positive Control	0	0	33	6	39
Total	259	974	3138	357	4728

113 classification, which is principally problematic with phages. It is worth noting that potentially prevalent
114 RNA viruses are omitted with our DNA-based sampling. The three most common bacterial phyla across
115 the world's cities ordered by the number of species observed were *Proteobacteria*, *Actinobacteria*, and
116 *Firmicutes*. To test for possible geographic bias in our data, we normalized the prevalence for each taxa
117 by the median prevalence within each city. The two normalization methods broadly agreed (Figure 2).

118 Despite their global prevalence, the core taxa are not uniformly abundant across all cities. Many
119 species exhibited a high standard deviation and kurtosis (calculated using Fisher's definition and normal
120 kurtosis of 0) than other species (Figure 2B). Furthermore, some species show distinctly high mean
121 abundance, often higher than the core species, but more heterogeneous global prevalence. For example,
122 *Salmonella enterica* is identified in less than half of all samples but is the 12th most abundant species
123 based on the fraction of DNA that can be ascribed to it. The most relatively abundant microbial species
124 was *Cutibacterium acnes* (Figure 2D) which had a comparatively stable distribution of abundance across
125 all samples; *Cutibacterium acnes* is known as a prominent member of the human skin microbiome. To
126 test for any biases arising from uneven geographic sampling, we measured the relative abundance of
127 each taxon by calculating the fraction of reads classified to each particular taxon, and compared the
128 raw distribution of abundance to the distribution of median abundance within each city (This process
129 is analogous to the one used for Figure 2C, Figure 2B); the two measures closely aligned. Also, an
130 examination of the positive and negative controls indicates that these results are not likely due to
131 contamination or batch effect (Supp. Figure S13). In total, we observed 31 core taxa (>97%), 1,145
132 sub-core taxa (70-97%) 2,466 peripheral taxa (<25%), and 4,424 taxa across all samples. We term the
133 set of all taxa observed *the urban panmicrobiome*.

134 To estimate the number of taxa present in our samples but which were not detected by our experi-
135 mental techniques, we performed a rarefaction analysis on the taxa that were identified. By estimating
136 the number of taxa identified for different numbers of samples, we see a diminishing trend (Figure 2D),
137 which indicates that at some point, the species in every new sample were likely already identified in a
138 previous one. Our rarefaction curve did not reach a plateau and, even after including all samples, it still
139 shows an expected marginal discovery rate of roughly 1 additional species for every 10 samples added
140 to the study. For clarity we note that this analysis only considers taxa already present in reference
141 databases, not newly discovered taxa (below). Despite the remaining unidentified taxa, we estimate
142 that most (80%) of the classifiable taxa in the urban microbiome could be identified with roughly 1,000
143 samples. However, as noted below, this new diversity is likely not evenly distributed across regions.

144 As humans are a major part of the urban environment, the DNA in our samples could be expected to
145 resemble commensal human microbiomes. To investigate this, we compared non-human DNA fragments
146 from our samples to a randomized set of 50 samples from 5 commensal microbiome sites in the Human
147 Microbiome Project (HMP) (Consortium et al., 2012) (stool, skin, airway, gastrointestinal tract, urogen-
148 ital tract). We used MASH to perform a k -mer based comparison of our samples vs. the selected HMP
149 samples, which showed a roughly uniform dissimilarity between MetaSUB samples and those from dif-
150 ferent human body sites (Figure 2E, Supp. Figure S2A B). Samples taken from surfaces that were likely
151 to have been touched more often by human skin, such as doorknobs, buttons, railings, and touchscreens,
152 were indeed more similar to human skin microbiomes than surfaces like bollards, windows, and the floor.
153 Given that a large fraction of DNA in our samples could not be classified and that a k -mer based compar-
154 ison did not find significant body-site specificity, it is possible that the unclassified DNA in our samples
155 is from novel taxa which are not human commensals. Of note, the taxonomic composition of our samples
156 do not closely resemble soil samples. We processed 28 metagenomic soil samples (Bahram et al., 2018)
157 using the same pipeline as the rest of the data and compared soil samples to our samples using MASH.
158 Our samples were very dissimilar from the soil samples (Figure 2F) even in comparison to human skin
159 microbiomes. This suggests that the unclassified DNA may represent heretofore uncharacterized taxa
160 that are not known commensals being shed into the environment.

161 We next estimated the fraction of sequences in our data that did not resemble sequences in known
162 reference databases. We took a subset of 10,000 reads from each sample and aligned these reads to
163 a number of reference databases using BLASTn (Altschul et al., 1990). We then identified reads that
164 mapped to sequences in the reference databases at 80%, 90%, and 95% Average Nucleotide Identity
165 (ANI) (Figure 2G). We used a broad set of databases for reference: RefSeq, NCBI's NT Environmental,
166 a large database of Metagenome Assembled Genomes (MAGs) from Pasolli et al. (2019), and MAGs from
167 MetaSUB itself (Section 2.4). At 80% ANI, the most permissive threshold, 34.6% of reads did not map
168 to any database while 47.3% of reads did not map or only mapped to MAGs from MetaSUB itself. This
169 mirrors results seen by previous urban microbiome works (Afshinnekoo et al., 2015; Hsu et al., 2016).

170 Next, we analyzed the fraction of sequences that aligned to these same databases by region. Sur-

171 prisingly, samples from Europe had the highest fraction of unaligned reads, followed by the middle east,
172 while samples from Sub Saharan Africa had the smallest fraction of unaligned reads (Supp. Figure
173 S1C). The proportion of reads aligned to each database did not vary significantly by region. We fur-
174 ther investigated the relationship between geography and sample composition. In ecology, an increasing
175 distance from the equator is associated with a decrease in taxonomic diversity (O'Hara et al., 2017).
176 The MetaSUB data recapitulates this result and identifies a significant decrease in taxonomic diversity
177 (though with significant noise, $p < 2e16$, $R^2 = 0.06915$) as a function of absolute latitude; samples are
178 estimated to lose 6.9672 species for each degree of latitude away from the equator (Supp. Figure S1A).
179 The effect of latitude on species diversity is not purely monotonic, since several cities have higher species
180 diversity than their latitude would predict. This is expected as latitude is only a rough predictor of a
181 city's climate. While this is an observation consistent with ecological theory, we note that our samples
182 are heavily skewed by the location of the target cities, as well as the prevalence of those cities in specific
183 latitude zones of the northern hemisphere.

184 2.1 Global Diversity Varies According to Covariates

185 Despite the core urban microbiome present in almost all samples, there was also geographic variation
186 in taxonomy and localization. We calculated the Jaccard distance between samples measured by the
187 presence and absence of species (which is robust to noise from relative abundance) and performed a
188 dimensionality reduction of the data using UMAP (Uniform Manifold Approximation and Projection,
189 McInnes et al. (2018)) for visualization (Figure 2A). Jaccard distance was correlated with distance based
190 on Jensen-Shannon Divergence (which accounts for relative abundance) and k -mer distance calculated by
191 MASH (which is based on the k -mer distribution in a sample, so cannot be biased by a database) (Supp.
192 Figure S10A, B, C). In principle, Jaccard distance could be influenced by read depth as low abundance
193 species drop below detection thresholds. However we expect this issue to be minor as the total number
194 of species identified stabilized at 100,000 reads (Supp. Figure S9B) compared to an average of 6.01M
195 reads per sample. Samples collected from North America and Europe were distinct from those collected
196 in East Asia, but the separation between other regions was less clear. A similar trend was found in an
197 analogous analysis based on functional pathways rather than taxonomy (Supp. Fig S5D), which indicates
198 geographic stratification of the metagenomes at both the functional and taxonomic levels. Subclusters
199 identified by UMAP roughly corresponded to city and climate but not surface type (Supp. Figure S5A,
200 B, C). These findings confirm and extend earlier analyses performed on a fraction of the MetaSUB data
201 which were run as a part of CAMDA Challenges in years 2017, 2018, and 2019 (camda.info).

202 We quantified the degree to which metadata covariates influence the taxonomic composition of our
203 samples using MAVRIC, a statistical tool to estimate the sources of variation in a count-based dataset
204 (Moskowitz and Greenleaf, 2018). We identified covariates which influenced the taxonomic composition
205 of our samples: city, population density, average temperature in June, region, elevation above sea-level,
206 surface type, surface material, elevation above or below ground and proximity to the coast. The most
207 important factor, which could explain 19% of the variation in isolation, was the city from which a sample
208 was taken followed by region which explained 11%. The other four factors ranged from explaining 2%
209 to 7% of the possible variation in taxonomy in isolation (Supp. Table S2). We note that many of
210 the factors were confounded with one another, so they can explain less diversity than their sum. One
211 metadata factor tested, the population density of the sampled city, had no significant effect on taxonomic
212 variation overall.

213 To quantify how the principle covariates, climate, continent, and surface material impacted the taxo-
214 nomic composition of samples, we performed a Principal Component Analysis (PCA) on our taxonomic
215 data normalized by proportion and identified principal components (PCs) which were strongly associated
216 with a metadata covariate in a positive or negative direction (PCs were centered so an average direction
217 indicates an association). We found that the first two PCs (representing 28.0% and 15.7% of the variance
218 of the original data, respectively) associated strongly with the city climate while continent and surface
219 material associate less strongly (Figure 2B).

220 Next, we tested whether geographic proximity (in km) of samples to one another had any effect on
221 the variation, since samples taken from nearby locations could be expected to more closely resemble one
222 another. Indeed, for samples taken in the same city, the average JSD (Jensen-Shannon distance) was
223 weakly predictive of the taxonomic distance between samples, with every increase of 1km in distance
224 between two samples representing an increase of 0.056% in divergence ($p < 2e16$, $R^2 = 0.01073$, Supp.
225 Figure S1B). This suggests a "neighborhood effect" for sample similarity analogous to the effect described
226 by Meyer et al. (2018), albeit a very minor one. To reduce bias that could be introduced by samples



Figure 2: Differences at global scale A) UMAP of taxonomic profiles based on Jaccard distance between samples. Colored by the region of origin for each sample. Axes are arbitrary and without meaningful scale. The color key is shared with panel B. B) Association of the first 25 principal components of sample taxonomy with climate, continent, and surface material. C) Distribution of major phyla, sorted by hierarchical clustering of all samples and grouped by continent. D) Distribution of high-level groups of functional pathways, using the same order as taxa (C). E) Distribution of AMR genes by drug class, using the same order as taxa (C). Note that MLS is macrolide-lincosamide-streptogramin.

227 taken from precisely the same object we excluded all pairs of samples within 1km of one another.

228 At a global level, we examined the prevalence and abundance of taxa and their functional profiles
229 between cities and continents. These data showed a fairly stable phyla distribution across samples, but
230 the relative abundance of these taxa is unstable (Figure 2C) with some continental trends. In contrast
231 to taxonomic variation, functional pathways were much more stable across continents, showing relatively
232 little variation in the abundance of high level categories (Figure 2D). This pattern may also be due to
233 the more limited range of pathway classes and their essential role in cellular function, in contrast to the
234 much more wide-ranging taxonomic distributions examined across metagenomes. Classes of antimicrobial
235 resistance were observed to vary by continent as well. Clusters of AMR classes were observed to occur
236 in groups of taxonomically similar samples (Figure 2E).

237 We quantified the relative variation of taxonomic and functional profiles by comparing the distribution
238 of pairwise distances in taxonomic and functional profiles. Both profiles were equivalently normalized
239 to give the probability of encountering a particular taxon or pathway. Taxonomic profiles have a mean
240 pairwise Jensen-Shannon Divergence (JSD) of 0.61 while pathways have a mean JSD of 0.099. The
241 distributions of distances are significantly different (Welch's *t*-test, unequal variances, $p < 2e - 16$). This
242 is consistent with observations from the Human Microbiome Project, where metabolic function varied
243 less than taxonomic composition (Consortium et al., 2012; Lloyd-Price et al., 2017) within samples from
244 a given body site.

245 2.2 Microbial Signatures Reveal Urban Characteristics

246 To facilitate characterization of novel sequences we created GeoDNA, a high-level web interface (Figure
247 3A) to search raw sequences against our dataset. Users can submit sequences to be processed against
248 a *k*-mer graph-based representation of our data. Query sequences are mapped to samples and a set of
249 likely sample hits is returned to the user. This interface will allow researchers to probe the diversity in
250 this dataset and rapidly identify the range of various genetic sequences.

251 We sought to determine whether a samples taxonomy reflected the environment in which it was
252 collected. To this end we trained a Random Forest Classifier (RFC) to predict a sample's city of origin
253 from its taxonomic profile. We trained an RFC with 100 components on 90% of the samples in our
254 dataset and evaluated its classification accuracy on the remaining 10%. We repeated this procedure with
255 multiple subsamples of our data at various sizes and with 5 replicates per size to achieve a distribution
256 (Fig. 3B). The RFC achieved 88% on held out data which compares favorably to the 7.01% that would
257 be achieved by a randomized classifier. These results from our RFC demonstrate that city specific
258 taxonomic signatures exist and can be predictive.

259 We expanded our analysis of environmental signatures in taxonomy to the prediction of features in
260 cities not present in our training set. To do this we collated a set of 7 features for each city: population,
261 surface material, elevation, proximity to the coast, population density, region, ave June temperature,
262 and Koppen climate classification. We trained a RFCs to predict each feature based on all samples that
263 were not taken from a given city then used the relevant RFC to predict the feature for samples from
264 the held out city and recorded the classification accuracy (Figure 3D). While not all features and cities
265 were equally predictable (in particular features for a number of British cities were roughly similar and
266 could be predicted effectively) in general the predictions exceeded random chance by a significant margin
267 (Supp. Figure S3A). This suggests that certain features of cities generate microbial signatures that are
268 present globally and distinct from city specific signatures. The successful geographic classification of
269 samples demonstrates distinct city-specific trends in the detected taxa, that may enable future forensic
270 biogeographical capacities.

271 However, unique, city-specific taxa are not uniformly distributed (Figure 3B). To quantify this, we
272 developed a score to reflect how endemic a given taxon is within a city, which reflects upon the forensic
273 usefulness of a taxon. We define the Endemicity Score (ES) of a taxa as term-frequency inverse document
274 frequency where the document consists of samples from some metadata defined group such as a city or
275 region. This score is designed to simultaneously reflect the chance that a taxon could identify a given
276 city and that that taxon could be found within the given city. A high ES for a taxon in a given city
277 could be evidence of the evolutionary advantage that the taxon has in a particular cities environment.
278 However, neutral evolution of microbes within a particular niche is also possible and the ES alone does
279 not distinguish between these two hypotheses.

280 Note that while the ES only considers taxa which are found in a city, a forensic classifier could also
281 take advantage of the absence of taxa for a similar metric. ES show a roughly bimodal distribution for
282 regions (Fig. 3C). Each region possesses a number of taxa with ES scores close to 1 and a slightly larger

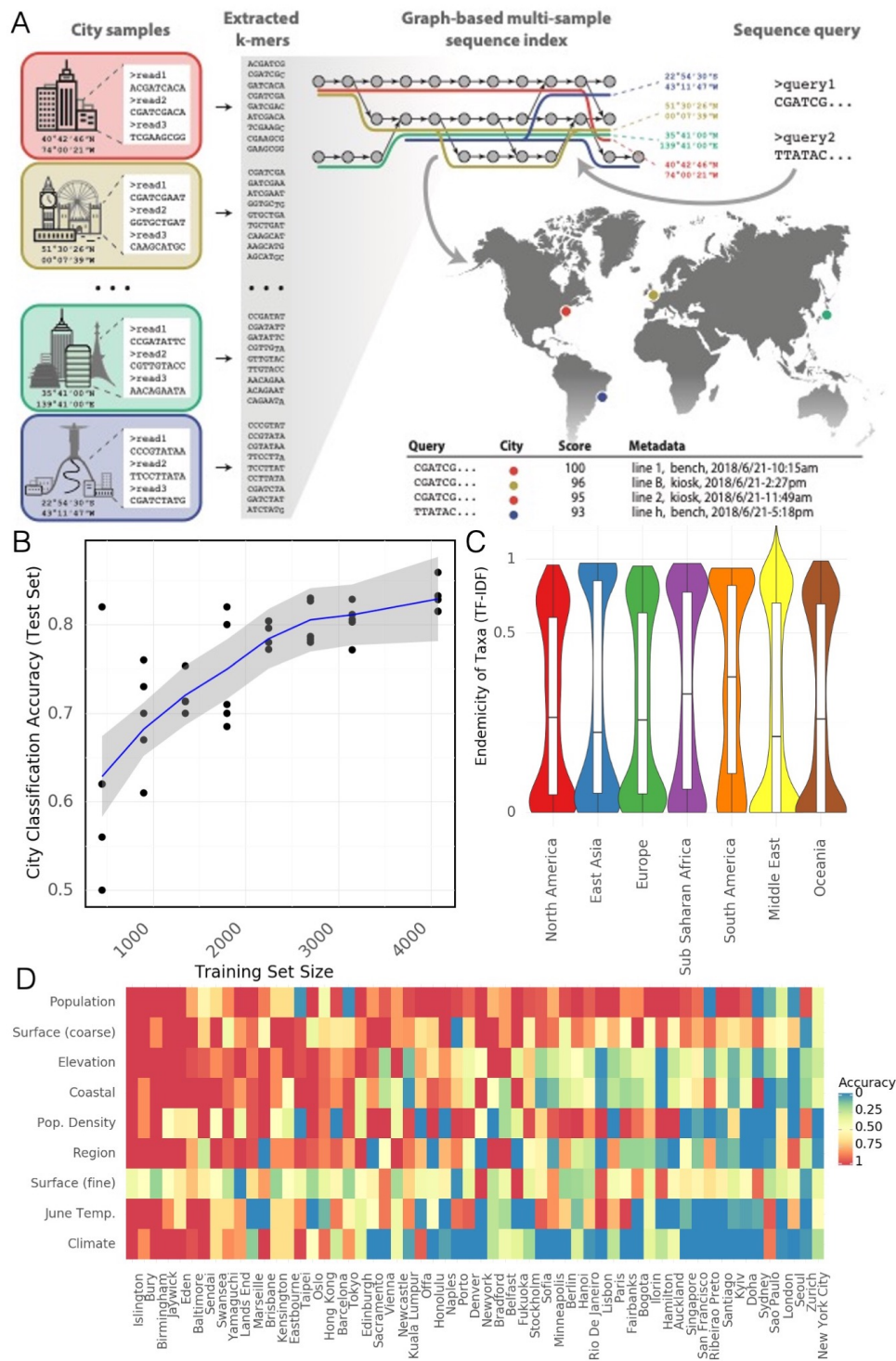


Figure 3: Microbial Signatures A) Schematic of GeoDNA representation generation – Raw sequences of individual samples for all cities are transformed into lists of unique *k*-mers (left). After filtration, the *k*-mers are assembled into a graph index database. Each *k*-mer is then associated with its respective city label and other informative metadata, such as geo-location and sampling information (top middle). Arbitrary input sequences (top right) can then be efficiently queried against the index, returning a ranked list of matching paths in the graph together with metadata and a score indicating the percentage of *k*-mer identity (bottom right). The geo-information of each sample is used to highlight the locations of samples that contain sequences identical or close to the queried sequence (middle right). B) Classification accuracy of a random forest model for assigning city labels to samples as a function of the size of training set. C) Distribution of Endemicity scores (term frequency inverse document frequency) for taxa in each region. D) Prediction accuracy of a random forest model for a given feature (rows) in samples from a city (columns) that was not present in the training set. Rows and columns sorted by average accuracy. Continuous features (e.g. Population) were discretized.

283 number close to 0 (note that ES is not bounded in $[0, 1]$). Some cities, like Offa (Nigeria), host many
284 unique taxa while others, like Zurich (Switzerland), host fewer endemic species (Supp. Figure S3B).
285 Large numbers of endemic species in a city may reflect geographic bias in sampling. However, some
286 cities from well sampled continents (e.g., Lisbon, Hong Kong) also host many endemic species which
287 would suggest that ES may indicate interchangeability and local pockets of microbiome variation for
288 some locations.

289 2.3 Antimicrobial Resistance Genes Form Distinct Clusters

290 Quantification of antimicrobial diversity and AMRs are key components of global antibiotic stewardship.
291 Yet, predicting antibiotic resistance from genetic sequences alone is challenging, and detection accuracy
292 depends on the class of antibiotics (i.e., some AMR genes are associated to main metabolic pathways
293 while others are uniquely used to metabolize antibiotics). As a first step towards a global survey of
294 antibiotic resistance in urban environments, we mapped reads to known antibiotic resistance genes,
295 using the MegaRES ontology and alignment software. We quantified their relative abundance using
296 reads/kilobase/million mapped reads (RPKM) for 20 classes of antibiotic resistance genes detected in
297 our samples (Figure 4A B). 2,210 samples had some sequence which were identified as belonging to an
298 AMR gene, but no consistent core set of genes was identified. The most common classes of antibiotic
299 resistance genes were for macrolides, lincosamides, streptogamines (MLS), and betalactams, yet the most
300 common class of antibiotic resistance genes, MLS was found in only 56% of the samples where AMR
301 sequence was identified.

302 Despite being relatively common, antibiotic resistance genes were universally in low abundance com-
303 pared to functional genes, with RPKM values for resistance classes typically ranging from 0.1 – 1 com-
304 pared to values of 10 - 100 for typical housekeeping genes (AMR classes contain many genes so RPKM
305 values may be lower than they would be for individual genes). In spite of the low abundance of the genes
306 themselves, some samples contained sequences from hundreds of distinct AMR genes. Clusters of high
307 AMR diversity were not evenly distributed across cities (Figure 4C). Some cities had more resistance
308 genes identified on average (15-20X) than others (e.g. Bogota) while other cities had bimodal distribu-
309 tions (e.g. San Francisco) where some samples had hundreds of genes while others very few. We note
310 that 99% of the cases where we detected an AMR genes had an average depth of 2.7x, indicating that
311 our global distribution would not dramatically change with altered read depth (Supp. Figure S6E).

312 As with taxa, AMR genes can be used to classify samples to cities - albeit with much less accuracy.
313 A random forest model analogous to the one trained to predict city classification from taxonomic profiles
314 was trained to predict from profiles of antimicrobial resistance genes. This model achieved 37.6% accuracy
315 on held out test data (Supp. Figure S6A). While poor for actual classification this accuracy far exceeds
316 the 7.01% that would be achieved by randomly assigning labels and indicates that there are possibly
317 weak, city specific signatures for antimicrobial resistance genes.

318 Multiple AMR genes can be carried on a single plasmid and ecological competition may cause mul-
319 tiple taxa in the same sample to develop antimicrobial resistance. As a preliminary analysis into these
320 phenomena we identified clusters of AMR genes that co-occurred in the same samples (Figure 4D).
321 We measured the Jaccard distance between all pairs of AMR genes found in at least 1% of samples and
322 performed agglomerative clustering on the resulting distance matrix. We identified three large clusters of
323 genes and numerous smaller clusters. Of note, these clusters often consist of genes from multiple classes
324 of resistance. At this point we do not posit a specific ecological mechanism for this co-occurrence, but
325 we note that the large clusters contain far more genes than are typically found on plasmids.

326 We performed a rarefaction analysis on the set of all resistance genes in the dataset, which we call
327 the “panresistome” (Figure (Supp. Figure S6B). Similar to the rate of detected species, the panresistome
328 also shows an open slope with an expected rate of discovery of 1 previously unobserved AMR gene per
329 10 samples. Given that AMR gene databases are rapidly expanding and that no AMR genes were found
330 in some samples, it is likely that future analyses will identify many more resistance genes in this data.

331 Additionally, AMR genes show a “neighbourhood” effect within samples that are geographically prox-
332 imal analogous to the effect seen for taxonomic composition (Supp. Figure S6C). Excluding samples
333 where no AMR genes were detected, the Jaccard distance between sets of AMR genes increases with
334 distance for pairs of samples in the same city. As with taxonomic composition. the overall effect is weak
335 and noisy, but significant.

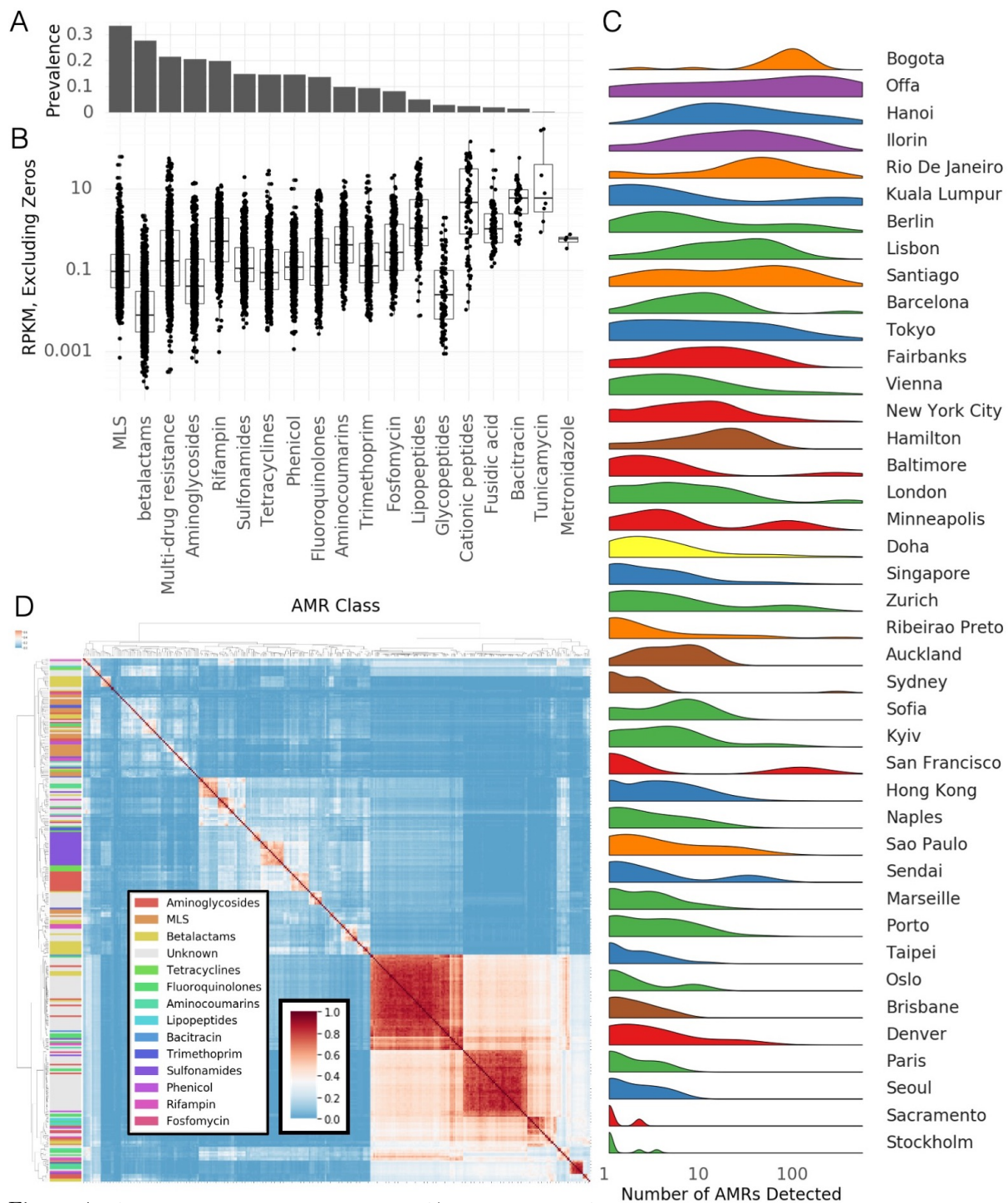


Figure 4: Antimicrobial Resistance Genes. A) Prevalence of AMR genes with resistance to particular drug classes. B) Abundance of AMR gene classes when detected, by drug class. C) Number of detected AMR genes by city. D) Co-occurrence of AMR genes in samples (Jaccard index) annotated by drug class.

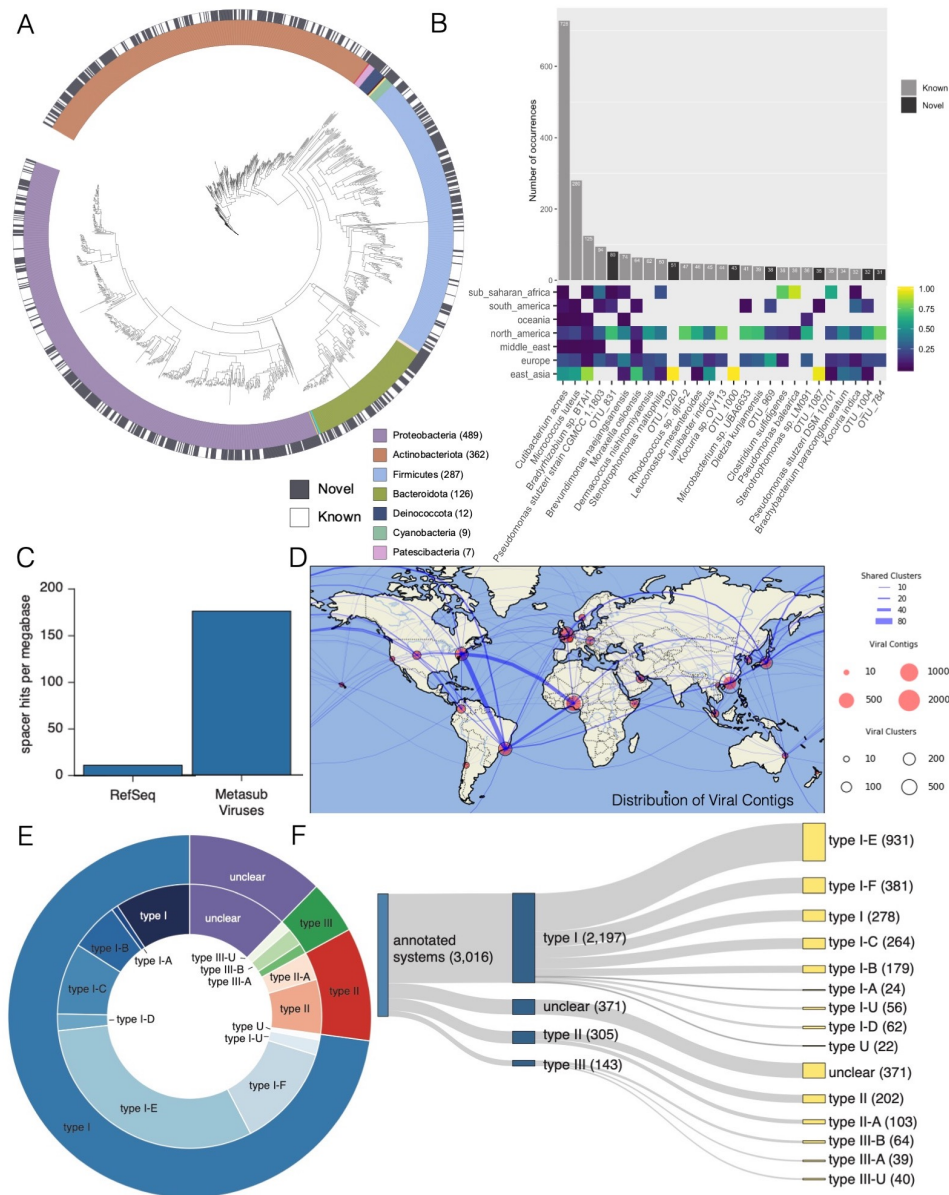


Figure 5: Novel Biology A) Taxonomic tree for Metagenome Assembled Genomes (MAGs) found in the MetaSUB data. Outer black and white ring indicates if the MAG matches a known species, inner ring indicates phyla of the MAG. B) Top: the number of samples where the most prevalent MAGs were found. Bottom: The regional breakdown of samples where the MAG was found. C) Mapping rate of CRISPR Spacers from MetaSUB data to viral genomes in RefSeq and viral genomes found in MetaSUB data. D) Geographic distribution of viral genomes found in MetaSUB data. E & F) Fractional breakdowns of identifiable CRISPR systems found in the MetaSUB data

336 2.4 Widespread Discovery of Novel Biology

337 To examine these samples for novel genetic elements, we assembled and identified Metagenome Assembled
338 Genomes (MAGs) for viruses, bacteria, and archaea and analyzed them with several algorithms. This
339 includes thousands of novel CRISPR arrays that reflect the microbial biology of the cities and 1,304
340 genomes from our data, of which 748 did not match any known reference genome within 95% average
341 nucleotide identity (ANI). 1302 of the genomes were classified as bacteria, and 2 as archaea. Bacterial
342 genomes came predominantly from four phyla: the Proteobacteria, Actinobacteria, Firmicutes, and
343 Bacteroidota. Novel bacteria were evenly spread across phyla (Figure 5A).

344 Assembled bacterial genomes were often identified in multiple samples. Several of the most prevalent
345 bacterial genomes were novel species (Figure 5B). Some assembled genomes, both novel and not, showed
346 regional specificity while others were globally distributed. The taxonomic composition of identifiable
347 genomes roughly matched the composition of the core urban microbiome (Section 2). The number of
348 identified bacterial MAGs was somewhat based on read depth and the sample count per city (Supp.
349 Figure S7A). The number of bacterial MAGs discovered in a city which did not match a known species
350 was closely correlated to the total number of bacterial MAGs discovered in that city (Supp. Figure S7B).
351 Bacterial MAGs were roughly evenly distributed geographically with the notable exception of Offa, which
352 had dramatically more novel bacterial species than other cities.

353 We investigated assembled contigs from our samples to identify 16,584 predicted uncultivated viral
354 genomes (UViGs). Taxonomic analysis of predicted UViGs to identify viral species yielded 2,009 clusters
355 containing a total of 6,979 UViGs and 9,605 singleton UViGs for a total of 11,614 predicted viral species.
356 Predicted viral species (Section ??) from samples collected within 10, 100 and 1000 kilometers of one
357 another were agglomerated to examine their planetary distribution at different scales (Figure ??C). At
358 any scale, most viral clusters appear to be weakly cosmopolitan; the majority of their members are found
359 at or near one location, with a few exceptions.

360 We compared the predicted species to known viral sequences in the JGI IMG/VR system, which
361 contains viral genomes from isolates, a curated set of prophages and 730k viral MAGs from other studies.
362 Of the 11,614 species discovered in our data 94.1% did not match any viral sequence in IMG/VR ([Paez-
363 Espino et al., 2019](#)) at the species level for a total of 10,928 novel viruses. We note that this number is
364 surprisingly high but was obtained using a conservative pipeline (99.6% precision) and corresponded well
365 with our identified CRISPR arrays (below). This suggests that urban microbiomes contain significant
366 diversity not observed in other environments.

367 Next, we attempted to identify possible bacterial and eukaryotic hosts for our predicted viral MAGs.
368 For the 686 species with similar sequences in IMG/VR, we projected known host information onto 2,064
369 MetaSUB viral MAGs. Additionally, we used CRISPR-Cas spacer matches in the IMG/M system to
370 assign possible hosts to a further 1,915 predicted viral species. Finally, we used a database of 20 million
371 metagenome-derived CRISPR spacers to provide further rough taxonomic assignments. Our predicted
372 viral hosts aligned with our taxonomic profiles, 41% of species in the core microbiome (Section 2) had
373 predicted viral-host interactions. Many of our viral MAGs were found in multiple locations (Figure 5D).
374 Many viruses were found in South America, North America and Africa. Viral MAGs in Japan often
375 corresponded to those in Europe and North America.

376 We identified 838,532 CRISPR arrays in our data of which 3,245 could be annotated for specific
377 systems. The annotated CRISPR arrays were principally type 1-E and 1-F but a number of type two
378 and three systems were identified as well (Figure 5E, F). A number of arrays had unclear or ambiguous
379 type assignment. Critically the spacers in our identified CRISPR arrays closely matched our predicted
380 viral MAGs. We aligned spacers to both our viral MAGs and all viral sequences in RefSeq. The total
381 fraction of spacers which could be mapped to our viral MAGS and RefSeq was similar (Supp. Figure
382 S7C) but the mapping rate to our viral MAGs dramatically exceeded the mapping rate to RefSeq (Figure
383 5C). We present this as additional evidence supporting these novel viral MAGs.

384 3 Discussion

385 MetaSUB is a global network of scientists and clinicians developing knowledge of urban microbiomes by
386 studying mass transit systems and hospitals within and between cities. We collected and sequenced 4,728
387 samples from 60 cities worldwide (Tables 1 and S1), constituting the first large scale metagenomic study
388 of the urban microbiome. We also identified species that are geographically constrained and showed that
389 these can be used to determine a samples city of origin (Section 2.1). Many of these species are associated
390 with commensal microbiomes from human skin and airways, but we observed that urban microbiomes are

391 nevertheless distinct from both human and soil microbiomes. Notably, no species from the *Bacteroidetes*,
392 a prominent group of human commensal organisms (Eckburg et al., 2005; Qin et al., 2010), was identified
393 in the core urban microbiome. We conclude that there is a consistent urban microbiome core (Figure
394 1, 2), which is supplemented by geographic variation (Figure 2) and microbial signatures based on the
395 specific attributes of a city (Figure 3). Our data also indicates that significant diversity remains to be
396 characterized and that novel taxa may be discovered in the data (Figure 5), that environmental factors
397 affect variation, and that sequences associated with AMR are globally widespread but not necessarily
398 abundant (Figure 4). In addition to these results, we present several ways to access and analyze our
399 data including interactive web based visualizations, search tools over raw sequence data, and high level
400 interfaces to computationally access results.

401 Unique taxonomic composition and association with covariates specific to the urban environment
402 suggest that urban microbiomes should be treated as ecologically distinct from both surrounding soil
403 microbiomes and human commensal microbiomes. Though these microbiomes undoubtedly interact
404 with the urban environment, they nonetheless represent distinct ecological niches with different genetic
405 profiles. While our metadata covariates were associated with the principal variation in our samples, they
406 do not explain a large proportion of the observed variance. It remains to be determined whether variation
407 is essentially a stochastic process or if a deeper analysis of our covariates proves more fruitful. We have
408 observed that less important principal components (roughly PCs 10-100) are generally less associated
409 with metadata covariates but that PCs 1-3 do not adequately describe the data alone. This is a pattern
410 that was observed in the human microbiome project as well, where minor PCs (such as our Figure 2B)
411 were required to separate samples from closely related body sites.

412 Much of the urban microbiome likely represents novel diversity as our samples contain a significant
413 proportion of unclassified DNA. This finding is comparable to many other metagenomic and microbiome
414 studies including other work done in subway environments (Afshinnekoo et al., 2015; Hsu et al., 2016),
415 airborne microbiomes (Yooseph et al., 2013), work done by the Earth Microbiome Project (Thompson
416 et al., 2017), and others. As noted in in Section ?? more sensitive methodology only marginally increases
417 the proportion of DNA that can be classified. We consider the DNA which would not be classified by
418 a sensitive technique to be true unclassified DNA and postulate that it may derive from novel genes or
419 species. Given that our samples did not closely resemble human commensal microbiomes or soil samples,
420 it is possible this represents novel urban DNA sequences.

421 Additionally, our discovery of a large number of novel viral sequences in our data suggests that there
422 are likely to be additional novel taxa from other domains. The fraction of predicted viral sequences which
423 belonged to previously unobserved taxa was particularly high in our study (94.1%) however taxonomic
424 associations of these viruses to observed microbial hosts suggests these results are not spurious. This
425 rate of discovery may prove prescient for novel taxa in other domains, and novel discovery of taxa may
426 help to reduce the large fraction of DNA which cannot currently be classified.

427 Many of the identified taxa are frequently implicated as infectious agents in a clinical setting including
428 specific *Staphylococcus*, *Streptococcus*, *Corynebacterium*, *Klebsiella* and *Enterobacter* species. There is
429 no clear indication that these species identified in the urban environment are pathogenic, and further in-
430 depth study is necessary to determine the clinical impact of urban microbiomes. This includes microbial
431 culture studies, specifically searching for virulence factors and performing strain-level characterization.
432 Seasonal variation also remains open to study as the majority of the samples collected here were from two
433 global City Sampling Days (June 21, 2016 and 2017). Further studies, some generating novel data, will
434 need to explore whether the core microbiome shifts over the course of the year, with particular interest
435 in the role of the microbiome in flu transmission (Cáliz et al., 2018; Korownyk et al., 2018).

436 As metagenomics and next-generation sequencing becomes more and more available for clinical (Wil-
437 son et al., 2019) and municipal use (Hendriksen et al., 2019), it is essential to contextualize the AMR
438 markers or presence of new species and strains within a global and longitudinal context. The most
439 common AMR genes were found for two classes of antibiotic: MLS and beta-lactams. MLS represents
440 macrolides, lincosamides and streptogramins, which are three groups of antibiotics with a mechanism
441 of action of inhibiting bacterial protein synthesis. Macrolides, with strong Gram-positive and limited
442 Gram-negative coverage, are prevalently used to treat upper respiratory, skin, soft tissue and sexually
443 transmitted infections amongst others. Beta-lactam antibiotics are a major class of antibiotics including
444 penicillins, cephalosporins, monobactams, carbapenems and carbacephems that are all used to treat a
445 wide array of infections. Antimicrobial resistance has surged due to the selection pressure of widespread
446 use of antibiotics and is now a global health issue plaguing communities and hospitals worldwide. Antimi-
447 crobrial resistance genes are thought to spread from a variety of sources including hospitals, agriculture
448 and water (Bougnom and Piddock, 2017; Klein et al., 2018). The antimicrobial classes particularly

449 impacted by resistance include beta-lactamases, glycopeptides and fluoroquinolones (Rice, 2012), all of
450 which we found antimicrobial resistance genes for across our samples. We found that there was uneven
451 distribution of AMR genes across cities. This could be the result of some of combination of different
452 levels of antibiotic use, differences in the urban geography between cities (population density, presence
453 of untreated wastewater etc), or reflect the background microbiome in different places in the world.
454 Techniques to estimate antibiotic resistance from sequencing data remain an area of intense research as
455 certain classes of AMR gene (ie. fluoroquinolones) are sensitive to small mutations and it is possible that
456 our methods may not fully reflect true resistance. Further research is needed to fully explore AMR genes
457 in the urban environment, including culture studies which directly measure the phenotype of resistance.

458 One of the challenges in the field of metagenomics of the built environment is dealing with low
459 biomass samples. Not only does it introduce the challenge of contamination (Kim et al., 2017) which
460 requires standardized sample preparation and the use of positive and negative controls, but there is
461 also the challenge in biases and data interpretation (McLaren et al., 2019). Metagenomic studies rely
462 on bioinformatics analyses that predict relative abundances of taxa, functional genes, antimicrobial
463 resistance genes, etc. When you have low biomass samples, these relative abundances may appear high
464 when their absolute abundance is in fact low when considering where the samples came from. However,
465 this is an inherent component of metagenomics that studies and examines microbiomes and communities
466 based on the metrics and measurements of relative abundances. There are important considerations to
467 be made from sample collection to bioinformatics analysis to ensure limited biases are introduced to a
468 study (McLaren et al., 2019). Moreover, the overall findings must be interpreted with the proper context
469 and scope of the experiment and samples collected.

470 In summary, this study presents a first molecular atlas of urban and mass-transit metagenomics from
471 across the world. By facilitating large scale epidemiological comparisons, it is a first critical step to-
472 wards quantifying the clinical role of environmental microbiomes and provides requisite data for tracking
473 changes in ecology or virulence. Moreover, in order to study the transmission of AMRs on a global scales
474 this dataset represents only focuses on some of the sources and vectors of the built environment. Indeed,
475 datasets from rural and suburban areas with livestock and farms, sewage from cities (Fresia et al., 2019;
476 Joseph et al., 2019), and other notable sources of AMRs need to be integrated together to truly capture
477 AMR mechanisms at the global scale (Singer et al., 2016; Thanner et al., 2016). Previous studies have
478 already demonstrated a role for precision clinical metagenomics in managing infectious disease and global
479 health (Afshinnkoo et al., 2017; Gardy and Loman, 2018; Ladner et al., 2019). As demonstrated by the
480 coronavirus disease 2019 (COVID-19) pandemic, as an atlas this data has the potential to aid physicians,
481 public health officers, government officials, and others in tracing, diagnosis, clinical decision making, and
482 policy within their communities.

483 3.1 Open Science

484 The MetaSUB dataset is built and organized for full accessibility to other researchers. This is consistent
485 with the concept of Open Science. Specifically, we built our study with the FAIR principles in mind:
486 Findable, Accessible, Interoperable and Reusable.

487 To make our study reproducible, we released an open source version-controlled pipeline called the
488 MetaSUB Core Analysis Pipeline (CAP). The CAP is intended to improve the reproducibility of our
489 findings by making it easy to apply a number of analyses consistently to a large dataset. This pipeline
490 includes all steps from extracting data from raw sequence data to producing refined results like taxonomic
491 and functional profiles. The CAP itself is principally composed of other open peer-reviewed scientific
492 tools, with only a few custom scripts for mundane tasks. Every tool in the CAP is open source with a
493 permissive license. The CAP is available as a docker container for easier installation in some instances
494 and all databases used in the CAP are available for public download. The CAP is versioned and includes
495 all necessary databases allowing researchers to replicate results. The CAP is not designed to produce
496 highly novel results but is meant to be a good practice agglomeration of open source tools.

497 However, the output of the CAP still consists of a number of different output formats with multiple
498 files for each sample. To make our results more reproducible and accessible, we have developed a program
499 to condense the outputs of the Core Analysis Pipeline into a condensed data-packet. This data packet
500 contains results as a series of Tidy-style data tables with descriptions. The advantage of this set-up is
501 that result tables for an entire dataset can be parsed with a single command in most high level analysis
502 languages like Python and R. This package also contains Python utilities for parsing and analyzing data
503 packets which streamlines most of the boilerplate tasks of data analysis. All development of the CAP
504 and data packet builder (Capalyzer) package is open source and permissively licensed.

505 In addition to general purpose data analysis tools essentially all analysis in this paper is available
506 as a series of Jupyter notebooks. Our hope is that these notebooks allow researchers to reproduce our
507 results, build upon our results in different contexts, and better understand precisely how we arrived at
508 our conclusions. By providing the exact source used to generate our analyses and figures, we also hope
509 to be able to quickly incorporate new data or correct any mistakes that might be identified.

510 For less technical purposes, we also provide web-based interactive visualizations of our dataset (typ-
511 ically broken into city-specific groups). These visualizations are intended to provide a quick reference
512 for major results as well as an exploratory platform for generating novel hypotheses and serendipitous
513 discovery. The web platform used, MetaGenScope, is open source, permissively licensed, and can be run
514 on a moderately powerful machine (though its output relies on results from the MetaSUB CAP).

515 Our hope is that by making our dataset open and easily accessible to other researchers the scientific
516 community can more rapidly generate and test hypotheses. One of the core goals of the MetaSUB
517 consortium is to build a dataset that benefits public health. As the project develops we want to make
518 our data easy to use and access for clinicians and public health officials who may not have computational
519 or microbiological expertise. We intend to continue to build tooling that supports these goals.

520 3.2 CAMDA

521 Since 2017 MetaSUB has partnered with the Critical Assessment of Massive Data Analysis (CAMDA)
522 camda.info, a full conference track at the Intelligent Systems for Molecular Biology (ISMB) Conference.
523 At this venue a subset of the MetaSUB data were released to the CAMDA community in the form
524 of annual challenge addressing the issue of geographically locating samples: ‘The MetaSUB Inter-City
525 Challenge’ in 2017 and ‘The MetaSUB Forensics Challenge’ in 2018 and 2019. In the latter challenge
526 the MetaSUB data has been complemented by data from EMP (Thompson et al., 2017) and other
527 studies (Delgado-Baquerizo et al., 2018; Hsu et al., 2016). This Open Science approach of CAMDA
528 has generated multiple interesting results and concepts relating to urban microbiomics, resulting in
529 several publications biologydirect.biomedcentral.com/articles/collections/camdaproc as well
530 as perspective manuscript about moving towards metagenomics in the intelligence (Mason-Buck et al.,
531 2020). The partnership is continued in 2020 with ‘The Metagenomic Geolocation Challenge’ where the
532 MetaSUB data has been complemented by the climate/weather data in order to construct multi-source
533 microbiome fingerprints and predict the originating ecological niche of the sample.

534 4 Data Access

535 Raw sequencing reads from this study contain significant amounts of human DNA and cannot yet be
536 made public. However, reads with the majority of human DNA filtered and low quality bases removed are
537 available for download from Wasabi (an Amazon S3 clone) with individual URLs located here: https://github.com/MetaSUB/metasub_utils. In addition to raw reads higher level results (e.g. taxonomic
538 profiles, functional pathways, etc.) are available in the MetaSUB data packet also available for download
539 from Wasabi. For instructional purposes we also provide a simplified data packet for teaching which
540 includes balanced numbers of samples from each city and completely filled metadata tables.

541 Interactive data visualizations are available on <https://pangea.gimmebio.com/contrib/metasub>,
542 <https://www.metagenscope.com> and GeoDNA, an interface to search query DNA sequences against
543 MetaSUB samples, is available at (dnaloc.ethz.ch/). MetaSUB data may be downloaded from <https://pangea.gimmebio.com>. MetaSUB metadata is available in the data-packet, on Pangea, or may
544 be downloaded from <https://github.com/MetaSUB/MetaSUB-metadata>. Programs used for analy-
545 sis of data may be found at https://github.com/MetaSUB/MetaSUB_CAP and <https://github.com/dcdanko/capalyzer>. Jupyter notebooks used to generate the figures and statistics in this study can be
546 found at https://www.github.com/MetaSUB/main_paper_figures. Additional tools and resources are
547 described here https://github.com/MetaSUB/bioinformatics_management.

551 5 Acknowledgement

552 DCD was supported by the Tri-Institutional Training Program in Computational Biology and Medicine
553 (CBM) funded by the NIH grant 1T32GM083937.

554 We thank GitHub for providing private repositories to the MetaSUB consortium at no cost.

555 We thank XSEDE and Philip Blood for their support of this project.

556 We would like to thank the Epigenomics and Genomics Core Facilities at Weill Cornell Medicine, fund-
557 ing from the Irma T. Hirsch and Monique Weill-Caulier Charitable Trusts, Bert L and N Kuggie Vallee
558 Foundation, the WorldQuant Foundation, Igor Tulchinsky, The Pershing Square Sohn Cancer Research
559 Alliance, NASA (NNX14AH50G, NNX17AB26G), the National Institutes of Health (R01ES021006,
560 R25EB020393, 1R21AI129851, 1R01MH117406), TRISH (NNX16AO69A:0107, NNX16AO69A:0061),
561 the NSF (1840275), the Bill and Melinda Gates Foundation (OPP1151054) and the Alfred P. Sloan Foun-
562 dation (G-2015-13964), Swiss National Science Foundation grant #407540_167331 “Scalable Genome
563 Graph Data Structures for Metagenomics and Genome Annotation” as part of Swiss National Research
564 Programme (NRP) 75 “Big Data”

565 Discovery of novel viral sequences was work conducted by the US Department of Energy Joint Genome
566 Institute, a DOE Office of Science User Facility, under contract number DE-AC02-05CH11231 and used
567 resources of the National Energy Research Scientific Computing Center, supported by the Office of
568 Science of the US Department of Energy.

569 MetaSUB Sweden was supported by Stockholm Health Authority (Region Stockholm) grant SLL
570 20160933 awarded to KIU.

571 MetaSUB Seoul was supported by the Institut Pasteur Korea (2015MetaSUB) and a National Re-
572 search Foundation of Korea (NRF) grant (NRF-2014K1A4A7A01074645, 2017M3A9G6068246).

573 Metasub Chile was supported by funding from CONICYT Fondecyt Iniciación grant 11140666 and
574 11160905, as well as funding from the Millennium Science Initiative of the Ministry of Economy, Devel-
575 opment and Tourism, Government of Chile.

576 MetaSUB Japan was supported by research funds from Keio University, the Yamagata prefectural
577 government and the City of Tsuruoka.

578 MetaSUB Austria and Ukraine acknowledge the bilateral AT-UA collaboration fund (WTZ:UA
579 02/2019; Ministry of Education and Science of Ukraine, UA:M/84-2019).

580 MetaSUB Ukraine was supported by research funds from Kyiv Academic Univeristy, Ministry of
581 Education and Science of Ukraine grant 0118U100290. MetaSUB Ukraine would like to express gratitude
582 to Kyiv Metro for the support of sampling days.

583 MetaSUB Barcelona was supported by the Spanish Ministry of Economy and Competitiveness, ‘Cen-
584 tro de Excelencia Severo Ochoa 2013-2017, the CERCA Programme / Generalitat de Catalunya, the “la
585 Caixa” Foundation, the CRG-Novartis-Africa mobility programme 2016 and TMB Director Eladio De
586 Miguel Sainz

587 Work in Colombia was partially funded by Colciencias (project No. 639677758300).

588 Work in Sao Paulo, Brazil was partially funded by CNPq (EDN - 309973/2015-5)

589 Sampling was carried out in compliance with regulations and permissions from local authorities
590 (Azienda Napoletana Mobilitàà s.p.a. in Naples, Italy; Régie des Transports Métropolitains in Marseille,
591 France; Transmilenio and ANLA permit 1484 in Bogotá, Colombia; Nigerian Railway Corporation (NRC)
592 [Ilorin and Offa Branch] and Kwara Express Transport)

593 We thank the many volunteers who made this study possible. Sara Abdul Majid, Natasha Abdullah,
594 Ait-hamlat Adel, Nayra Aguilar Rojas, Affifah Saadah Ahmad Kassim, Faisal S Al-Quaddoomi, Alex
595 Alexiev, Muhammad Al-Fath Amran, Watson Andrew, Harilanto Andrianjakarivony, Álvaro Aranguren,
596 Carme Arnan, Freddy Asenjo, Juliette Auvinet, Nuria Aventin, Erdenetsetseg Batdelger, François Baudon,
597 Carla Bello, Médine Benchouaia, Hannah Benisty, Anne-Sophie Benoiston, Diego Benítez, Juliana Bernardes,
598 Tristan Bitard-Feildel, Lucie Bittner, Guillaume Blanc, Julia Boeri, Kevin Bolzli, Alexia Bordigoni, Ciro
599 Borrelli, Sonia Bouchard, Jean-Pierre Bouly, Alessandra Breschi, Alan Briones, Aszia Burrell, Alina Bu-
600 tova, Dayana Calderon, Angela Cantillo, Miguel Carbajo, Katerine Carrillo, Laurie Casalot, Sofia Castro,
601 Jasna Chalangal, Starr Chatziefthimiou, Francisco Chavez, Allaeddine Chettouh, Erika Cifuentes, Sylvie
602 Collin, Romain Conte, Flavia Corsi, Cecilia N Cossio, Bruno D’Alessandro, Ophélie Da Silva, Katherine
603 E Dahlhausen, Natalie R Davidson, Eleonora De Lazzari, Stéphane Delmas, Chloé Dequeker, Alexandre
604 Desert, Clara N. Dias, Valeriia Dotsenko, Cassie L Ettinger, Emile Faure, Fazlina Fauzi, Aubin Fleiss,
605 Juan Carlos Forero, Mathilde Garcia, Catalina Garcia, Sonia L Ghose, Liliana Godoy, Andrea Gon-
606 zalez, Camila Gonzalez-Poblete, Charlotte Greselle, Sophie Guasco, Nika Gurianova, Sebastien Halary,
607 Eric Helfrich, Aliaksei Holik, Chiaki Homma, Michael Huber, Stephanie Hyland, Andrea Hässig, Roland
608 Häusler, Nathalie Hüsser, Badamnyambuu Iderzorig, Mizuki Igarashi, Shino Ishikawa, Sakura Ishizuka,
609 Kohei Ito, Sota Ito, Tomoki Iwashiro, Marisano James, Marianne Jaubert, Marie-Laure Jerier, Guilla-
610 uame Jospin, Nao Kato, Inderjit Kaur, Akash Keluth Chavan, Mahshid Khavari, Maryna Korshevniuk,
611 Jonas Krebs, Andrii Kuklin, Antonietta La Stora, Juliana Lago, Elodie Laine, Olha Lakhneko, Ger-
612 ardo de Lamotte, Romain Lannes, Madeline Leahy, Vincent Lemaire, Dagmara Lewandowska, Manon
613 Loubens, Olexandr Lykhenko, Salah Mahmoud, Nataalka Makogon, Dimitri Manoir, German Marchan-

614 don, Natalia Marciniak, Vincent Matthys, Arif Asyraf Md Supie, Irène Mauricette Mendy, Roy Meoded,
615 Mathilde Mignotte, Ryusei Miura, Kunihiko Miyake, Maria D Moccia, Mauricio Moldes, Jennifer Mo-
616 linet, Orgil-Erdene Molomjamts, Mario Moreno, Maureen Muscat, Cristina Muñoz, Francesca Nadalin,
617 Dorottya Nagy-Szkal, Ashanti Narce, Hiba Naveed, Thomas Neff, Wan Chiew Ng, Elsy Ngwa, Agier
618 Nicolas, Pierre Nicolas, Abdollahi Nika, Javier Quilez Oliete, Nils Ordioni, Mitsuki Ota, Francesco Oteri,
619 Yuya Oto, Coral Pardo-Este, Young-Ja Park, Jananan Pathmanathan, Manuel Perez, Melissa P Pizzi,
620 María Gabriela Portilla, Leonardo Posada, Catherine E. Pugh, Kyrylo Pyrshev, Sreya Ray Chaudhuri,
621 Hubert Rehrauer, Renee Richer, Paula Rodríguez, Paul Roldán, Sandra Roth, Maria Ruiz, Mariia Ry-
622 bak, Ikuto Saito, Yoshitaka Saito, Khaliun Sanchir, Kai Sasaki, Kaisei Sato, Masaki Sato, Ryo Sato,
623 Seisuke Sato, Yuma Sato, Oli Schacher, Christian Schori, Felipe Sepulveda, Juan C Severyn, Sarah
624 Shalaby, Hikaru Shirahata, Jordana M Silva, Gwenola Simon, Kasia Sluzek, Rebecca Smith, Yuya Sono-
625 hara, Nicolas Sprinsky, Stefan G Stark, Chisato Suzuki, Sora Takagi, Kou Takahashi, Naoya Takahashi,
626 Tomoki Takeda, Soma Tanaka, Andrea Tassinari, Eunice Thambiraja, Antonin Thiébaud, Takumi To-
627 gashi, Yuto Togashi, Anna Tomaselli, Itsuki Tomita, Nora C Toussaint, Takafumi Tsurumaki, Yelyzaveta
628 Tymoshenko, Mariko Usui, Sophie Vacant, Laura E Vann, Jhovana L Velasco Flores, Fabienne Velter,
629 Riccardo Vicedomini, Tomoro Warashina, Ayuki Watanabe, Tina Wunderlin, Olena Yemets, Tetiana
630 Yeskova, Shusei Yoshikawa, Stas Zubenko.

631 **6 Methods**

632 **6.1 Metadata Collection and Cleaning**

633 Metadata from individual cities was collected from a standardized form and set of fields. The principle
634 fields collected were the location of sampling, the material being sampled, the type of object being
635 sampled, the elevation above or below ground, and the station or line where the sample was collected.
636 However, several cities were unable to use the provided apps for various reasons and submitted their
637 metadata as separate spreadsheets. Additionally, certain metadata features, such as those related to
638 sequencing and quality control, were added after initial sample collection.

639 To collate various metadata sources, we built a publicly available program which assembled a large
640 master spreadsheet with consistent sample UUIDs. After assembling the originally collected data at-
641 tributes we added normalized attributes based on the original metadata to account for surface material,
642 control status, and features of individual cities. A full description of ontologies used is provided as part
643 of the collating program.

644 **6.2 Sample Collection and Preparation**

645 To obtain a comprehensive picture of microbial communities within a sample it is essential to choose
646 a sampling method which absorbs and preserves biological materials during sampling, transport and
647 storage until DNA extraction. The effectiveness of a swab may be influenced by a number of factors,
648 including most importantly the material of the swab tip affecting the rate at which bacteria are absorbed
649 during the sampling process. Furthermore, the design of the transport tube and DNA preserving liquids
650 affect the integrity of the material during transport. Finally, the amount of background contamination
651 identified for different products should be taken into account.

652 **6.3 Swab Comparisons**

653 In this study we have benchmarked various types of swabs and DNA preservative tubes, including Copan
654 Liquid Amies Elution Swab (ESwab, Copan Diagnostics, Cat.:480C) referred to as 'copan swab' and
655 Isohelix Swabs (Mini-Swab, Isohelix Cat.:MS-02) referred to as 'isohelix swabs', which were combined
656 with 2D Thermo Scientific™ Matrix™ storage tubes (3741-WP1D-BR/Matrix 1.0 ml/EA) referred to
657 as 'matrix tube', which have been prefilled with the preservative liquid Zymo Shield Zymo DNA/RNA
658 Shield™ (R1100-250) referred to as 'Zymo shield'. Copan swabs contain a transport medium for sample
659 preservation. After samples were collected with Copan swabs they were transported at room temperature
660 and stored at -80C until DNA extraction. Isohelix swabs have been stored in matrix tubes containing
661 400µl Zymo shield preservative. Matrix tubes were also transported at room temperature and stored at
662 -80C until DNA extraction. We tested the absorption strength of Copan and Isohelix swabs for various
663 biological and surface materials encountered when sampling subway stations. For a designated sampling
664 area of an office desk, a Isohelix swabs were moistened by submerging the swab for a few seconds in

665 preservative media. The desk area was then swabbed for 3 min. Results were compared to sampling
666 with copan swabs, which were similarly used to swab the area for 3 min.

667 6.4 Sampling Protocol

668 A standard operating procedure (SOP) was developed for the sample collection to be followed by all
669 members of the MetaSUB consortium participating in CSD. This protocol was adapted from work by
670 Afshinnekoo et al. (2015). The goal was to standardize as much of the sampling procedure and ensure
671 high quality control across the various cities and sampling teams. Thus it was recommended that teams
672 collect samples from surfaces that are present throughout most subway and transit stations and systems
673 around the world. These included ticket kiosks, turnstiles, railings, seats or benches, etc. Some cities
674 had to adapt the SOP according to their city especially if they did not have a subway system and were
675 collecting samples from other transit systems. However, the vast majority of sampling teams collected
676 samples from these surfaces. Moreover, a significant amount of metadata was recorded throughout sample
677 collections to ensure as much information regarding the samples was captured. All cities also developed
678 sampling plans for their collections and submitted them for review to have swabs sent to them, this was
679 to ensure consistency across the various sites.

680 All principal investigators and MetaSUB city leaders were trained in the sampling instructions and
681 this training was further disseminated to the respective sampling teams to ensure consistent and quality
682 control sampling. Swabbers were instructed to put on gloves before each sample collection. The swab
683 was dipped in the preservative medium to be pre-moistened before collection and sampling was timed to
684 3 minutes to ensure highest yield. Other key points in training included ensure highest surface area was
685 used for collection (i.e. swab entire bench, not just one area) and avoiding any areas that appeared wet,
686 contaminated, and not consistent with a subway surface. Any other observations or important notes
687 during sample collection that could add more context to data analysis and interpretation were recorded
688 on the notes section of the metadata collection apps.

689 There were some changes between CSD2016 and CSD2017 sampling protocols that are important to
690 note. First, the swab was changed from Copan (CSD16) to Isohelix (CSD17) this was after the results of
691 benchmarking work comparing the swabs and ensuring we are optimizing the amount of DNA collected
692 from swabbing these surfaces. Moreover a barcoding system was set in place in CSD17 to improve
693 metadata collection and sample tracking compared to the CSD ID system utilized in CSD16 collection
694 (CSD-City Code-00XYZ).

695 6.4.1 In-Lab controls CSD2016

696 As positive lab control we used 30 μ l ZymoBiOMICS Microbial Community standard (Catalog #D6300),
697 which we added to an empty sterile urine cap, followed by swabbing with Copan Liquid Amies Elution
698 Swab (ESwab, Copan Diagnostics, Cat.:480C) for 1.5min / 3 minutes. As negative (background) lab
699 control we used 50 μ l of the final resuspension buffer (MoBio PowerSoil®DNA Isolation Kit, Cat.:12888-
700 100), which we have added to an empty sterile urine cup followed by swabbing for 3 min (Fig.S1).
701 Furthermore, the working space has been swabbed for 1.5 min / 3 min before and after treatment with
702 10% bleach (Fig. S2) to test for background contamination rates. To identify the background levels of
703 biological material in the air at sample areas, a Copan swab has been held for 1.5 min - 3 min in the
704 air. To estimate the source and amount of contamination in commercial swab and tube products used
705 for MetaSUB, we tested all consumables in triplicates in the sterilized hood (UV light and 10% bleach
706 wiped with ethanol).

707 6.4.2 DNA Extraction from Isohelix swabs using ZymoBiomics 96 MagBead

708 The Isohelix swab head and the entire 400 μ l of DNA/RNA Shield-solubilized sample were transferred
709 into ZR BashingBead Lysis Tubes (0.1 & 0.5 mm) (Cat# S6012-50) to which an additional 600 μ l of
710 DNA/RNA Shield was added. Mechanical lysis using bead beating was performed on a maximum of 18
711 samples simultaneously using the Scientific Industries Vortex-Genie 2 with Horizontal-(24) Microtube
712 Adapter (Cat # SI-0236 and SI-H524) at maximum power for 40 minutes. The resulting lysate (400 μ l)
713 was transferred to NuncTM 96-Well Polypropylene DeepWell Storage Plates (Cat # 278743), followed
714 by DNA extraction using the ZymoBIOMICS 96 MagBead Kit (Lysis Tubes) (Catalog # D4308) on the
715 Hamilton Star according to manufacturer instructions.

716 6.4.3 DNA extraction from Copan swabs using MoBio PowerSoil®DNA

717 Droplets in the Copan Liquid Amies Elution Swab tube (ESwab, Copan Diagnostics, Cat.:480C (<http://goo.gl/8a9uCP>)) were spun down at 300rpm/1min. Next, the swab pad was transferred to a Mo-
718 Bio PowerSoil®DNA vial containing beads using sterile scissors, which we sterilized by flaming with
719 100% ethanol. The remaining 400-500µl Copan Amies liquid has been transferred into an Eppendorf
720 tube and centrifuged at full speed to collect bacteria and debris in a pellet. The pellet was finally
721 transferred to the same MoBio PowerSoil®DNA vial also containing the corresponding swab pad. Mo-
722 Bio PowerSoil®DNA Isolation Kit, Cat.:12888-100 (<https://goo.gl/65rcn2>) was used according to
723 manufacturer's instructions except for the following modifications:
724

725 Both swab and pellet have been re-suspended with 135µl C1 buffer (MoBio PowerSoil®DNA). Sample
726 homogenization was performed using either TissueLyser II (Qiagen) with 2 cycles of 3 minutes at 30Hz
727 (<https://goo.gl/hBg8Lb>), or using the Vortex-Genie 2 (Vortex Catalog #13000-V1-24) adaptor and
728 vortex at maximum speed for 10 minutes. The sequencing centers in Stockholm and Shanghai used
729 different procedures for homogenization. Stockholm used a method based on MPI FASTPREP, while
730 Shanghai added 0.6 grams of 100-micron zirconium-silica beads to 2ml tubes containing the swab pad
731 and the media, followed by bead beating for 1 min. The eluted samples have been additionally purified
732 and concentrated by Beckmann Coulter Agencourt AMPure XP (Cat.:A63881) purification (1.8X) and
733 eluted into 12µl - 50µl elution buffer. Subsequently, DNA was quantified using Qubit® dsDNA HS
734 Assay (Catalog #Q32854).

735 6.4.4 DNA extraction using Promega Maxwell

736 We added 300µl Promega Maxwell Lysis buffer and 30µl Promega Maxwell Proteinase K to Copan swab
737 heads or Isohelix swab heads and transferred the swabs back to their respective collection tube. For lysis
738 the sample tubes containing the swabs and the lysis mixture were incubated in a water bath at 54C for
739 30min. Following lysis, Copan swab heads were cut off their stem using sterile scissors and transferred
740 into a filter tube (Promega V4745). The filter containing the swab was placed into a 2ml Eppendorf tube
741 and spun down at full speed for 2min. This step is necessary since the Copan swab material consists of a
742 foam, which harbors the main liquid containing the extracted DNA. Next, the eluate has been combined
743 with the corresponding sample tube media and added to the first well of the cartridge (Maxwell® RSC
744 Buccal Swab kit AS1640). Cartridges were processed using the Maxwell® RSC Instrument (AS4500)
745 following the manufacturer's default instructions. Extracted DNA was eluted in 50µl Promega Elution
746 Buffer and stored at -80C.

747 The matrix tubes containing the Isohelix swabs and the lysis buffer have been vortexed at full speed for
748 one minute. The Isohelix swab head material is a non-porous material, which allows for easy collection of
749 the lysate. We transferred the lysate to the first cartridge of the Maxwell® RSC Blood DNA KitAS1400
750 using syringes (BD 3 mL Syringes with 18G x 1.5" Luer Lok Tip Blunt Fill Needles) and ran the Promega
751 Maxwell using the Blood program according to manufacturer's instructions. Samples were subsequently
752 eluted in 50µl elution buffer and stored at -80C.

753 Pilot samples collected in Barcelona and Stockholm were prepared for NGS analysis using QIAGEN
754 QIAseq FX DNA Library Kit. Samples from CSD2017 and CSD2018 have been prepped at HudsonAlpha
755 Genome Center described by [Afshinnkoo et al. \(2015\)](#).

756 6.5 Quality Control

757 6.5.1 Sequencing quality

758 We measured sequencing quality based on 5 metrics: number of reads obtained from a sample, GC
759 content, Shannon's entropy of k -mers, post PCR Qubit score, and recorded DNA concentration before
760 PCR. The number of reads in each sample was counted both before and after quality control, we used
761 the number of reads after quality control for our results though the difference was slight. GC content
762 was estimated from 100,000 reads in each sample after low quality DNA and human reads had been
763 removed. Shannon's entropy of k -mers was estimated from 10,000 reads taken from each samples. PCR
764 Qubit score and DNA concentration are described in the wet lab methods.

765 6.5.2 Sequencing quality scores show expected trends

766 We measured sequencing quality based on 5 metrics: number of reads obtained from a sample, GC
767 content (taken after removing human reads), Shannon's entropy of k -mers (from 10,000 reads sampled

768 from each sample), post PCR Qubit score, and recorded DNA concentration before PCR. We observed
769 good separation of negative and positive controls based on both PCR Qubit and k -mer entropy (Supp.
770 Figure S14). Distributions of DNA concentration and the number of reads were as expected. GC content
771 was broadly distributed for negative controls while positive controls were tightly clustered, expected since
772 positive controls have a consistent taxonomic profile. Comparing the number of reads before and after
773 quality control did not reveal any major outliers.

774 6.5.3 Batch effect appears minimal

775 A major concern for this low-biomass studies and large-scale studies are batch effects. The median flowcell
776 used in our study contained samples from 3 cities and 2 continents. However, two flowcells covered 18
777 cities from 5 or 6 continents respectively. When samples from these flowcells were plotted using UMAP
778 (see Section 2.1 for details) the major global trends we described were recapitulated (Supp. Figure
779 S15A). Further, when plotting samples by PCR qubit and k -mer entropy (the two metrics that most
780 reliably separated our positive and negative controls) and overlaying the flowcell used to sequence each
781 sample only one outlier flowcell was identified and this flowcell was used to sequence a large number of
782 background control samples (Supp. Figure S15B). Plots of the number of reads against city of origin and
783 surface material (Supp. Figure S15C & D) showed a stable distribution of reads across cities. Analogous
784 plots of PCR Qubit scores were less stable than the number of reads but showed a clear drop for control
785 samples (Supp. Figure S15E & F). These results led us to conclude that batch effects are likely to be
786 minimal.

787 6.5.4 Strain Contamination

788 We used BLASTn to align nucleotide assemblies from case samples to control samples. We used a
789 threshold of 8,000 base pairs and 99.99% identity as a minimum to consider two sequences homologous.
790 This threshold was chosen to be sensitive without solely capturing conserved regions. We identified all
791 connected groups of homologous sequences and found approximate taxonomic identifications by aligning
792 contigs to NCBI-NT using BLASTn searching for 90% nucleotide identity over half the length of the
793 longest contig in each group.

794 6.5.5 Strain contamination is rare or absent

795 Despite good separation of positive and negative controls (see Section 6.5.1) we identified several species
796 in our negative controls which were also identified as prominent taxa in the data-set as a whole (See
797 Section 2). Our dilemma was that a microbial species that is common in the urban environment might
798 also reasonably be expected to be common in the lab environment. In general, negative controls had
799 lower k -mer complexity, fewer reads, and lower post PCR Qubit scores than case samples and no major
800 flowcell specific species were observed. Similarly, positive control samples were not heavily contami-
801 nated. These results suggest samples are high quality but do not systematically exclude the possibility
802 of contamination.

803 Previous studies have reported that microbial species whose relative abundance is negatively cor-
804 related with DNA concentration may be contaminants. We observed a number of species that were
805 negatively correlated with DNA concentration (Supp. Figure S13A) but this distribution followed the
806 same shape (but had a greater magnitude) as a null distribution of uniformly randomly generated rela-
807 tive abundances (Supp. Figure S13B) leading us to conclude that negative correlation may simply be a
808 statistical artifact. We also plotted correlation with DNA concentration against each species mean rela-
809 tive abundance across the entire data-set (Supp. Figure S13C). Species that were negatively correlated
810 with DNA concentration were clearly more abundant than uncorrelated species, this suggests that there
811 may be a jackpot effect for prominent species in samples with lower concentrations of DNA but is not
812 generally consistent with contamination.

813 We analyzed the total complexity of case samples in comparison to control samples. Case samples
814 had a significantly higher taxonomic diversity (Supp. Figure S12A) than any type of negative control
815 sample. We also compared the confidence of taxonomic assignments to control assignments for prominent
816 taxa (Supp. Figure S12B) using the number of unique marker k -mers to compare assignments. We found
817 that case samples had more and higher quality assignments than could be found in controls. One species,
818 *Bradyrhizobium sp. BTAi1*, was not clearly better in case samples than controls but in this case we were
819 able to assemble genomes for this species in several unique samples so we feel it is ambiguous.

820 Finally, we compared assemblies from negative controls to assemblies from our case samples searching
821 for regions of high similarity that could be from the same microbial strain. We reasoned that uncontam-
822 inated samples may contain the same species as negative controls but were less likely to contain identical
823 strains. Only 137 case samples were observed to have any sequence with high similarity to an assem-
824 bled sequence from a negative control (8,000 base pairs minimum of 99.99% identity). The identified
825 sequences were principally from *Bradyrhizobium* and *Cutibacterium*. Since these genera are core taxa
826 (See Section 2) observed in nearly every sample but high similarity was only identified in a few samples,
827 we elected not to remove species from these genera from case samples.

828 6.5.6 K-Mer Based Analyses

829 We generated 31-mer profiles for raw reads using Jellyfish. All k -mers that occurred at least twice in
830 a given sample were retained. We also generated MASH sketches from the non-human reads of each
831 sample with 10 million unique minimizers per sketch.

832 We calculated the Shannon's entropy of k -mers by sampling 31-mers from a uniform 10,000 reads per
833 sample. Shannon's entropy of taxonomic profiles was calculated using the CAPalyzer package (Section
834 4).

835 6.5.7 K-Mer based metrics correlate with taxonomic metrics

836 We found clear correlations between three pairwise distance metrics (Supp. Figure S10A, B, C): k -mer
837 based Jaccard distance (MASH), taxonomic Jaccard distance, and taxonomic Jensen-Shannon diver-
838 gence. This suggests that taxonomic variation reflects meaningful variation in the underlying sequence
839 in a sample.

840 We also compared alpha diversity metrics (Supp. Figure S10D): Shannon entropy of k -mers, and
841 Shannon entropy of taxonomic profiles. As with pairwise distances these metrics were correlated though
842 noise was present. This noise may reflect sub-species taxonomic variation in our samples.

843 6.5.8 Sequence Preprocessing

844 Sequence data were processed with AdapterRemoval (v2.17, [Schubert et al. \(2016\)](#)) to remove low quality
845 reads and reads with ambiguous bases. Subsequently reads were aligned to the human genome (hg38,
846 including alternate contigs) using Bowtie2 (v2.3.0, fast preset, [Langmead and Steven L Salzberg \(2013\)](#)).
847 Read pairs where both ends mapped to the human genome were separated from read pairs where neither
848 mate mapped. Read pairs where only one mate mapped were discarded. Hereafter, we refer to the read
849 sets as human reads and non-human reads.

850 6.5.9 Unmapped DNA is not similar to any known sequence

851 A large proportion of the reads in our samples were not mapped to any references sequences. There
852 are three major reasons why a fragment of DNA would not be classified in our analysis 1) The DNA
853 originated from a non-human and non-microbial species which would not be present in the databases
854 we used for classification 2) Our classifier (KrakenUniq) failed to classify a DNA fragment that was in
855 the database due to slight mismatch 3) The DNA fragment is novel and not represented in any existing
856 database. Explanations (1) and (2) are essentially drawbacks of the database and computational model
857 used, and we can quantify them by mapping reads using a more sensitive aligner to a larger database,
858 such as BLASTn ([Altschul et al., 1990](#)), or ensemble methods for analysis ([McIntyre et al., 2017](#)). To
859 estimate the proportion of reads which could be assigned, we took 10k read subsets from each sample
860 and mapped these to a set of large database using BLASTn (see 2 for details). This resulted in 34.6%
861 reads which could not be mapped to any external database compared to 41.3% of reads mapped using
862 our approach with KrakenUniq. We note that our approach to estimate the fraction of reads that could
863 be classified using BLASTn does not account for hits to low quality taxa which would ultimately be
864 discarded in our pipeline, and so represents a worst-case comparison. Explanation (3) is altogether more
865 interesting and we refer to this DNA as true unclassified DNA. In this analysis we do not seek to quantify
866 the origins of true unclassified DNA except to postulate that it may derive from novel species as have
867 been identified in other similar studies.

868 6.6 Computational Analysis

869 6.6.1 Taxonomic Analysis

870 We generated taxonomic profiles by processing non-human reads with KrakenUniq (v0.3.2 [Breitwieser](#)
871 [et al. \(2018\)](#)) using a database based on all draft and reference genomes in NCBI/RefSeq Microbial (bac-
872 teria/archaea, fungi and virus) ca. March 2017. KrakenUniq was selected because its high performance,
873 as it has been demonstrated to be comparable or having higher sensitivity than the best tools identified
874 in a recent benchmarking study ([McIntyre et al. \(2017\)](#)) on the same comparative dataset. In addition,
875 KrakenUniq allows for tunable specificity and identifies k -mers that are unique to particular taxa in a
876 database. Reads are broken into k -mers and searched against this database. Finally, the taxonomic
877 makeup of a sample is given by identifying the taxa with the greatest leaf to ancestor weight.

878 KrakenUniq reports the number of unique marker k -mers assigned to each taxon, as well as the total
879 number of reads, the fraction of available marker k -mers found, and the mean copy number of those
880 k -mers. We found that requiring more k -mers to identify a species resulted in a roughly linear decrease
881 in the total number of species identified without a plateau or any other clear point to set a threshold
882 (Supp. Figure S9A). In an ongoing but unpublished clinical study we have used a threshold of 512
883 marker k -mers to accurately recapitulate the results of culturing while identifying few species which were
884 not cultured. Since false positives are less problematic in the current study than in a clinical study and
885 because we could use our large number of samples as a partially orthogonal confirmation we chose less
886 strict thresholds for KrakenUniq in this study.

887 At a minimum we required three reads assigned to a taxa with 64 unique marker k -mers. This setting
888 captures a group of taxa with low abundance but reasonable (~ 10 -20%) coverage of the k -mers in their
889 marker set (Supp. Figure S9C). However, this also allows for a number of taxa with very high (105)
890 duplication of the identified marker k -mers and very few k -mers per read which we believe is biologically
891 implausible (Supp. Figure S9D). We filtered these taxa by applying a further filter which required that
892 the number of reads not exceed $\frac{10}{25}$ times the number of unique k -mers, unless the set of unique k -mers
893 was saturated ($> 90\%$ completeness). We include a full list of all taxonomic calls from all samples
894 including diagnostic values for each call. We do not attempt to classify reads below the species level in
895 this study.

896 We further evaluated prominent taxonomic classifications for sequence complexity and genome cov-
897 erage. For each microbe evaluated we calculated two indices generated using a random subset of 152
898 samples: the average topological entropy of reads assigned to the microbe and the Gini-coefficient of read
899 positions on the microbial genome. For brevity we refer to these as *mean sequence entropy* (MSE) and
900 *coverage equality* (CE). The formula for topological entropy of a DNA sequence is described by [Koslicki](#)
901 [\(2011\)](#). Values close to 0 correspond to low-complexity sequences and values near 1 are high complexity.
902 In this work we use a word size of 3 with an overall sequence length of 64 since this readily fits into
903 our reads. To find the MSE of a microbial classification we take the arithmetic mean of the topological
904 entropy of all reads that map to a given microbial genome in a sample. The Gini-coefficient is a classic
905 economic measure of income inequality. We repurpose it here to evaluate the evenness of read coverage
906 over a microbial classification. Reads mapping to a microbial genome are assigned to a contiguous 10kbp
907 bin and the Gini-coefficient of all bins is calculated. Like MSE, the Gini-coefficient is bounded in $[0, 1]$.
908 Lower values indicate greater inequality, very low values indicate that a taxon may be misidentified from
909 conserved and near conserved regions. We downloaded one representative genome per species evaluated
910 and mapped all reads from samples to using Bowtie2 (sensitive-local preset). Indices were processed
911 from alignments using a custom script. Species classifications with an average MSE less than 0.75 or CE
912 less than 0.1 were flagged.

913 To determine relative abundance of taxa where applicable we rarefied samples to 100,000 classified
914 reads, computed the proportion of reads assigned to each taxon, and took the distribution of values from
915 all samples. This was the minimum number of reads sufficient to maintain taxonomic richness (Supp.
916 Figure S9B). We chose sub-sampling (sometimes referred to as rarefaction in the literature) based on the
917 study by [Weiss et al. \(2017\)](#), showing that sub-sampling effectively estimates relative abundance. Note
918 that we use the term prevalence to describe the fraction of samples where a given taxon is found at any
919 abundance and we use the term relative abundance to describe the fraction of DNA in a sample from a
920 given taxon.

921 We compared our samples to metagenomic samples from the Human Microbiome Project and a
922 metagenomic study of European soil samples using MASH ([Ondov et al., 2016](#)), a fast k -mer based
923 comparison tool. We built MASH sketches from all samples with 10 million unique k -mers to ensure
924 a sensitive and accurate comparison. We used MASH's built-in Jaccard distance function to generate

925 distances between our samples and HMP samples. We then took the distribution of distances to each
926 particular human commensal community as a proxy for the similarity of our samples to a given human
927 body site.

928 We also compared our samples to HMP and soil samples using taxonomic profiles generated by
929 MetaPhlAn v2.0 (Segata et al., 2012). We generated taxonomic profiles from non-human reads using
930 MetaPhlAn v2.0 and found the cosine similarity between all pairs of samples.

931 We used the Microbe Directory (Shaaban et al., 2018) to annotate taxonomic calls. The Microbe
932 Directory is a hand curated, machine readable, database of functional annotations for 5,000 microbial
933 species.

934 6.6.2 Functional Analysis

935 We analyzed the metabolic functions in each of our samples by processing non-human reads with HU-
936 MAnN2 (Franzosa et al., 2018). We aligned all reads to UniRef90 using DIAMOND (v0.8.36, (Buchfink
937 et al., 2014)) and used HUMAnN2 to produce estimate of pathway abundance and completeness. We
938 filtered all pathways that were less than 50% covered in a given sample but otherwise took the reported
939 pathway abundance as is after relative abundance normalization (using HUMAnN2's attached script).

940 High level categories of functional pathways were found by grouping positively correlated pathways
941 and manually annotating resulting clusters.

942 6.7 Assembly and Plasmid Annotations

943 All samples were assembled using metaSPAdes (v3.8.1 Nurk et al. (2017)) with default settings. Assem-
944 bled scaffolds of at least 1,500bp of length were annotated using PlasFlow (v1.1 Krawczyk et al. (2018))
945 using default settings. PlasFlow predicts whether a contig is likely from a chromosome or a plasmid and
946 gives a rough taxonomic annotation. Predicting which sequences are from plasmids is a difficult problem
947 and some annotations may be incorrect.

948 6.7.1 Analysis of Antimicrobial Resistance Genes

949 We generated profiles of antimicrobial resistance genes using MegaRes (v1.0.1, Lakin et al. (2017)). To
950 generate profiles from MegaRes, we mapped non-human reads to the MegaRes database using Bowtie2
951 (v2.3.0, very-sensitive presets, Langmead and Steven L Salzberg (2013)). Subsequently, alignments
952 were analyzed using ResistomeAnalyzer (commit 15a52dd github.com/cdeanj/resistomeanalyzer)
953 and normalized by total reads per sample and gene length to give RPKMs. MegaRes includes an ontology
954 grouping resistance genes into gene classes, AMR mechanisms, and gene groups. AMR detection remains
955 a difficult problem and we note that detection of a homologous sequence to a known AMR gene does
956 not necessarily imply an equivalent resistance in our samples. Currently, the gold standard for detecting
957 AMR is via culturing.

958 Known AMR genes can come from gene families with homologous regions of sequence. To reduce
959 spurious mapping from gene homology we used BLASTn to align all MegaRes AMR genes against
960 themselves. We considered any connected group of genes with an average nucleotide identity of 80%
961 across 50% of the gene length as a set of potentially confounded genes. We collapsed all such groups
962 into a single pseudo-gene with the mean abundance of all constituent genes. Before clustering genes we
963 removed all genes which were annotated as requiring SNP verification to predict resistance.

964 In addition to MegaRes we mapped non-human reads from all samples to the amino acid gene se-
965 quences in the Comprehensive Antibiotic Resistance Database (McArthur et al., 2013) using DIAMOND.
966 While we do not use this analysis explicitly in this study we provide the results as a data table.

967 Assembled contigs were annotated for AMR genes using metaProdigal (Hyatt et al., 2010), HMMER3
968 (Eddy, 2011), and ResFam (Gibson et al., 2015) as described by Rahman et al. (2018). All predicted
969 gene annotations with an e-value higher than 10^{-10} were discarded.

970 6.7.2 Beta Diversity

971 Inter-sample (beta) diversity was measured by using Jaccard distances. We note that Jaccard distances
972 do not use relative abundance information. Matrices of Jaccard distances were produced using built in
973 SciPy functions treating all elements greater than 0 as present. Hierarchical clustering (average linkage)
974 was performed on the matrix of Jaccard distances using SciPy (<https://www.scipy.org/>).

975 Dimensionality reduction of taxonomic and functional profiles was performed using UMAP (McInnes
976 et al., 2018) on the matrix of Jaccard distances with 100 neighbours (UMAP-learn package, random
977 seed of 42). We did not use Principal Component Analysis as a preprocessing step before UMAP as it
978 is sometimes done for high dimensional data.

979 6.7.3 Alpha Diversity

980 Intra-sample (alpha) diversity was measured by using Species Richness and Shannon's Entropy. We
981 took species richness as the total number of detected species in a sample after rarefaction to 1 million
982 reads. Shannon's entropy is robust to sample read depth and accounts for the relative size of each
983 group in diversity estimation. Shannon's entropy is typically defined as $H = -\sum a_i \log_2 a_i$ where a_i is the
984 relative abundance of taxon i in the sample. For alpha diversity based on k -mers or pathways, we simply
985 substitute the relative abundance of a species for the relative abundance of the relevant type of object.

986 6.7.4 GeoDNA Sequence Search

987 For building the sequence graph index, each sample was processed with KMC (version 3, [1]) to convert
988 the reads in FASTA format into lists of k -mer counts, using different values of k ranging from 13 to 19 in
989 increments of 2. All k -mers that contained the character "N" or occurred in a sample less than twice were
990 removed. For each value of k , we built a separate index, consisting of a labeled de Bruijn graph, using an
991 implicit representation of the complete graph and a compressed label representation based on Multiary
992 Binary Relation Wavelet Trees (Multi-BRWT). For further details, we refer to the manuscript [2]. To
993 build the index, for each sample the KMC k -mer count lists were transformed into de Bruijn graphs, from
994 which path covers in the form of contig sets were extracted and stored as intermediate FASTA files. The
995 contig sets of each sample were then transformed into annotation columns (one column per sample) by
996 mapping them onto an implicit complete de Bruijn graph of order k . All annotation columns were then
997 merged into a joint annotation matrix and transformed into Multi-BRWT format. Finally, the topology
998 of the Multi-BRWT representation was optimized by relaxing its internal tree arity constraints to allow
999 for a maximum arity of 40.

1000 6.8 Novel Biology

1001 6.9 Identifying Bacteria and Archaea

1002 **Metagenomic Assembly and Binning** All samples were re-assembled with metaSPAdes (v3.10.1
1003 Nurk et al., 2017); generated contigs with length <1000nt were excluded from further analysis. Remaining
1004 contigs were binned with MetaBAT2 (v2.12.1 Kang et al. (2019)) with default parameters, resulting in
1005 14,080 bins. As MetaBAT2 uses contig abundance (mean base coverage) in its analysis, we mapped reads
1006 back to their respective contigs via Bowtie2 (v2.3.4.1 Langmead and Steven L Salzberg (2013)) with the
1007 flags `-local -very-sensitive-local` to provide accurate coverage metrics. Draft genome quality was assessed
1008 via CheckM (v1.0.13 Parks et al. (2015)) lineage_wf workflow with default parameters. Using the
1009 strategy proposed by Parks et al. (2018) we filtered bins by quality score, defined as $QS = completeness -$
1010 $5 * contamination$; bins with $QS < 50$ were removed from consideration. The remaining 6,107 bins were
1011 labeled by quality based on the MIMAG standard (Bowers et al. (2018)), with some modification: 1,448
1012 high quality (completeness >90%, contamination <5%, strain heterogeneity <0.5%) bins, 4,532 medium
1013 quality (completeness >50%, contamination <5%) bins, all others low quality. Bins of at least medium
1014 quality were selected as acceptable MAGs (5,980 total).

1015 **MAG Dereplication** OTUs (MAG representatives) were chosen with a two-step clustering strategy.
1016 Single-linkage clustering formed primary clusters of MAGs based on Mash ANI (v2.1.1), with intra-cluster
1017 identity at 90%. Though Mash ANI can be inaccurate for potentially incomplete genomes (Olm et al.
1018 (2017)), we can leverage the technique's speed for the many pairwise comparisons needed in this granular
1019 step. Within primary clusters, MAGs were compared pairwise by a more accurate whole-genome ANI
1020 (gANI) via dnadiff (v1.3) from MUMmer (v3.23 Kurtz et al. (2004)). Secondary, more refined clusters
1021 were grouped based on gANI using average-linkage hierarchical clustering from the R package dendextend
1022 (v1.12.0 Galili (2015)). A gANI cut-off of 95% resulted in 1,304 representative OTUs.

1023 **OTU to Reference Genome Matching** OTUs were compared against reference genomes from Ref-
1024 Seq (release 96 from November 2019, complete bacterial and archaeal genomes only, with “Exclude
1025 anomalous” and “Exclude derived from surveillance project” applied) as well as the full Integrated Gut
1026 Genomes (IGG) dataset (v1.0 [Nayfach et al. \(2019\)](#)); 23,790 representative genomes). A MinHash sketch
1027 was created for each reference genome via Mash (v2.1.1) with default parameters to find Mash distances
1028 and select candidate “best matches” from each reference database. Then, dnadiff (v1.3) was used to
1029 further quantify differences between each OTU and its best match from either database. ANI between
1030 OTUs and their matches was found as “M-to-M AvgIdentity” in the query report column (ANI 95% over
1031 60% OTU sequence qualified as a match).

1032 **OTU Taxonomic Assignment** OTUs were placed into a bacterial or archaeal reference tree (based
1033 on the Genome Database Taxonomy, GTDB) and then assigned taxonomic classifications using GTDB-
1034 Tk (v1.0.2 [Chaumeil et al. \(2019\)](#)). GTDB-Tk relies on 120 bacterial and 122 archaeal marker genes;
1035 domain assignment is chosen based on domain-specific marker content of the OTU sequence. Using the
1036 GTDB-Tk placements, we built an OTU-only bacterial phylogeny with FastTree (v2.1.10 [Price et al.](#)
1037 [\(2010\)](#)). The tree was visualized using iTOL (v5.5 [Letunic and Bork \(2019\)](#)).

1038 6.9.1 Viral Discovery

1039 We followed the protocol described by [Paez-Espino et al. \(2017\)](#). Briefly, we used an expanded and
1040 curated set of viral protein families (VPFs) as bait in combination with recommended filtering steps to
1041 identify 16,584 UViGs directly from all MetaSUB metagenomic assemblies greater than 5kb. Then, the
1042 UViGs were clustered with the content of the IMG/VR system (a total of over 730k viral sequences
1043 including isolate viruses, prophages, and UViGs from all kind of habitats). The clustering step relied on
1044 a sequence-based classification framework (based on 95% sequence identity across 85% of the shortest
1045 sequence length) followed by the markov clustering (mcl). This approach yielded 2,009 viral clusters
1046 (ranging from 2-611 members) and 9,605 singletons (or viral clusters of 1 member), sequences that failed
1047 to cluster with any sequence from the dataset or the references from IMG/VR, resulting in a total of
1048 11,614 vOTUs. We define viral species from vOTUs as sequences sharing at least 95% identity over 85%
1049 of their length. Out of this total MetaSUB viral diversity, only 686 vOTUs clustered with any known
1050 viral sequence in IMG/VR.

1051 6.9.2 Identifying Host Virus Interactions

1052 We used two computational methods to reveal putative host-virus connections ([Paez-Espino et al., 2016a](#)).
1053 (1) For the 686 vOTUs that clustered with viral sequences from the IMG/VR system, we projected the
1054 known host information to all the members of the group (total of 2,064 MetaSUB UViGs). (2) We used
1055 bacterial/archaeal CRISPR-Cas spacer matches (from the IMG/M 1.1 million isolate spacer database) to
1056 the UViGs (allowing only for 1 SNP over the whole spacer length) to assigned a host to 1,915 MetaSUB
1057 vOTUs. Additionally, we also used a database of over 20 million CRISPR-Cas spacers identified from
1058 metagenomic contigs from the IMG/M system with taxonomy assigned. Since some of these spacers may
1059 derive from short contigs these results should be interpreted with caution.

1060 6.9.3 CRISPR Array Detection and Annotation

1061 Using CRISPRCasFinder the MetaSUB database was investigated to predict CRISPR arrays and an-
1062 notate them with their corresponding predicted type based on CRISPR-Cas genes in their vicinity.
1063 CRISPRCasFinder was run with default parameters, “-so” and “-cas” options to identify cas genes. The
1064 precision and recall of the virus detection was 99.6% and 37.5% respectively, as previously reported by
1065 ([Paez-Espino et al., 2016b](#)).

1066 CRISPR-Cas types were assigned to arrays based on detected cas genes within a 10 kilobases vicinity.
1067 Cases where CRISPRCasFinder associated several cas genes of contradicting CRISPR-Cas types with
1068 the same CRISPR array were regarded as unclear annotation. This procedure yielded 838,532 predicted
1069 CRISPR arrays (with additional CRISPR arrays predicted with default parameters for PILER-CR), of
1070 which, 3,245 CRISPR arrays had unambiguous annotation, resulting in 43,656 unique spacers queried
1071 against genomic databases using BLASTN.

1072 **6.10 Organisms/BLAST Databases**

1073 In order to associate detected spacers within defined groups (plasmids, prophages, viruses) four different
1074 genomic databases were aggregated to be searched with BLASTN. The aggregated database consisted
1075 of IMG/VR, PHASTER, and PLSDB alongside bacterial and archaeal genomic sequences from the
1076 National Center for Biotechnology Information (NCBI). All database downloads were made on the 28th
1077 January 2020. Detected and annotated spacers were searched against the databases mentioned above
1078 using BLASTN with the following additional arguments, which correspond to the default parameters of
1079 CRISPRTarget: word_size=7, evalue=1, gapopen=10, gapextend=2, penalty=-1, reward=1.

1080 **6.11 MetaSUB Genomic Database and Statistical Analysis**

1081 Genomic data was acquired from the MetaSUB database and matched by sample names to the corre-
1082 sponding metadata downloaded from the MetaSUB-metadata github repository (<https://github.com/MetaSUB/MetaSUB>
1083 metadata). All data derived from MetaSUB and the subsequent steps described above was then analysed
1084 using Python 3.6. Python packages plotly, matplotlib and seaborn where used for plotting as well as pan-
1085 das to create and manage dataframes. The heatmap is clustered by Euclidean distance on the columns.
1086

References

- 1087
- 1088 Afshinnkoo, E., Chou, C., Alexander, N., Ahsanuddin, S., Schuetz, A. N., and Mason, C. E. (2017).
1089 Precision metagenomics: Rapid metagenomic analyses for infectious disease diagnostics and public
1090 health surveillance. *Journal of Biomolecular Techniques*, 28(1):40–45.
- 1091 Afshinnkoo, E., Meydan, C., Chowdhury, S., Jaroudi, D., Boyer, C., Bernstein, N., Maritz, J. M.,
1092 Reeves, D., Gandara, J., Chhangawala, S., Ahsanuddin, S., Simmons, A., Nessel, T., Sundaresh, B.,
1093 Pereira, E., Jorgensen, E., Kolokotronis, S. O., Kirchberger, N., Garcia, I., Gandara, D., Dhanraj, S.,
1094 Nawrin, T., Saletore, Y., Alexander, N., Vijay, P., Hénaff, E. M., Zumbo, P., Walsh, M., O'Mullan,
1095 G. D., Tighe, S., Dudley, J. T., Dunaif, A., Ennis, S., O'Halloran, E., Magalhaes, T. R., Boone, B.,
1096 Jones, A. L., Muth, T. R., Paolantonio, K. S., Alter, E., Schadt, E. E., Garbarino, J., Prill, R. J.,
1097 Carlton, J. M., Levy, S., and Mason, C. E. (2015). Geospatial Resolution of Human and Bacterial
1098 Diversity with City-Scale Metagenomics. *Cell Systems*, 1(1):72–87.
- 1099 Allen, H. K., Moe, L. A., Rodbumrer, J., Gaarder, A., and Handelsman, J. (2009). Functional metage-
1100 nomics reveals diverse β -lactamases in a remote alaskan soil. *The ISME journal*, 3(2):243–251.
- 1101 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. (1990). Altschul et al.. 1990.
1102 Basic Local Alignment Search Tool.pdf.
- 1103 Bahram, M., Hildebrand, F., Forslund, S. K., Anderson, J. L., Soudzilovskaia, N. A., Bodegom, P. M.,
1104 Bengtsson-Palme, J., Anslan, S., Coelho, L. P., Harend, H., Huerta-Cepas, J., Medema, M. H., Maltz,
1105 M. R., Mundra, S., Olsson, P. A., Pent, M., Pölme, S., Sunagawa, S., Ryberg, M., Tedersoo, L., and
1106 Bork, P. (2018). Structure and function of the global topsoil microbiome.
- 1107 Bougnom, B. P. and Piddock, L. J. (2017). Wastewater for Urban Agriculture: A Significant Factor in
1108 Dissemination of Antibiotic Resistance.
- 1109 Bowers, R. M., Kyrpides, N. C., Stepanauskas, R., Harmon-Smith, M., Doud, D., Reddy, T. B. K.,
1110 Schulz, F., Jarett, J., Rivers, A. R., Eloë-Fadrosch, E. A., Tringe, S. G., Ivanova, N. N., Copeland,
1111 A., Clum, A., Becraft, E. D., Malmstrom, R. R., Birren, B., Podar, M., Bork, P., Weinstock, G. M.,
1112 Garrity, G. M., Dodsworth, J. A., Yooseph, S., Sutton, G., Glöckner, F. O., Gilbert, J. A., Nelson,
1113 W. C., Hallam, S. J., Jungbluth, S. P., Ettema, T. J. G., Tighe, S., Konstantinidis, K. T., Liu, W. T.,
1114 Baker, B. J., Rattei, T., Eisen, J. A., Hedlund, B., McMahon, K. D., Fierer, N., Knight, R., Finn, R.,
1115 Cochrane, G., Karsch-Mizrachi, I., Tyson, G. W., Rinke, C., Lapidus, A., Meyer, F., Yilmaz, P., Parks,
1116 D. H., Eren, A. M., Schriml, L., Banfield, J. F., Hugenholtz, P., and Woyke, T. (2018). Corrigendum:
1117 Minimum information about a single amplified genome (MISAG) and a metagenome-assembled genome
1118 (MIMAG) of bacteria and archaea. *Nat. Biotechnol.*, 36(7):660.
- 1119 Breitwieser, F. P., Baker, D. N., and Salzberg, S. L. (2018). KrakenUniq: confident and fast metagenomics
1120 classification using unique k-mer counts. *Genome biology*, 19(1):198.
- 1121 Brito, I. L., Yilmaz, S., Huang, K., Xu, L., Jupiter, S. D., Jenkins, A. P., Naisilisili, W., Tamminen, M.,
1122 Smillie, C. S., Wortman, J. R., Birren, B. W., Xavier, R. J., Blainey, P. C., Singh, A. K., Gevers, D.,
1123 and Alm, E. J. (2016). Mobile genes in the human microbiome are structured from global to individual
1124 scales. *Nature*, 535(7612):435–439.
- 1125 Brooks, B., Olm, M. R., Firek, B. A., Baker, R., Thomas, B. C., Morowitz, M. J., and Banfield, J. F.
1126 (2017). Strain-resolved analysis of hospital rooms and infants reveals overlap between the human and
1127 room microbiome. *Nature communications*, 8(1):1–7.
- 1128 Buchfink, B., Xie, C., and Huson, D. H. (2014). Fast and sensitive protein alignment using DIAMOND.
- 1129 Cáliz, J., Triadó-Margarit, X., Camarero, L., and Casamayor, E. O. (2018). A long-term survey unveils
1130 strong seasonal patterns in the airborne microbiome coupled to general and regional atmospheric
1131 circulations. *Proceedings of the National Academy of Sciences*, 115(48):12229–12234.
- 1132 Chaumeil, P. A., Mussig, A. J., Hugenholtz, P., and Parks, D. H. (2019). GTDB-Tk: a toolkit to classify
1133 genomes with the Genome Taxonomy Database. *Bioinformatics*.
- 1134 Consortium, T. H. M. P., Human, T., Project, M., Consortium, T. H. M. P., Human, T., and Project, M.
1135 (2012). Structure, function and diversity of the healthy human microbiome. *Nature*, 486(7402):207–14.

- 1136 Cooley, J. D., Wong, W. C., Jumper, C. A., and Straus, D. C. (1998). Correlation between the prevalence
1137 of certain fungi and sick building syndrome. *Occupational and Environmental Medicine*, 55(9):579–584.
- 1138 Delgado-Baquerizo, M., Oliverio, A. M., Brewer, T. E., Benavent-González, A., Eldridge, D. J., Bardgett,
1139 R. D., Maestre, F. T., Singh, B. K., and Fierer, N. (2018). A global atlas of the dominant bacteria
1140 found in soil. *Science*, 359(6373):320–325.
- 1141 Eckburg, P. B., Mian, M. F., Surette, M. G., Bienenstock, J., Forsythe, P., and Sargent, M. (2005).
1142 Diversity of the Human Intestinal Microbial Flora. *Science*, 308(5728):1635–1638.
- 1143 Eddy, S. R. (2011). Accelerated profile HMM searches. *PLoS Computational Biology*, 7(10).
- 1144 Franzosa, E. A., McIver, L. J., Rahnavard, G., Thompson, L. R., Schirmer, M., Weingart, G., Lipson,
1145 K. S., Knight, R., Caporaso, J. G., Segata, N., and Huttenhower, C. (2018). Species-level functional
1146 profiling of metagenomes and metatranscriptomes. *Nature methods*, 15(11):962–968.
- 1147 Fresia, P., Antelo, V., Salazar, C., Giménez, M., D’Alessandro, B., Afshinnekoo, E., Mason, C., Gonnet,
1148 G. H., and Iraola, G. (2019). Urban metagenomics uncover antibiotic resistance reservoirs in coastal
1149 beach and sewage waters. *Microbiome*, 7(1).
- 1150 Galili, T. (2015). dendextend: an R package for visualizing, adjusting and comparing trees of hierarchical
1151 clustering. *Bioinformatics*, 31(22):3718–3720.
- 1152 Gardy, J. L. and Loman, N. J. (2018). Towards a genomics-informed, real-time, global pathogen surveil-
1153 lance system.
- 1154 Gibson, M. K., Forsberg, K. J., and Dantas, G. (2015). Improved annotation of antibiotic resistance
1155 determinants reveals microbial resistomes cluster by ecology. *ISME Journal*, 9(1):207–216.
- 1156 Gilbert, J. A. and Stephens, B. (2018). Microbiology of the built environment.
- 1157 Hendriksen, R. S., Munk, P., Njage, P., van Bunnik, B., McNally, L., Lukjancenko, O., Röder, T.,
1158 Nieuwenhuijse, D., Pedersen, S. K., Kjeldgaard, J., Kaas, R. S., Clausen, P. T. L. C., Vogt, J. K.,
1159 Leekitcharoenphon, P., van de Schans, M. G. M., Zuidema, T., de Roda Husman, A. M., Rasmussen,
1160 S., Petersen, B., Bego, A., Rees, C., Cassar, S., Coventry, K., Collignon, P., Allerberger, F., Rahube,
1161 T. O., Oliveira, G., Ivanov, I., Vuthy, Y., Sopheak, T., Yost, C. K., Ke, C., Zheng, H., Baisheng,
1162 L., Jiao, X., Donado-Godoy, P., Coulibaly, K. J., Jergović, M., Hrenovic, J., Karpíšková, R., Villacis,
1163 J. E., Legesse, M., Eguale, T., Heikinheimo, A., Malania, L., Nitsche, A., Brinkmann, A., Saba,
1164 C. K. S., Kocsis, B., Solymosi, N., Thorsteinsdottir, T. R., Hatha, A. M., Alebouyeh, M., Morris,
1165 D., Cormican, M., O’Connor, L., Moran-Gilad, J., Alba, P., Battisti, A., Shakenova, Z., Kiiyukia, C.,
1166 Ng’eno, E., Raka, L., Avsejenko, J., Bērziņš, A., Bartkevics, V., Penny, C., Rajandas, H., Parimannan,
1167 S., Haber, M. V., Pal, P., Jeunen, G.-J., Gemell, N., Fashae, K., Holmstad, R., Hasan, R., Shakoor,
1168 S., Rojas, M. L. Z., Wasyl, D., Bosevska, G., Kochubovski, M., Radu, C., Gassama, A., Radosavljevic,
1169 V., Wuertz, S., Zuniga-Montanez, R., Tay, M. Y. F., Gavačová, D., Pastuchova, K., Truska, P., Trkov,
1170 M., Esterhuyse, K., Keddy, K., Cerdà-Cuéllar, M., Pathirage, S., Norrgren, L., Örn, S., Larsson,
1171 D. G. J., Heijden, T. V. d., Kumburu, H. H., Sanneh, B., Bidjada, P., Njanpop-Lafourcade, B.-M.,
1172 Nikiema-Pessinaba, S. C., Levent, B., Meschke, J. S., Beck, N. K., Van, C. D., Phuc, N. D., Tran,
1173 D. M. N., Kwenda, G., Tabo, D.-a., Wester, A. L., Cuadros-Orellana, S., Amid, C., Cochrane, G.,
1174 Sicheritz-Ponten, T., Schmitt, H., Alvarez, J. R. M., Aidara-Kane, A., Pamp, S. J., Lund, O., Hald,
1175 T., Woolhouse, M., Koopmans, M. P., Vigre, H., Petersen, T. N., Aarestrup, F. M., and project
1176 consortium, T. G. S. S. (2019). Global monitoring of antimicrobial resistance based on metagenomics
1177 analyses of urban sewage. *Nature Communications*, 10(1):1124.
- 1178 Hoch, J. M., Rhodes, M. E., Shek, K. L., Dinwiddie, D., Hiebert, T. C., Gill, A. S., Salazar Estrada,
1179 A., Griffin, K. L., Palmer, M. I., and McGuire, K. L. (2019). Soil microbial assemblages are linked to
1180 plant community composition and contribute to ecosystem services on urban green roofs. *Front. Ecol.*
1181 *Evol.* 7: 198. doi: 10.3389/fevo.
- 1182 Hsu, T., Joice, R., Vallarino, J., Abu-Ali, G., Hartmann, E. M., Shafquat, A., DuLong, C., Baranowski,
1183 C., Gevers, D., Green, J. L., Morgan, X. C., Spengler, J. D., and Huttenhower, C. (2016). Urban
1184 Transit System Microbial Communities Differ by Surface Type and Interaction with Humans and the
1185 Environment. *mSystems*, 1(3):e00018–16.

- 1186 Hyatt, D., Chen, G. L., LoCascio, P. F., Land, M. L., Larimer, F. W., and Hauser, L. J. (2010). Prodigal:
1187 Prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics*, 11.
- 1188 Joseph, S. M., Battaglia, T., Maritz, J. M., Carlton, J. M., and Blaser, M. J. (2019). Longitudinal
1189 comparison of bacterial diversity and antibiotic resistance genes in new york city sewage. *MSystems*,
1190 4(4):e00327–19.
- 1191 Joyner, J. L., Kerwin, J., Deeb, M., Lozefski, G., Paltseva, A., Prithiviraj, B., McLaughlin, J., Cheng,
1192 Z., Groffman, P., and Muth, T. R. (2019). Green infrastructure design influences urban soil bacteria
1193 communities. *Frontiers in microbiology*, 10:982.
- 1194 Kang, D. D., Li, F., Kirton, E., Thomas, A., Egan, R., An, H., and Wang, Z. (2019). MetaBAT
1195 2: an adaptive binning algorithm for robust and efficient genome reconstruction from metagenome
1196 assemblies. *PeerJ*, 7:e7359.
- 1197 Kang, K., Ni, Y., Li, J., Imamovic, L., Sarkar, C., Kobler, M. D., Heshiki, Y., Zheng, T., Kumari, S.,
1198 Wong, J. C. Y., Archana, A., Wong, C. W. M., Dingle, C., Denizen, S., Baker, D. M., Sommer, M.
1199 O. A., Webster, C. J., and Panagiotou, G. (2018). The Environmental Exposures and Inner- and
1200 Intercity Traffic Flows of the Metro System May Contribute to the Skin Microbiome and Resistome.
1201 *Cell Reports*, 24(5):1190–1202.e5.
- 1202 Kim, D., Hofstaedter, C. E., Zhao, C., Mattei, L., Tanes, C., Clarke, E., Lauder, A., Sherrill-Mix,
1203 S., Chehoud, C., Kelsen, J., et al. (2017). Optimizing methods and dodging pitfalls in microbiome
1204 research. *Microbiome*, 5(1):52.
- 1205 Klein, E. Y., Van Boeckel, T. P., Martinez, E. M., Pant, S., Gandra, S., Levin, S. A., Goossens, H.,
1206 and Laxminarayan, R. (2018). Global increase and geographic convergence in antibiotic consumption
1207 between 2000 and 2015. *Proceedings of the National Academy of Sciences*, 115(15):E3463–E3470.
- 1208 Korownyk, C., Liu, F., and Garrison, S. (2018). Population level evidence for seasonality of the human
1209 microbiome. *Chronobiology International*, 35(4):573–577.
- 1210 Koslicki, D. (2011). Topological entropy of DNA sequences. *Bioinformatics*, 27(8):1061–1067.
- 1211 Krawczyk, P. S., Lipinski, L., and Dziembowski, A. (2018). PlasFlow: predicting plasmid sequences in
1212 metagenomic data using genome signatures. *Nucleic Acids Research*, 46(6):e35–e35.
- 1213 Kurtz, S., Phillippy, A., Delcher, A. L., Smoot, M., Shumway, M., Antonescu, C., and Salzberg, S. L.
1214 (2004). Versatile and open software for comparing large genomes. *Genome Biol.*, 5(2):R12.
- 1215 Ladner, J. T., Grubaugh, N. D., Pybus, O. G., and Andersen, K. G. (2019). Precision epidemiology for
1216 infectious disease control.
- 1217 Lakin, S. M., Dean, C., Noyes, N. R., Dettenwanger, A., Ross, A. S., Doster, E., Rovira, P., Abdo, Z.,
1218 Jones, K. L., Ruiz, J., Belk, K. E., Morley, P. S., and Boucher, C. (2017). MEGARes: An antimicrobial
1219 resistance database for high throughput sequencing. *Nucleic Acids Research*, 45(D1):D574–D580.
- 1220 Langmead and Steven L Salzberg (2013). Bowtie2. *Nature methods*, 9(4):357–359.
- 1221 Lax, S., Sangwan, N., Smith, D., Larsen, P., Handley, K. M., Richardson, M., Guyton, K., Krezalek,
1222 M., Shogan, B. D., Defazio, J., et al. (2017). Bacterial colonization and succession in a newly opened
1223 hospital. *Science translational medicine*, 9(391):eaah6500.
- 1224 Letunic, I. and Bork, P. (2019). Interactive Tree Of Life (iTOL) v4: recent updates and new develop-
1225 ments. *Nucleic Acids Res.*, 47(W1):W256–W259.
- 1226 Leung, M. H., Wilkins, D., Li, E. K., Kong, F. K., and Lee, P. K. (2014). Indoor-air microbiome in an
1227 urban subway network: diversity and dynamics. *Appl. Environ. Microbiol.*, 80(21):6760–6770.
- 1228 Lloyd-Price, J., Mahurkar, A., Rahnavard, G., Crabtree, J., Orvis, J., Hall, A. B., Brady, A., Creasy,
1229 H. H., McCracken, C., Giglio, M. G., McDonald, D., Franzosa, E. A., Knight, R., White, O., and
1230 Huttenhower, C. (2017). Strains, functions and dynamics in the expanded Human Microbiome Project.
1231 *Nature*, 550(7674):61–66.

- 1232 Maritz, J. M., Ten Eyck, T. A., Alter, S. E., and Carlton, J. M. (2019). Patterns of protist diversity
1233 associated with raw sewage in new york city. *The ISME journal*, 13(11):2750–2763.
- 1234 Martínez, J. L. (2008). Antibiotics and antibiotic resistance genes in natural environments. *Science*,
1235 321(5887):365–367.
- 1236 Mason-Buck, G., Graf, A., Elhaik, E., Robinson, J., Pospiech, E., Oliveira, M., Moser, J., Lee, P. K. H.,
1237 Githae, D., Ballard, D., Bromberg, Y., Casimiro-Soriguer, C. S., Dhungel, E., Ahn, T.-H., Kawulok,
1238 J., Loucera, C., Ryan, F., Walker, A. R., Zhu, C., Mason, C. E., Amorim, A., Syndercombe Court,
1239 D., Branicki, W., and Łabaj, P. (2020). DNA based methods in intelligence - moving towards metage-
1240 nomics. *Preprints*, page 2020020158.
- 1241 McArthur, A. G., Waglechner, N., Nizam, F., Yan, A., Azad, M. A., Baylay, A. J., Bhullar, K., Canova,
1242 M. J., De Pascale, G., Ejim, L., Kalan, L., King, A. M., Koteva, K., Morar, M., Mulvey, M. R.,
1243 O'Brien, J. S., Pawlowski, A. C., Piddock, L. J., Spanogiannopoulos, P., Sutherland, A. D., Tang, I.,
1244 Taylor, P. L., Thaker, M., Wang, W., Yan, M., Yu, T., and Wright, G. D. (2013). The comprehensive
1245 antibiotic resistance database. *Antimicrobial Agents and Chemotherapy*, 57(7):3348–3357.
- 1246 McInnes, L., Healy, J., Saul, N., and Großberger, L. (2018). UMAP: Uniform Manifold Approximation
1247 and Projection. *Journal of Open Source Software*, 3(29):861.
- 1248 McIntyre, A. B., Ounit, R., Afshinnekoo, E., Prill, R. J., Hénaff, E., Alexander, N., Minot, S. S., Danko,
1249 D., Foux, J., Ahsanuddin, S., Tighe, S., Hasan, N. A., Subramanian, P., Moffat, K., Levy, S., Lonardi,
1250 S., Greenfield, N., Colwell, R. R., Rosen, G. L., and Mason, C. E. (2017). Comprehensive benchmarking
1251 and ensemble approaches for metagenomic classifiers. *Genome Biology*, 18(1).
- 1252 McLaren, M. R., Willis, A. D., and Callahan, B. J. (2019). Consistent and correctable bias in metage-
1253 nomic sequencing measurements. *BioRxiv*, page 559831.
- 1254 MetaSUB International Consortium. Mason, C., Afshinnekoo, E., Ahsannudin, S., Ghedin, E., Read,
1255 T., Fraser, C., Dudley, J., Hernandez, M., Bowler, C., Stolovitzky, G., Chernonetz, A., Gray, A.,
1256 Darling, A., Burke, C., ?abaj, P. P., Graf, A., Noushmehr, H., Moraes, S., Dias-Neto, E., Ugalde, J.,
1257 Guo, Y., Zhou, Y., Xie, Z., Zheng, D., Zhou, H., Shi, L., Zhu, S., Tang, A., Ivankovi?, T., Siam, R.,
1258 Rascovan, N., Richard, H., Lafontaine, I., Baron, C., Nedunuri, N., Prithiviraj, B., Hyat, S., Mehr,
1259 S., Banihashemi, K., Segata, N., Suzuki, H., Alpuche Aranda, C. M., Martinez, J., Christopher Dada,
1260 A., Osuolale, O., Oguntoyinbo, F., Dybwad, M., Oliveira, M., Fernandes, A., Oliveira, M., Fernandes,
1261 A., Chatziefthimiou, A. D., Chaker, S., Alexeev, D., Chuvelev, D., Kurilshikov, A., Schuster, S., Siwo,
1262 G. H., Jang, S., Seo, S. C., Hwang, S. H., Ossowski, S., Bezdán, D., Udekwu, K., Udekwu, K., Lungj-
1263 dahl, P. O., Nikolayeva, O., Sezerman, U., Kelly, F., Metrustry, S., Elhaik, E., Gonnet, G., Schriml,
1264 L., Mongodin, E., Huttenhower, C., Gilbert, J., Hernandez, M., Vayndorf, E., Blaser, M., Schadt,
1265 E., Eisen, J., Beitel, C., Hirschberg, D., Schriml, L., and Mongodin, E. (2016). The Metagenomics
1266 and Metadesign of the Subways and Urban Biomes (MetaSUB) International Consortium inaugural
1267 meeting report. *Microbiome*, 4(1):24.
- 1268 Meyer, K. M., Memiaghe, H., Korte, L., Kenfack, D., Alonso, A., and Bohannan, B. J. (2018). Why do
1269 microbes exhibit weak biogeographic patterns? *ISME Journal*, 12(6):1404–1413.
- 1270 Moskowitz, D. M. and Greenleaf, W. J. (2018). Nonparametric analysis of contributions to variance in
1271 genomics and epigenomics data. *bioRxiv*.
- 1272 Nayfach, S., Shi, Z. J., Seshadri, R., Pollard, K. S., and Kyrpides, N. C. (2019). New insights from
1273 uncultivated genomes of the global human gut microbiome. *Nature*, 568(7753):505–510.
- 1274 Neiderud, C. J. (2015). How urbanization affects the epidemiology of emerging infectious diseases. *African
1275 Journal of Disability*, 5(1).
- 1276 Nicolaou, N., Siddique, N., and Custovic, A. (2005). Allergic disease in urban and rural populations:
1277 Increasing prevalence with increasing urbanization. *Allergy: European Journal of Allergy and Clinical
1278 Immunology*, 60(11):1357–1360.
- 1279 Nurk, S., Meleshko, D., Korobeynikov, A., and Pevzner, P. A. (2017). MetaSPAdes: A new versatile
1280 metagenomic assembler. *Genome Research*, 27(5):824–834.

- 1281 O'Hara, N. B., Reed, H. J., Afshinnkoo, E., Harvin, D., Caplan, N., Rosen, G., Frye, B., Woloszynek, S.,
1282 Ounit, R., Levy, S., Butler, E., and Mason, C. E. (2017). Metagenomic characterization of ambulances
1283 across the USA. *Microbiome*, 5(1):125.
- 1284 Olm, M. R., Brown, C. T., Brooks, B., and Banfield, J. F. (2017). dRep: a tool for fast and accurate ge-
1285 nomic comparisons that enables improved genome recovery from metagenomes through de-replication.
1286 *ISME J*, 11(12):2864–2868.
- 1287 Ondov, B. D., Treangen, T. J., Melsted, P., Mallonee, A. B., Bergman, N. H., Koren, S., and Phillippy,
1288 A. M. (2016). Mash: fast genome and metagenome distance estimation using MinHash. *Genome*
1289 *biology*, 17(1):132.
- 1290 Paez-Espino, D., Eloë-Fadrosh, E. A., Pavlopoulos, G. A., Thomas, A. D., Huntemann, M., Mikhailova,
1291 N., Rubin, E., Ivanova, N. N., and Kyrpides, N. C. (2016a). Uncovering Earth's virome. *Nature*,
1292 536(7617):425–430.
- 1293 Paez-Espino, D., Eloë-Fadrosh, E. A., Pavlopoulos, G. A., Thomas, A. D., Huntemann, M., Mikhailova,
1294 N., Rubin, E., Ivanova, N. N., and Kyrpides, N. C. (2016b). Uncovering earth's virome. *Nature*,
1295 536(7617):425–430.
- 1296 Paez-Espino, D., Pavlopoulos, G. A., Ivanova, N. N., and Kyrpides, N. C. (2017). Nontargeted virus
1297 sequence discovery pipeline and virus clustering for metagenomic data. *Nature Protocols*, 12(8):1673–
1298 1682.
- 1299 Paez-Espino, D., Roux, S., Chen, I. M. A., Palaniappan, K., Ratner, A., Chu, K., Huntemann, M.,
1300 Reddy, T. B., Pons, J. C., Llabrés, M., Eloë-Fadrosh, E. A., Ivanova, N. N., and Kyrpides, N. C.
1301 (2019). IMG/VR v.2.0: An integrated data management and analysis system for cultivated and
1302 environmental viral genomes. *Nucleic Acids Research*, 47(D1):D678–D686.
- 1303 Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P., and Tyson, G. W. (2015). CheckM:
1304 assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes.
1305 *Genome Res.*, 25(7):1043–1055.
- 1306 Parks, D. H., Rinke, C., Chuvochina, M., Chaumeil, P. A., Woodcroft, B. J., Evans, P. N., Hugenholtz,
1307 P., and Tyson, G. W. (2018). Author Correction: Recovery of nearly 8,000 metagenome-assembled
1308 genomes substantially expands the tree of life. *Nat Microbiol*, 3(2):253.
- 1309 Pasolli, E., Asnicar, F., Manara, S., Zolfo, M., Karcher, N., Armanini, F., Beghini, F., Manghi, P.,
1310 Tett, A., Ghensi, P., et al. (2019). Extensive unexplored human microbiome diversity revealed by over
1311 150,000 genomes from metagenomes spanning age, geography, and lifestyle. *Cell*, 176(3):649–662.
- 1312 Price, M. N., Dehal, P. S., and Arkin, A. P. (2010). FastTree 2—approximately maximum-likelihood trees
1313 for large alignments. *PLoS ONE*, 5(3):e9490.
- 1314 Qin, J., Li, R., Raes, J., and Arumugam, M. (2010). A human gut microbial gene catalogue established
1315 by metagenomic sequencing: Commentary.
- 1316 Rahman, S. F., Olm, M. R., Morowitz, M. J., and Banfield, J. F. (2018). Machine Learning Leverag-
1317 ing Genomes from Metagenomes Identifies Influential Antibiotic Resistance Genes in the Infant Gut
1318 Microbiome. *mSystems*, 3(1).
- 1319 Rice, L. B. (2012). Mechanisms of resistance and clinical relevance of resistance to β -lactams, glycopep-
1320 tides, and fluoroquinolones. In *Mayo Clinic Proceedings*, volume 87, pages 198–208. Elsevier.
- 1321 Ritchie, H. and Roser, M. (2020). Urbanization. *Our World in Data*.
1322 <https://ourworldindata.org/urbanization>.
- 1323 Schubert, M., Lindgreen, S., and Orlando, L. (2016). AdapterRemoval v2: Rapid adapter trimming,
1324 identification, and read merging. *BMC Research Notes*, 9(1).
- 1325 Segata, N., Waldron, L., Ballarini, A., Narasimhan, V., Jousson, O., and Huttenhower, C. (2012).
1326 Metagenomic microbial community profiling using unique clade-specific marker genes. *Nature Methods*,
1327 9(8):811.

- 1328 Shaaban, H., Westfall, D. A., Mohammad, R., Danko, D., Bezdán, D., Afshinnekoo, E., Segata, N.,
1329 and Mason, C. E. (2018). The Microbe Directory: An annotated, searchable inventory of microbes'
1330 characteristics. *Gates Open Research*, 2:3.
- 1331 Singer, A. C., Shaw, H., Rhodes, V., and Hart, A. (2016). Review of antimicrobial resistance in the
1332 environment and its relevance to environmental regulators. *Frontiers in microbiology*, 7:1728.
- 1333 Thanner, S., Drissner, D., and Walsh, F. (2016). Antimicrobial resistance in agriculture. *MBio*,
1334 7(2):e02227–15.
- 1335 Thompson, L. R., Sanders, J. G., McDonald, D., Amir, A., Ladau, J., Locey, K. J., Prill, R. J., Tripathi,
1336 A., Gibbons, S. M., Ackermann, G., Navas-Molina, J. A., Janssen, S., Kopylova, E., Vazquez-Baeza,
1337 Y., Gonzalez, A., Morton, J. T., Mirarab, S., Xu, Z. Z., Jiang, L., Haroon, M. F., Kanbar, J., Zhu, Q.,
1338 Song, S. J., Kosciulek, T., Bokulich, N. A., Lefler, J., Brislawn, C. J., Humphrey, G., Owens, S. M.,
1339 Hampton-Marcell, J., Berg-Lyons, D., McKenzie, V., Fierer, N., Fuhrman, J. A., Clauset, A., Stevens,
1340 R. L., Shade, A., Pollard, K. S., Goodwin, K. D., Jansson, J. K., Gilbert, J. A., and Knight, R. (2017).
1341 A communal catalogue reveals Earth's multiscale microbial diversity. *Nature*, 551(7681):457–463.
- 1342 United Nations (2016). Political declaration of the high-level meeting of the General Assembly on
1343 antimicrobial resistance. Technical report.
- 1344 United Nations (2018). World Urbanization Prospects: The 2018 Revision. Key facts. Technical report.
- 1345 Van Boeckel, T. P., Brower, C., Gilbert, M., Grenfell, B. T., Levin, S. A., Robinson, T. P., Teillant, A.,
1346 and Laxminarayan, R. (2015). Global trends in antimicrobial use in food animals. *Proceedings of the
1347 National Academy of Sciences of the United States of America*, 112(18):5649–54.
- 1348 Venter, H., Henningsen, M. L., and Begg, S. L. (2017). Antimicrobial resistance in healthcare, agriculture
1349 and the environment: the biochemistry behind the headlines. *Essays in biochemistry*, 61(1):1–10.
- 1350 Weiss, S., Xu, Z. Z., Peddada, S., Amir, A., Bittinger, K., Gonzalez, A., Lozupone, C., Zaneveld, J. R.,
1351 Vázquez-Baeza, Y., Birmingham, A., Hyde, E. R., and Knight, R. (2017). Normalization and microbial
1352 differential abundance strategies depend upon data characteristics. *Microbiome*, 5(1):27.
- 1353 Wilson, M. R., Sample, H. A., Zorn, K. C., Arevalo, S., Yu, G., Neuhaus, J., Federman, S., Stryke,
1354 D., Briggs, B., Langelier, C., Berger, A., Douglas, V., Josephson, S. A., Chow, F. C., Fulton, B. D.,
1355 DeRisi, J. L., Gelfand, J. M., Naccache, S. N., Bender, J., Dien Bard, J., Murkey, J., Carlson, M.,
1356 Vespa, P. M., Vijayan, T., Allyn, P. R., Campeau, S., Humphries, R. M., Klausner, J. D., Ganzon,
1357 C. D., Memar, F., Ocampo, N. A., Zimmermann, L. L., Cohen, S. H., Polage, C. R., DeBiasi, R. L.,
1358 Haller, B., Dallas, R., Maron, G., Hayden, R., Messacar, K., Dominguez, S. R., Miller, S., and Chiu,
1359 C. Y. (2019). Clinical Metagenomic Sequencing for Diagnosis of Meningitis and Encephalitis. *New
1360 England Journal of Medicine*, 380(24):2327–2340.
- 1361 Yooseph, S., Andrews-Pfannkoch, C., Tenney, A., McQuaid, J., Williamson, S., Thiagarajan, M., Bami,
1362 D., Zeigler-Allen, L., Hoffman, J., Goll, J. B., et al. (2013). A metagenomic framework for the study
1363 of airborne microbial communities. *PLoS One*, 8(12):e81862.
- 1364 Zhu, Y.-G., Gillings, M., Simonet, P., Stekel, D., Banwart, S., and Penuelas, J. (2017). Microbial mass
1365 movements. *Science*, 357(6356):1099–1100.

1366 7 Contributing Members of the MetaSUB Consortium

1367 Marcos Abraao, Muhammad Afaq, Ireen Alam, Gabriela E Albuquerque, Kalyn Ali, Lucia E Alvarado-
1368 Arnez, Sarh Aly, Jennifer Amachee, Maria G. Amorim, Majelia Ampadu, Nala An, Núria Andreu So-
1369 mavilla, Michael Angelov, Verónica Antelo, Catharine Aquino, Mayra Arauco Livia, Luiza F Araujo,
1370 Jenny Arevalo, Lucia Elena Alvarado Arnez, Fernanda Arredondo, Matthew Arthur, Sadaf Ayaz, Silva
1371 Baburyan, Abd-Manaaf Bakere, Katrin Bakhil, Thais F. Bartelli, Kevin Becher, Joseph Benson, Denis
1372 Bertrand, Silvia Beurmann, Christina Black, Brittany Blyther, Bazartseren Boldgiv, Gabriela P Branco,
1373 Christian Brion, Paulina Buczanska, Catherine M Burke, Irvind Buttar, Jalia Bynoe, Sven Bönigk, Kari
1374 O Bøifot, Hiram Caballero, Alessandra Carbone, Anais Cardenas, Ana V Castro, Ana Valeria B Castro,
1375 Astred Castro, Simone Cawthorne, Jonathan Cedillo, Salama Chaker, Allison Chan, Anastasia I Chas-
1376 api, Gregory Chem, Jenn-Wei Chen, Michelle Chen, Xiaoqing Chen, Ariel Chernomoretz, Daisy Cheung,
1377 Diana Chicas, Hira Choudhry, Carl Chrispin, Kianna Ciaramella, Jake Cohen, David A Coil, Colleen
1378 Conger, Ana F. Costa, Delisia Cuebas, Aaron E Darling, Pujita Das, Lucinda B Davenport, Laurent
1379 David, Gargi Dayama, Paola F De Sessions, Chris K Deng, Monika Devi, Felipe S Dezem, Sonia Dorado,
1380 LaShonda Dorsey, Steven Du, Alexandra Dutan, Naya Eady, Stephen Eduard Boja Ruiz, Jonathan A
1381 Eisen, Miar Elaskandrany, Lennard Epping, Juan P Escalera-Antezana, Iqra Faiz, Luce Fan, Nadine
1382 Farhat, Kelly French, Skye Felice, Laís Pereira Ferreira, Gabriel Figueroa, Denisse Flores, Marcos AS
1383 Fonseca, Jonathan Foon, Aaishah Francis, Pablo Fresia, Jacob Friedman, Jaime J Fuentes, Josephine
1384 Galipon, Laura Garcia, Annie Geiger, Samuel M Gerner, Dao Phuong Giang, Matías Giménez, Donato
1385 Giovannelli, Dedan Githae, Samantha Goldman, Gaston H Gonnet, Juana Gonzalez, Irene González
1386 Navarrete, Tranette Gregory, Felix Hartkopf, Arya Hawkins-Zafarnia, Nur Hazlin Hazrin-Chong, Tam-
1387 era Henry, Samuel Hernandez, David Hess-Homeier, Yui Him Lo, Lauren E Hittle, Nghiem Xuan Hoan,
1388 Irene Hoxie, Elizabeth Humphries, Shaikh B Iqbal, Riham Islam, Sharah Islam, Takayuki Ito, Tomislav
1389 Ivankovic, Sarah Jackson, JoAnn Jacobs, Esmeralda Jiminez, Ayantu Jinfessa, Takema Kajita, Amrit
1390 Kaur, Fernanda de Souza Gomes Kehdy, Vedbar S Khadka, Shaira Khan, Michelle Ki, Gina Kim, Hyung
1391 Jun Kim, Sangwan Kim, Ryan J King, Kaymisha Knights, Ellen Koag, Nadezhda Kobko-Litskevitch,
1392 Giuseppe KoLoMonaco, Michael Kozhar, Nanami Kubota, Sheelta S Kumar, Lawrence Kwong, Rachel
1393 Kwong, Ingrid Lafontaine, Manolo Laiola, Isha Lamba, Hyunjung Lee, Lucy Lee, Yunmi Lee, Emily
1394 Leong, Marcus H Y Leung, Chenhao Li, Weijun Liang, Moses Lin, Yan Ling Wong, Priscilla Lisboa,
1395 Anna Litskevitch, Tracy Liu, Sonia Losim, Jennifer Lu, Simona Lysakova, Gustavo Adolfo Malca Salas,
1396 Denisse Maldonado, Krizzy Mallari, Tathiane M Malta, Maliha Mamun, Yuk Man Tang, Sonia Mari-
1397 novic, Brunna Marques, Nicole Mathews, Yuri Matsuzaki, Madelyn May, Elias McComb, Adiel Melamed,
1398 Wayne Menary, Ambar Mendez, Katterinne N Mendez, Irene Meng, Ajay Menon, Mark Menor, Nancy
1399 Merino, Cem Meydan, Karishma Miah, Tanja Miketic, Eric Minwei Liu, Wilson Miranda, Athena Mit-
1400 sios, Natasha Mohan, Mohammed Mohsin, Karobi Moitra, Laura Molina, Eftar Moniruzzaman, Sookwon
1401 Moon, Isabelle de Oliveira Moraes, Maritza S Mosella, Maritza S Mosella, Josef W Moser, Christopher
1402 Mozsary, Amanda L Muehlbauer, Oasima Muner, Muntaha Munia, Naimah Munim, Tatjana Mustac,
1403 Kaung Myat San, Areeg Naeem, Mayuko Nakagawa, Masaki Nasu, Bryan Nazario, Narasimha Rao
1404 Nedunuri, Aida Nesimi, Aida Nesimi, Gloria Nguyen, Hosna Noorzi, Avigdor Nosrati, Houtan Noush-
1405 mehr, Diana N. Nunes, Kathryn O'Brien, Niamh B O'Hara, Gabriella Oken, Rantimi A Olawoyin, Kiara
1406 Olmeda, Itunu A Oluwadare, Tolulope Oluwadare, Jenessa Orpilla, Jacqueline Orrego, Melissa Ortega,
1407 Princess Osma, Israel O Osuolale, Oluwatosin M Osuolale, Rachid Ounit, Christos A Ouzounis, Sub-
1408 hamitra Pakrashi, Rachel Paras, Andrea Patrignani, Ante Peros, Sabrina Persaud, Anisia Peters, Robert
1409 A Petit III, Adam Phillips, Lisbeth Pineda, Alketa Plaku, Alma Plaku, Brianna Pompa-Hogan, Max
1410 Priestman, Bharath Prithiviraj, Sambhawa Priya, Phanthira Pugdeethosal, Benjamin Pulatov, Angelika
1411 Pupiec, Tao Qing, Saher Rahiel, Savlatjon Rahmatulloev, Kannan Rajendran, Aneisa Ramcharan, Adan
1412 Ramirez-Rojas, Shahryar Rana, Prashanthi Ratnanandan, Timothy D Read, Hugues Richard, Alexis
1413 Rivera, Michelle Rivera, Alessandro Robertiello, Courtney Robinson, Anyelic Rosario, Kaitlan Russell,
1414 Timothy Ryan Donahoe, Krista Ryon, Thais S Sabedot, Thais S Sabedot, Mahfuza Sabina, Cecilia
1415 Salazar, Jorge Sanchez, Ryan Sankar, Paulo Thiago de Souza Santos, Zulena Saravi, Thomas Saw Aung,
1416 Thomas Saw Aung, Nowshin Sayara, Steffen Schaaf, Anna-Lena M Schinke, Ralph Schlapbach, Jason R
1417 Schriml, Felipe Segato, Marianna S. Serpa, Heba Shaaban, Maheen Shakil, Hyenah Shim, Yuh Shiwa,
1418 Shaleni K Singh, Eunice So, Camila Souza, Jason Sperry, Kiyoshi Suganuma, Hamood Suliman, Jill
1419 Sullivan, Jill Sullivan, Fumie Takahara, Isabella K Takenaka, Anyi Tang, Emilio Tarcitano, Mahdi Taye,
1420 Alexis Terrero, Andrew M Thomas, Sade Thomas, Masaru Tomita, Xinzhao Tong, Jennifer M Tran,
1421 Catalina Truong, Stefan I Tsonev, Kazutoshi Tsuda, Michelle Tuz, Carmen Urgiles, Brandon Valentine,
1422 Hitler Francois Vasquez Arevalo, Valeria Ventorino, Patricia Vera-Wolf, Sierra Vincent, Renee Vivancos-

1423 Koopman, Andrew Wan, Cindy Wang, Samuel Weekes, Xiao Wen Cai, Johannes Werner, David Westfall,
1424 Lothar H Wieler, Michelle Williams, Silver A Wolf, Brian Wong, Tyler Wong, Hyun Woo Joo, Rasheena
1425 Wright, Ryota Yamanaka, Jingcheng Yang, Hirokazu Yano, George C Yeh, Tsoi Ying Lai, Laraib Zafar,
1426 Amy Zhang, Shu Zhang, Yang Zhang, Yuanting Zheng,

1427 8 Supplemental Materials

Table S1: Sample Counts

Region	project city	Pilot	CSD16	CSD17	Other	Total
Control	Background Control	0.0	40	0	0.0	40
	Lab Control	0.0	20	6	0.0	26
	Positive Control	0.0	33	6	0.0	39
East Asia	Region Total	26.0	1297	0	34.0	1357
	Hanoi	0.0	16	0	0.0	16
	Hong Kong	0.0	712	0	12.0	724
	Kuala Lumpur	0.0	30	0	0.0	30
	Sendai	0.0	32	0	0.0	32
	Seoul	0.0	80	0	12.0	92
	Shanghai	0.0	0	0	10.0	10
	Singapore	0.0	192	0	0.0	192
	Taipei	0.0	94	0	0.0	94
	Tokyo	26.0	132	0	0.0	158
	Yamaguchi	0.0	9	0	0.0	9
	Europe	Region Total	310.0	939	1	177.0
Barcelona		99.0	0	0	25.0	124
Belfast		0.0	5	0	0.0	5
Berlin		55.0	1	0	0.0	56
Birmingham		0.0	5	1	0.0	6
Bradford		0.0	4	0	0.0	4
Bury		0.0	6	0	0.0	6
Eastbourne		0.0	6	0	0.0	6
Eden		0.0	5	0	0.0	5
Edinburgh		0.0	6	0	0.0	6
Islington		0.0	5	0	0.0	5
Jaywick		0.0	6	0	0.0	6
Kensington		0.0	6	0	0.0	6
Kyiv		0.0	97	0	0.0	97
Lands End		0.0	5	0	0.0	5
Lisbon		60.0	0	0	28.0	88
London		0.0	534	0	0.0	534
Marseille		96.0	16	0	0.0	112
Naples		0.0	16	0	0.0	16
Newcastle		0.0	5	0	0.0	5
Oslo		0.0	16	0	12.0	28
Paris		0.0	16	0	0.0	16
Porto	0.0	0	0	112.0	112	
Sofia	0.0	16	0	0.0	16	
Stockholm	0.0	62	0	0.0	62	
Swansea	0.0	6	0	0.0	6	
Vienna	0.0	16	0	0.0	16	
Zurich	0.0	79	0	0.0	79	
Middle East	Region Total	100.0	15	0	0.0	115
	Doha	100.0	15	0	0.0	115

Table S1: Sample Counts Cont.

continent	project city	Pilot	CSD16	CSD17	Other	Total
North America	Region Total	284.0	371	276	28.0	959
	Baltimore	0.0	23	0	0.0	23
	Denver	24.0	23	0	0.0	47
	Fairbanks	141.0	0	0	0.0	141
	Mexico City	0.0	0	0	10.0	10
	Minneapolis	0.0	16	0	0.0	16
	New York City	103.0	279	276	0.0	658
	Sacramento	16.0	0	0	18.0	34
	San Francisco	0.0	30	0	0.0	30
Oceania	Region Total	94.0	32	0	0.0	126
	Auckland	16.0	0	0	0.0	16
	Brisbane	0.0	16	0	0.0	16
	Hamilton	16.0	0	0	0.0	16
	Honolulu	0.0	16	0	0.0	16
	Sydney	62.0	0	0	0.0	62
South America	Region Total	44.0	199	68	20.0	331
	Bogota	17.0	0	0	0.0	17
	Montevideo	0.0	0	0	20.0	20
	Ribeirao Preto	0.0	93	0	0.0	93
	Rio De Janeiro	0.0	77	68	0.0	145
	Santiago	27.0	0	0	0.0	27
	Sao Paulo	0.0	29	0	0.0	29
Sub Saharan Africa	Region Total	116.0	192	0	0.0	308
	Ilorin	90.0	134	0	0.0	224
	Offa	26.0	58	0	0.0	84

Table S2: Covariate Variance. The sample variance that can be explained by each factor, in isolation.

Factor	Variance Explained
City	19%
City Population Density	0%
City Ave June Temp	4%
City Elevation	2%
Coastal City	1%
Surface Material	4%
Koppen Climate Classification	8%
Setting	2%
Above/Below Ground	7%
Continent	11%

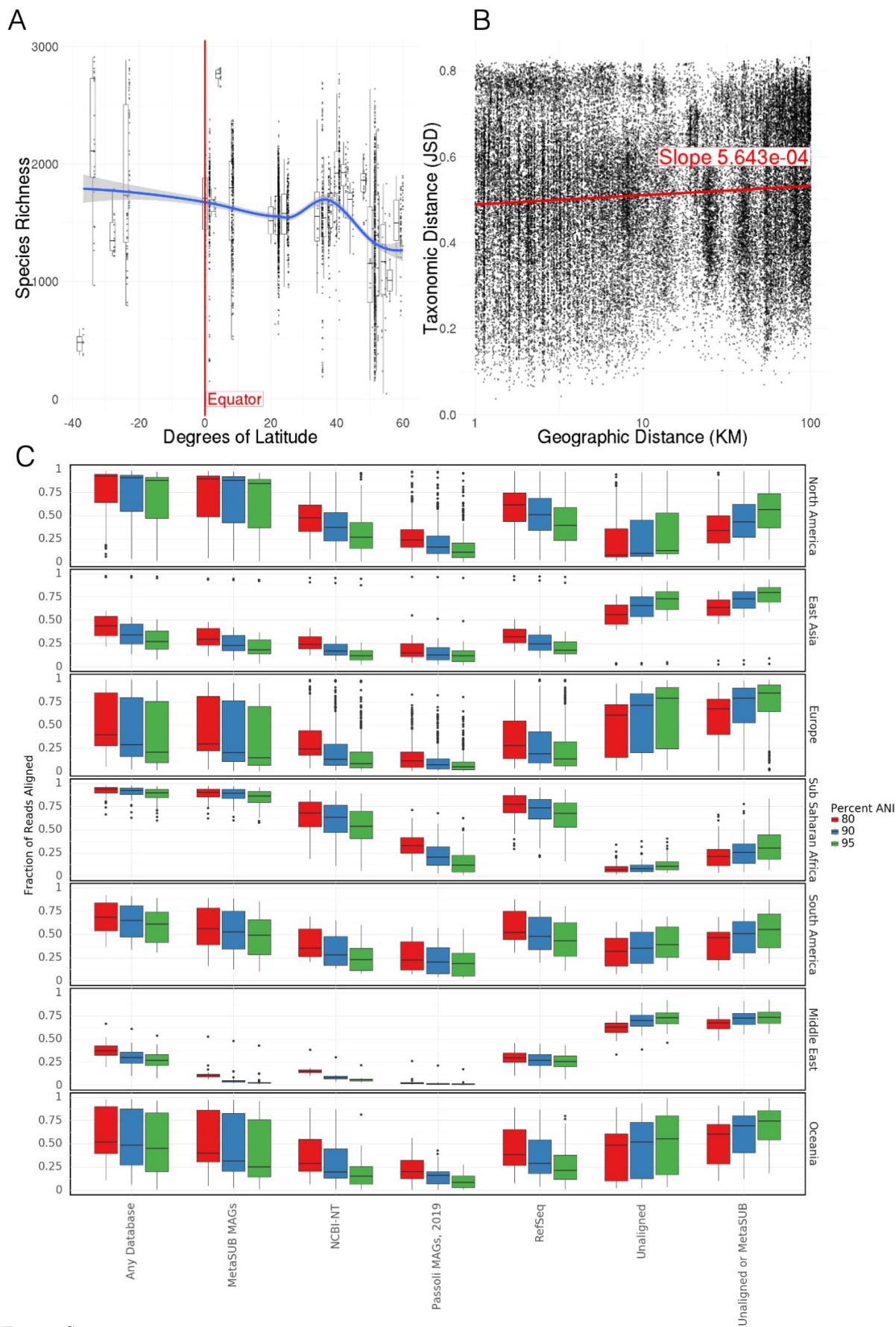


Figure S1: Ecological relationships with taxa. A) Correlation between species richness and latitude. Richness decreases significantly with latitude B) Neighbourhood effect. Taxonomic distance weakly correlates with geographic distance within cities. C) Fraction of reads assigned to different databases by BLAST for each region, at different levels of average nucleotide identity

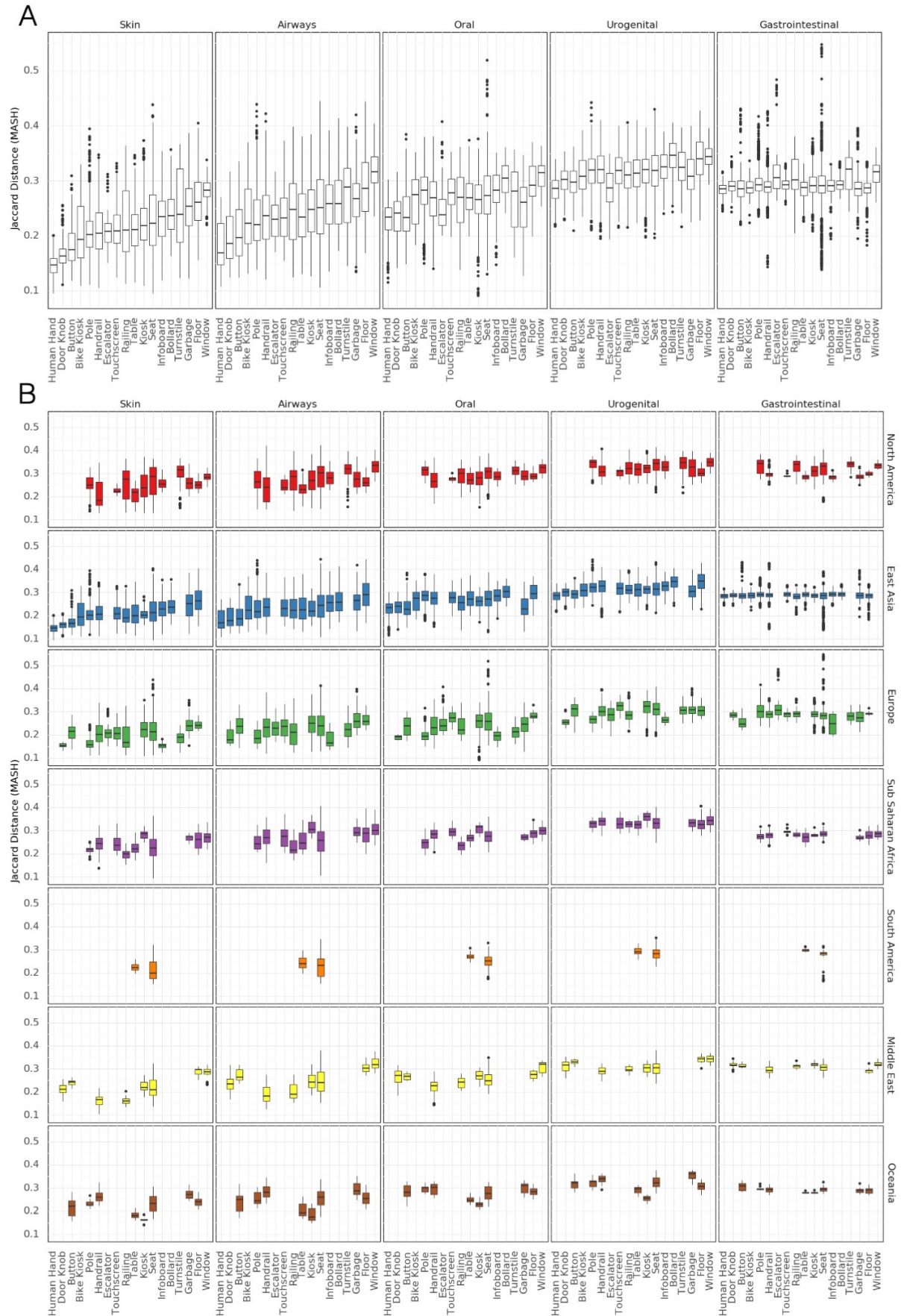


Figure S2: Comparison to Human Microbiome Project. A) Jaccard similarity of MASH indices to HMP samples for different surface types. B) Jaccard similarity of MASH indices to HMP samples for different surface types by region.

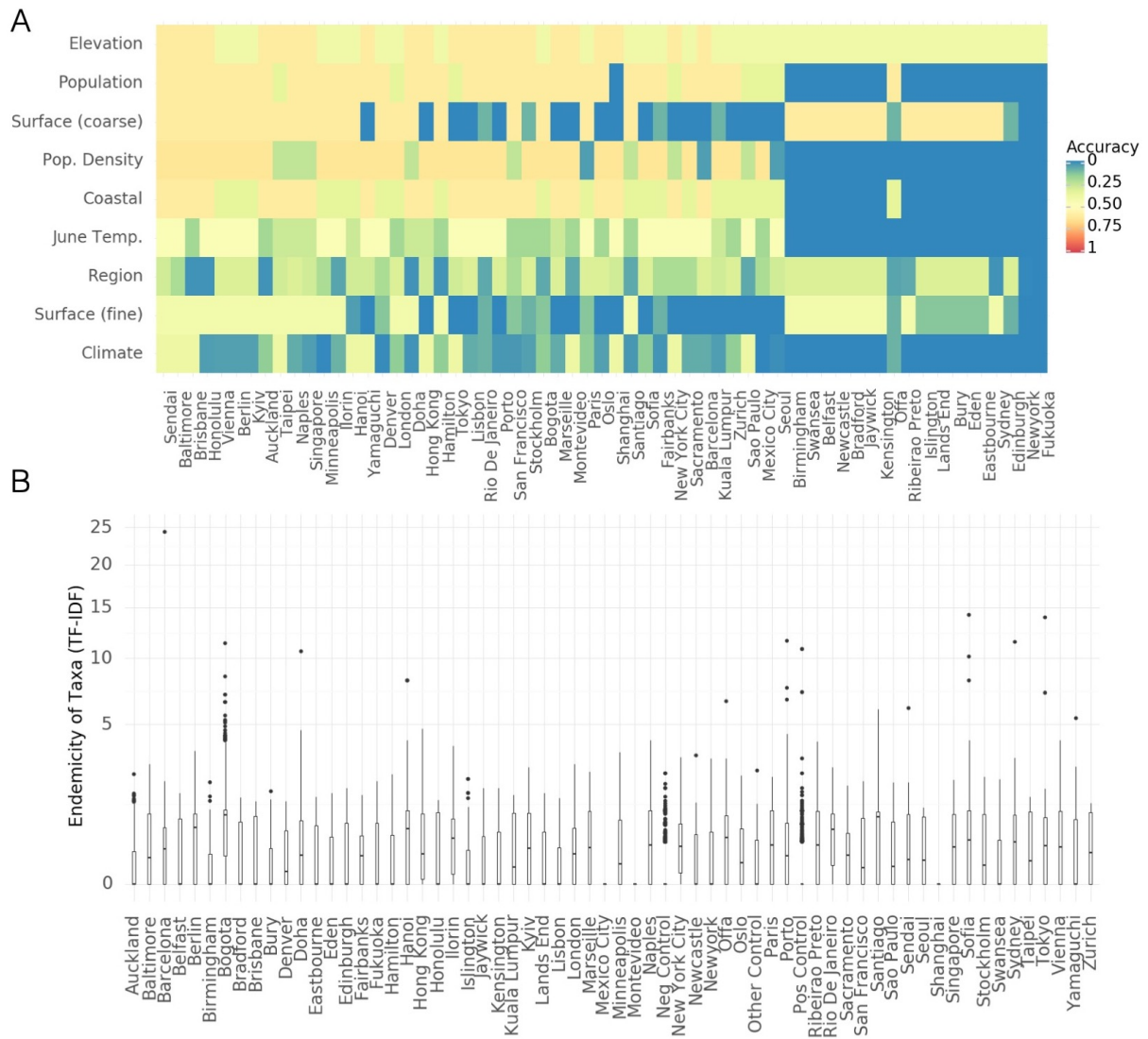


Figure S3: Microbial Signatures, supplemental. A) Classification accuracy that would be achieved by a random model predicting features (rows) for held out cities (columns) B) Endemicity Score (Term Frequency Inverse Document Frequency) for taxa in cities

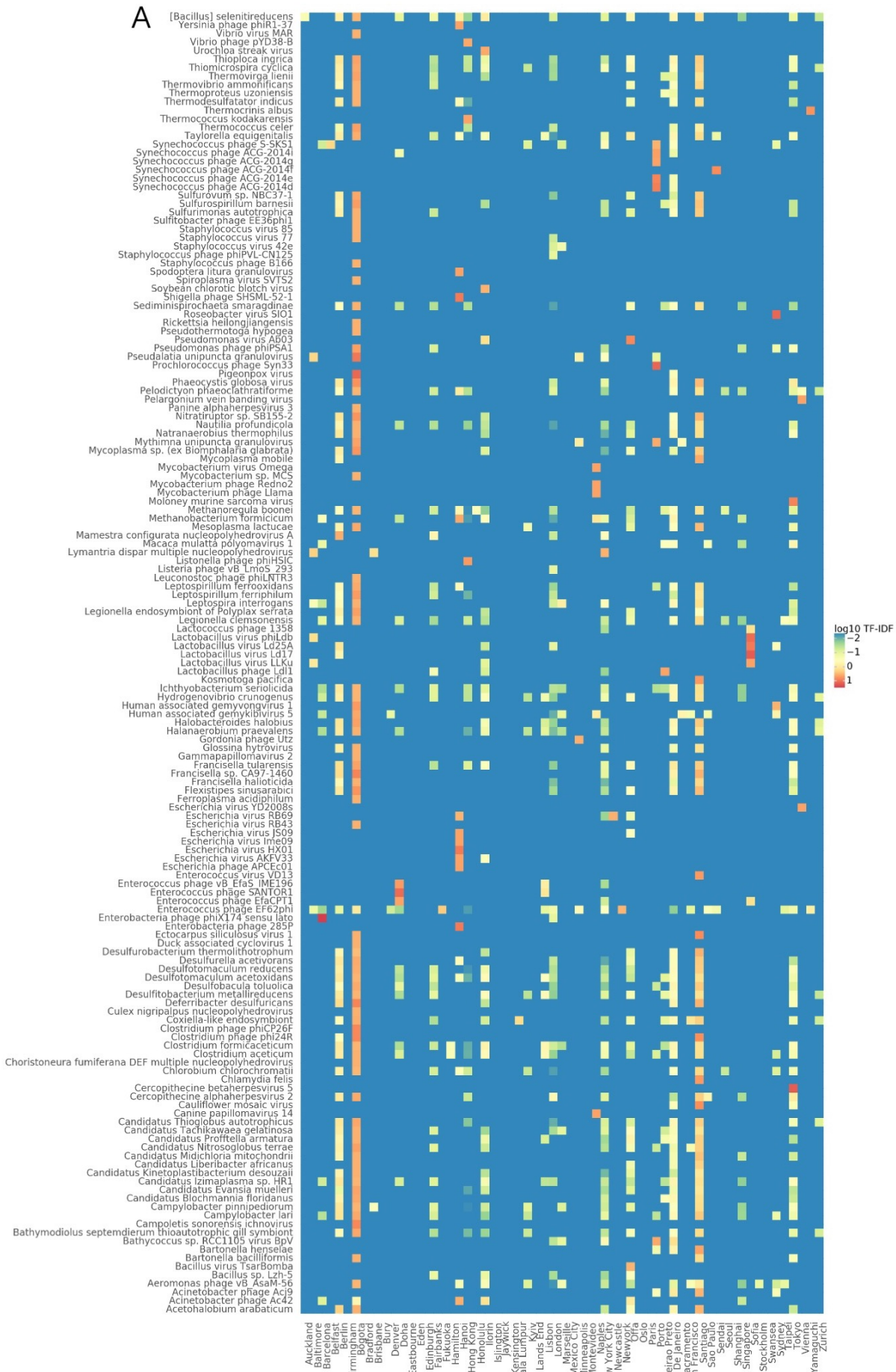


Figure S4: Endemicity scores of particular taxa. A) Heatmap showing the endemicity scores (term-frequency inverse document frequency) for taxa in different cities. This table is filtered to show only taxa with high endemicity scores in at least one city.

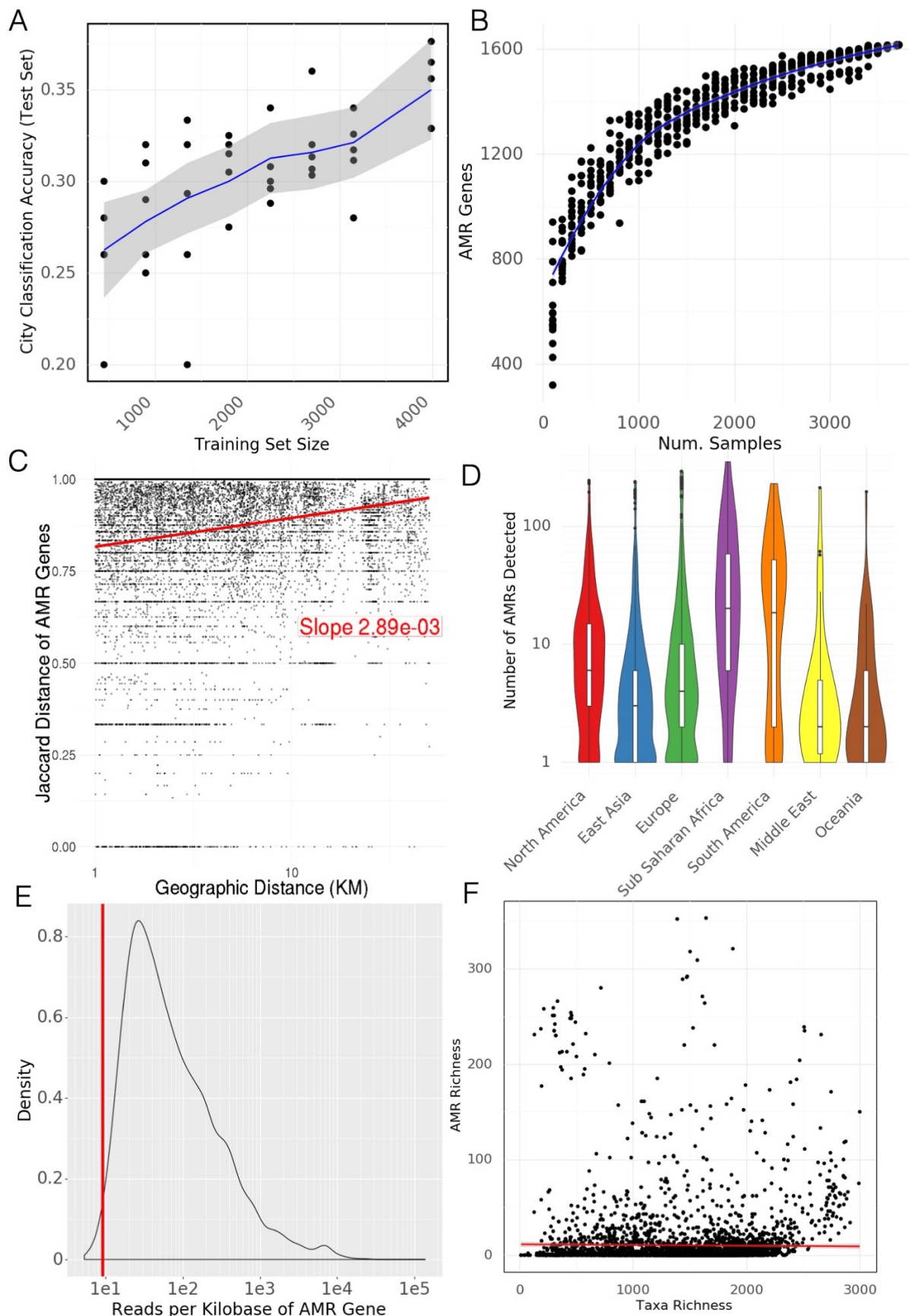


Figure S6: Antimicrobial Resistance Genes, supplemental. A) Classification accuracy of a random forest model predicting city labels for held out samples from antimicrobial resistance genes. B) Rarefaction analysis of antimicrobial resistance genes. Curve does not flatten suggesting we would identify more AMR genes with more samples. C) Neighbourhood effect. Jaccard distance of AMR genes weakly correlates with geographic distance within cities. D) Number of AMR genes detected for samples in each region. E) Distribution of reads per gene (normalized by kilobases of gene length) for AMR gene calls. The vertical red line indicates that 99% of AMR genes have more than 9.06 reads per kilobase and would still be called at a lower read depth.

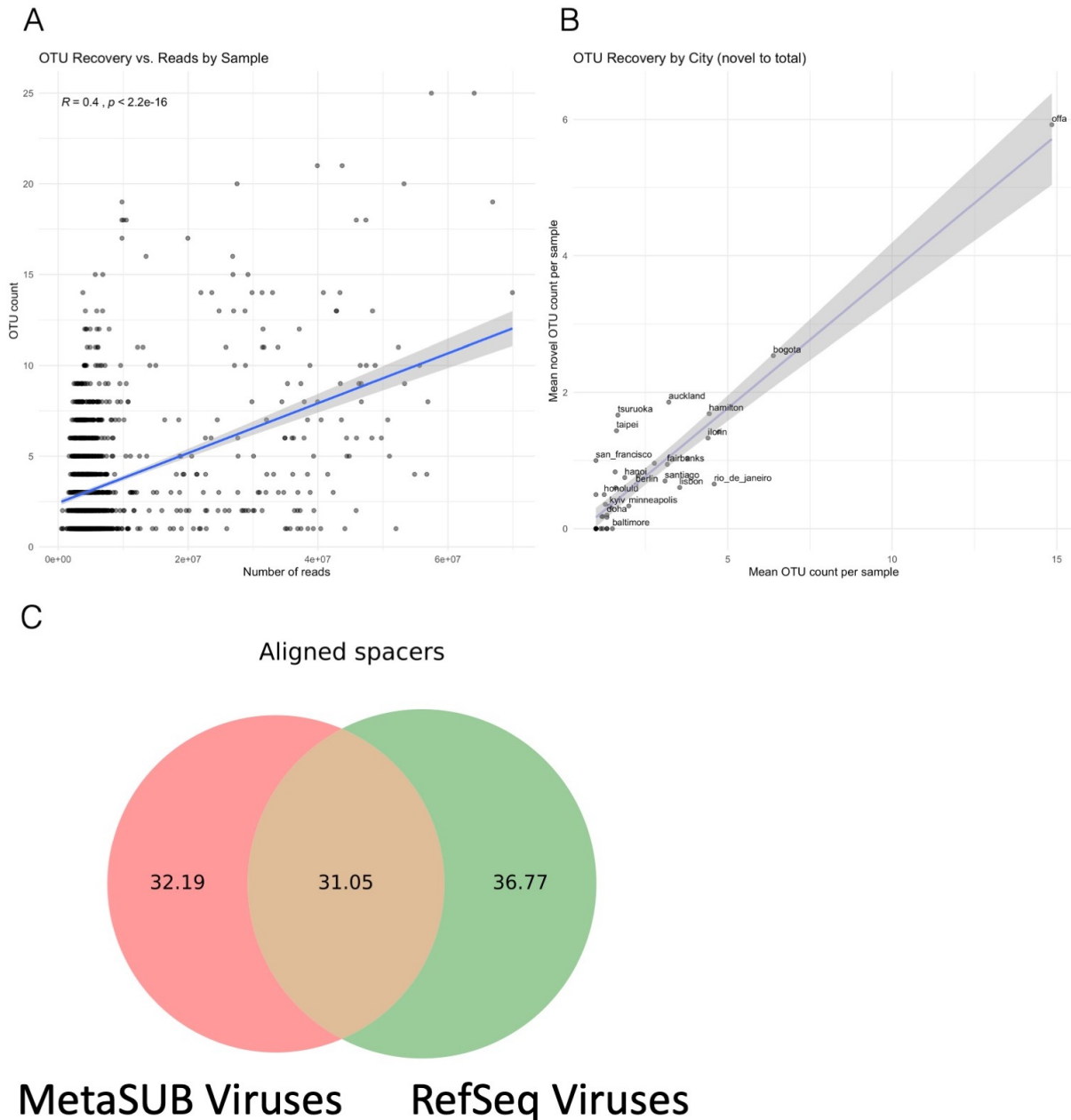


Figure S7: Novel biology, supplemental. A) Relation of read depth to the number of identified bacterial Metagenome Assembled Genomes (MAGs) in a sample. B) Discovery rate for bacterial MAGs in each city. C) Total fraction of CRISPR spacers aligned to MetaSUB viral MAGs and viral genomes in RefSeq.

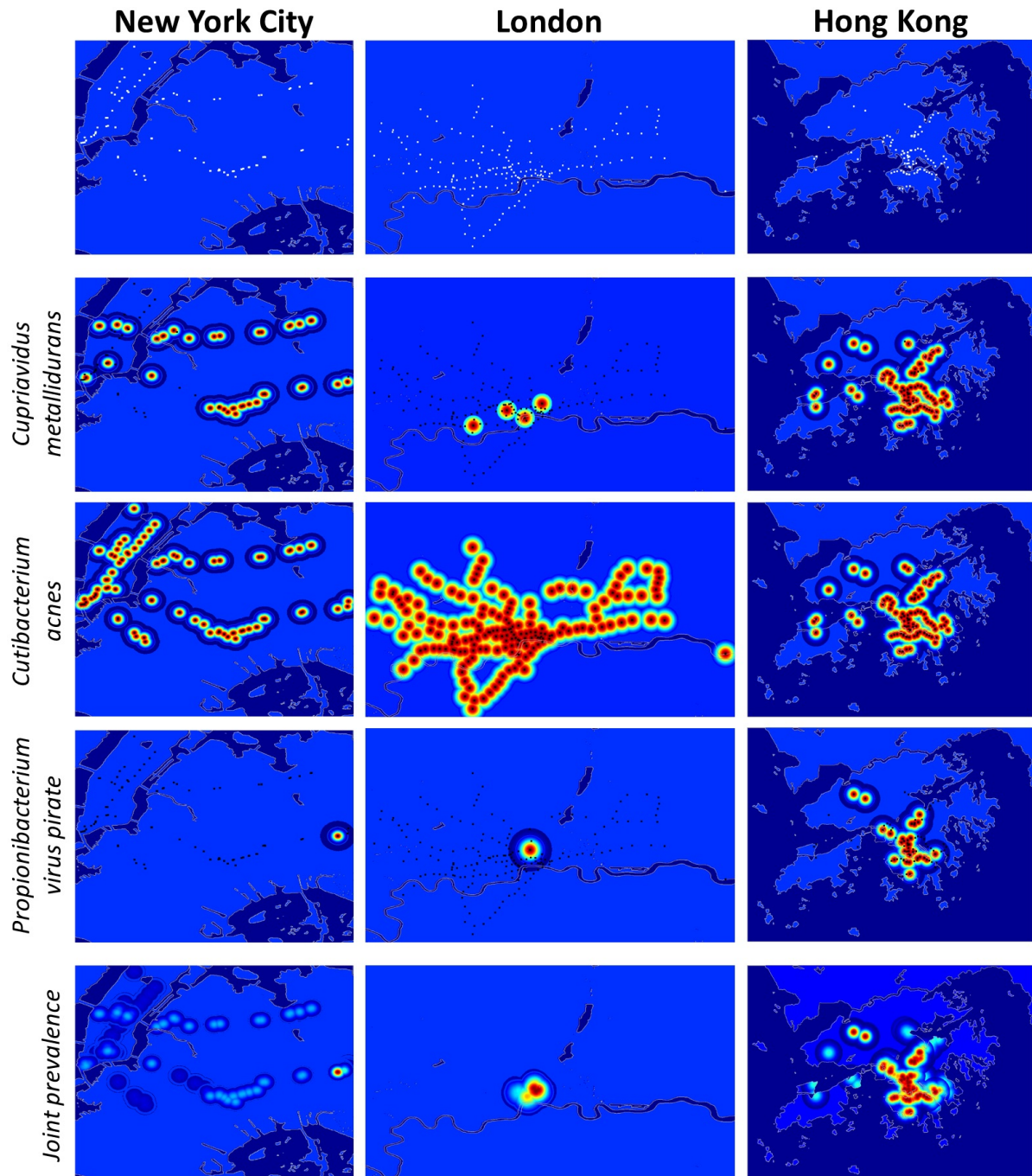


Figure S8: Example Geographic taxonomic Distributions. Distributions of taxa were estimated by fitting Gaussian distributions to sampling locations where the taxa was found with standard deviations based on the geographic distance between observations. Top Row) Sampling sites in three major cities Rows 2-4) Estimated distribution of different example species in major cities Row 5) Estimated distribution of three species together in major cities

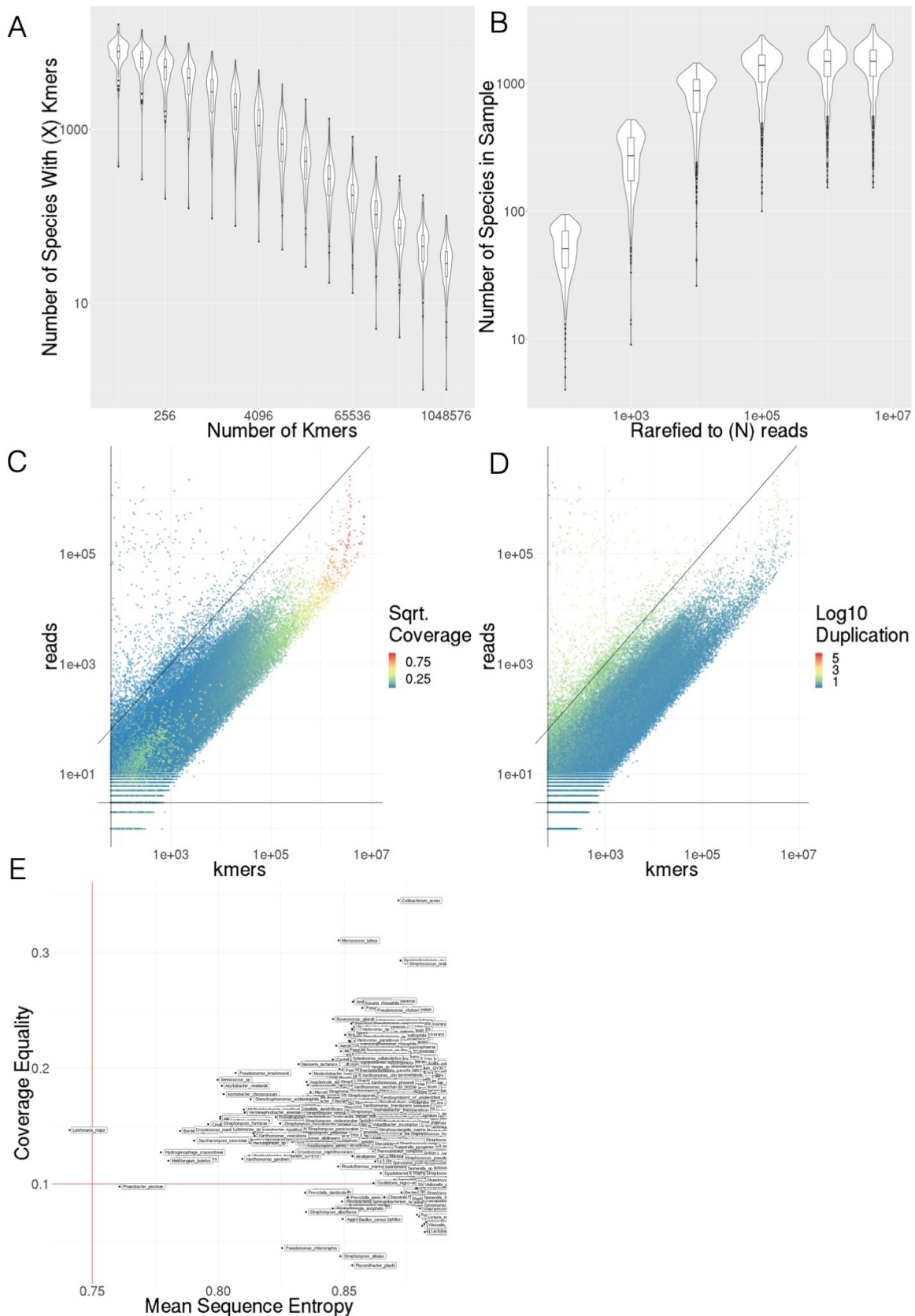


Figure S9: A) Number of species detected as k -mer threshold increases for 100 randomly selected samples B) Number of species detected as number of sub-sampled reads increase C) k -mer counts compared to number of reads for species level annotations in 100 randomly selected samples, colored by coverage of marker k -mer set D) k -mer counts compared to number of reads for species level annotations in 100 randomly selected samples, colored by average duplication of k -mers E) Comparison of Mean Sequence Entropy and Coverage Equality for core and sub-core taxa. Thresholds are shown by red lines.

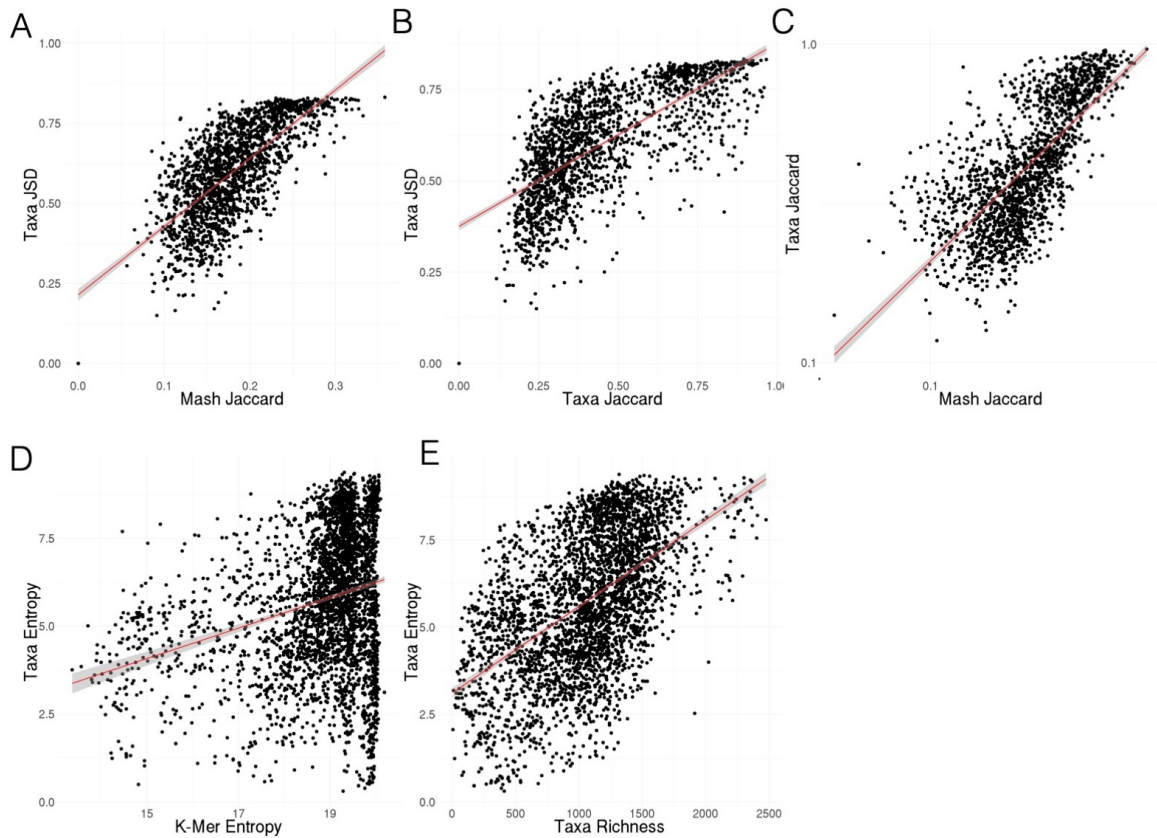


Figure S10: A) Jensen-Shannon Divergence of taxonomic profiles vs MASH Jaccard distance of k -mers B) Jensen-Shannon Divergence of taxonomic profiles vs Jaccard distance of taxonomic profiles. C) Jaccard distance of taxonomic profiles vs MASH Jaccard distance of k -mers D) Shannon's Entropy of taxonomic profiles vs Shannon's Entropy of k -mers E) Taxonomic richness (number of species) vs Shannon's Entropy of taxonomic profiles

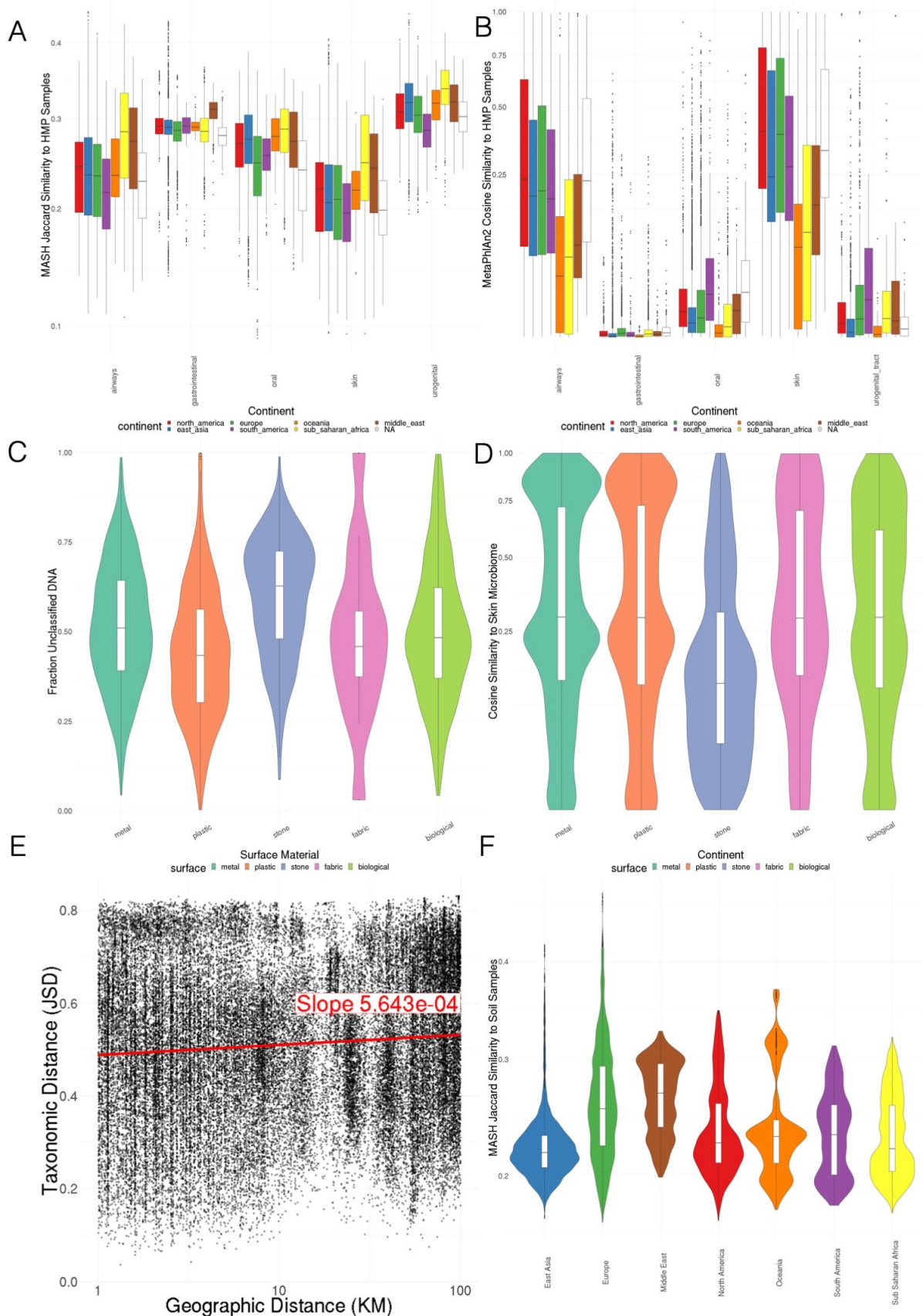


Figure S11: A) MASH *k*-mer Jaccard similarity to representative HMP samples, colored by continent B) MetaPhlan v2.0 cosine similarity to representative HMP samples, colored by continent C) Fraction unclassified DNA by surface material D) Cosine similarity to MetaPhlan v2.0 skin microbiome profile by surface E) Jensen-Shannon distance between pairs of taxonomic profiles vs Geographic Distance F) MASH *k*-mer Jaccard similarity to representative soil samples, colored by continent

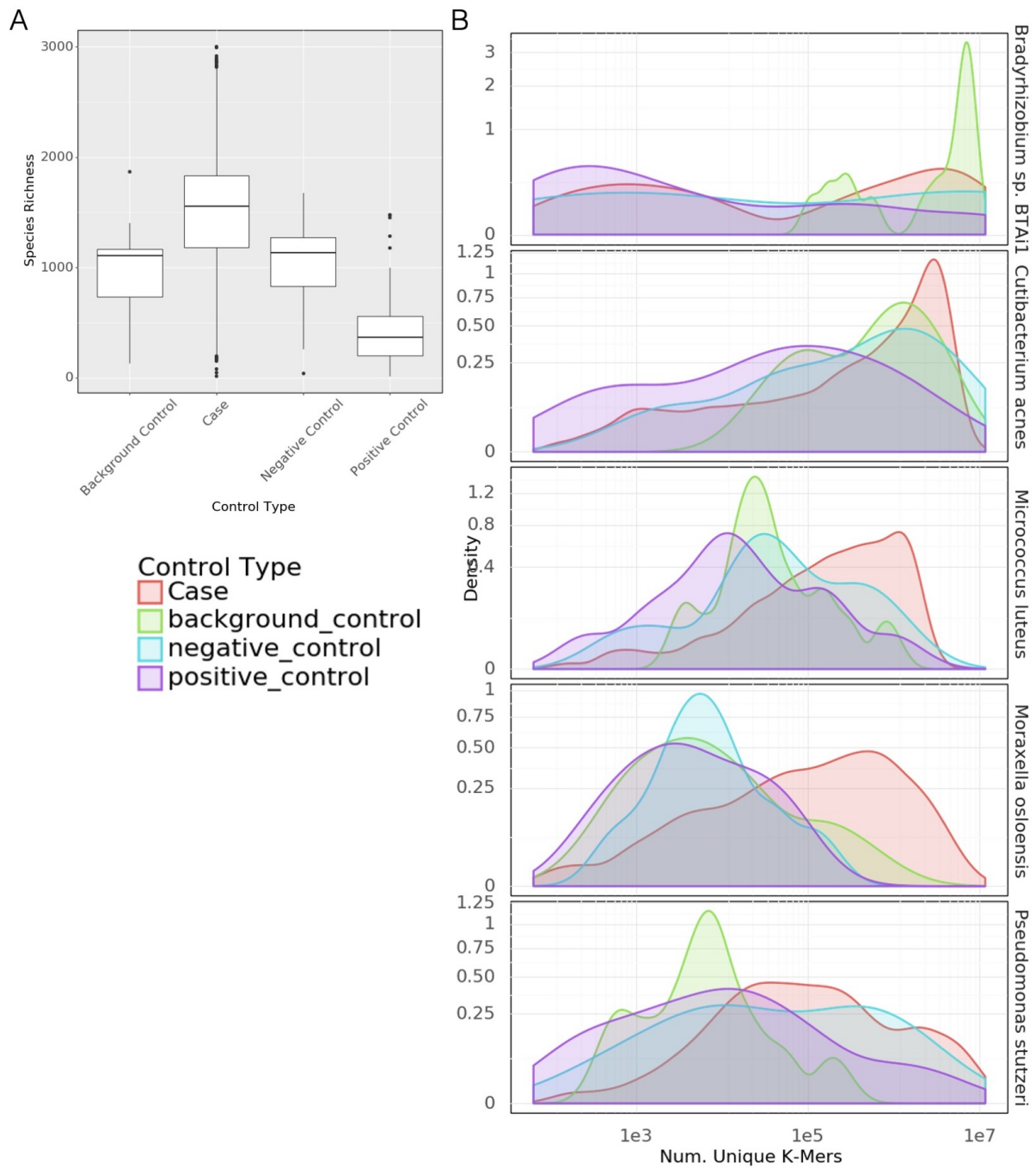


Figure S12: A) Taxonomic Richness in Cases vs. Types of Controls B) Distributions of k -mer counts in control types vs cases for 5 most abundant taxa. k -mer count is a marker of assignment confidence.

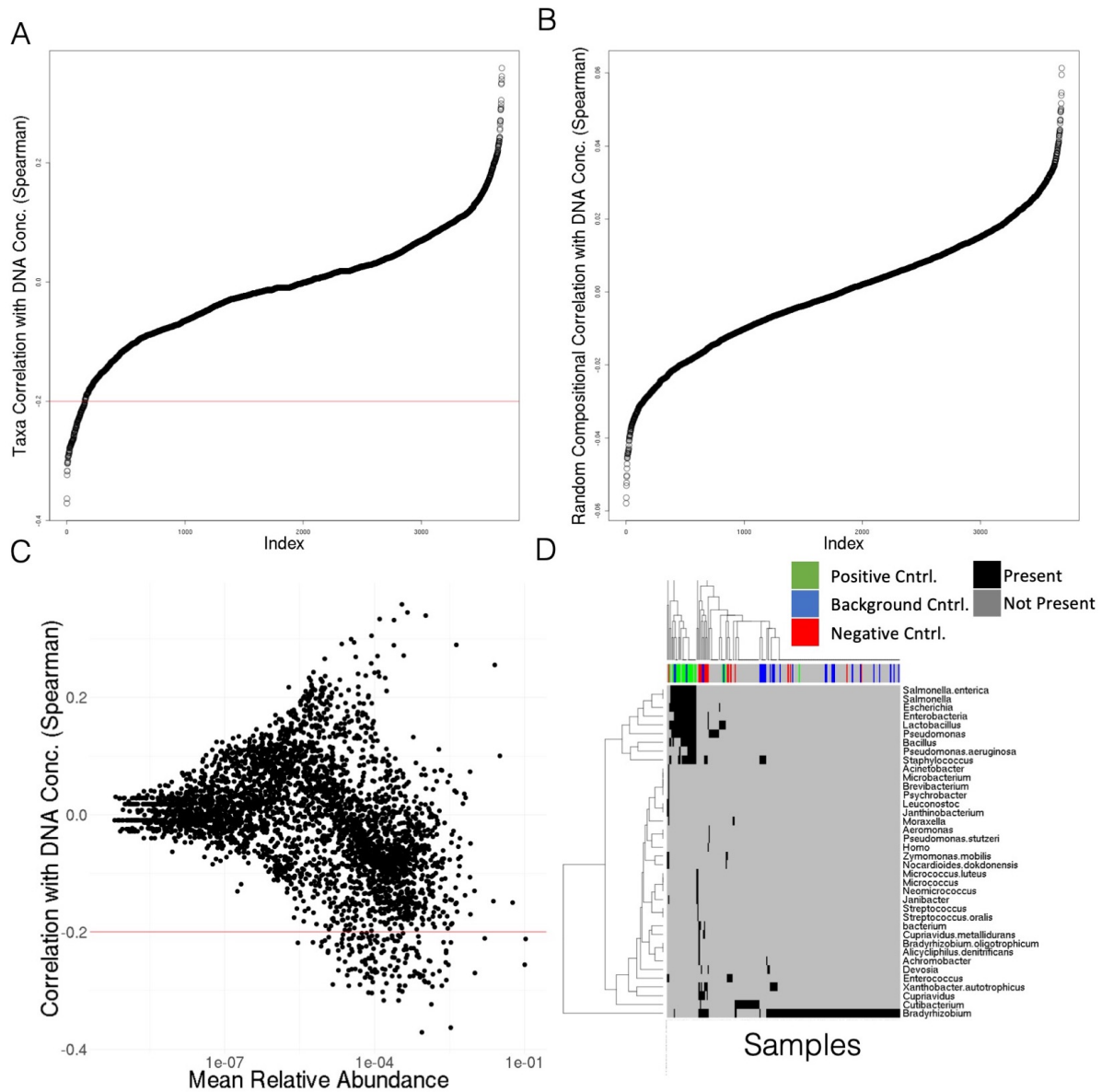


Figure S13: A) Correlation of taxonomic (species) relative abundances with DNA concentration B) Correlation of randomly generated compositional vectors with DNA concentration. Note the same shape but lower magnitude C) Correlation of taxa with DNA Concentration vs the mean relative abundance of that taxa D) Presence (black) absence (grey) heatmap of taxa found in controls and other samples. Colored bar at top, red are negative controls, blue are background, green are positive. Case samples with homology are grey. Case samples without homology to control sequences are not shown.

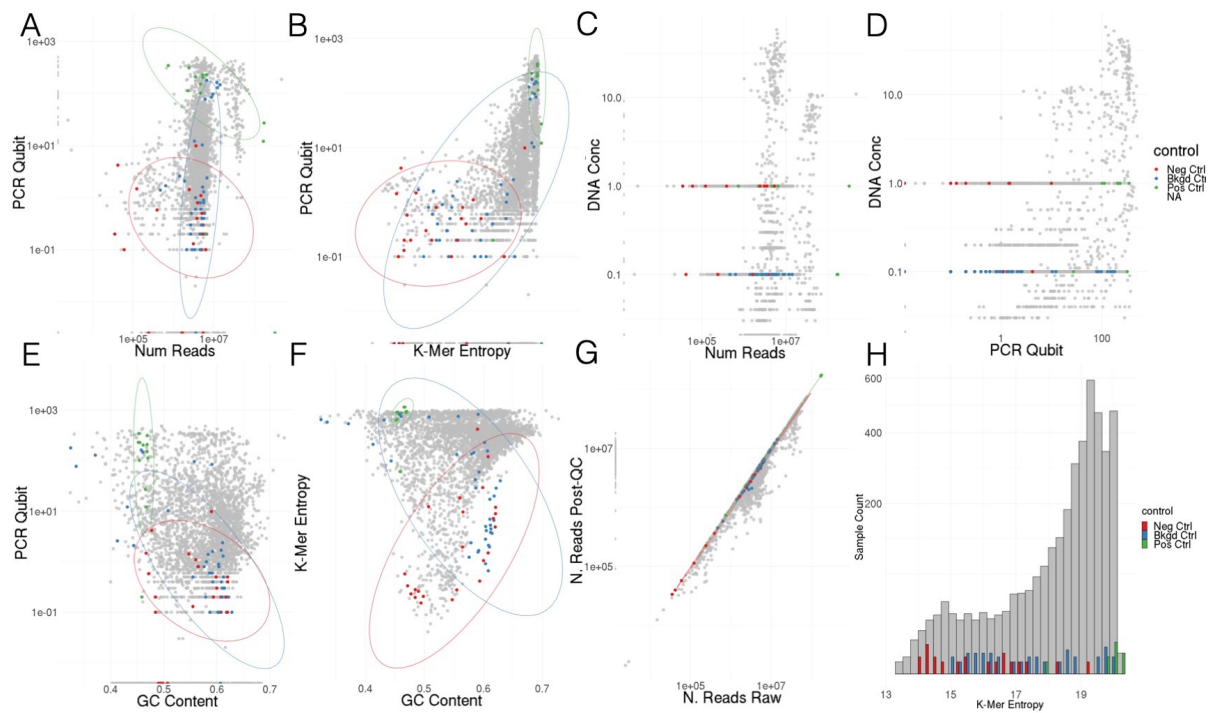


Figure S14: Comparisons of different sequencing quality control metrics with controls marked. A-F) Comparisons of the raw reads, PCR Qubit scores, manually recorded DNA concentrations, *k*-mer Shannon entropy, and GC fraction of quality controlled reads G) Comparison of read counts before and after quality control but before human reads were removed H) Histogram showing the number of samples with different *k*-mer entropies.

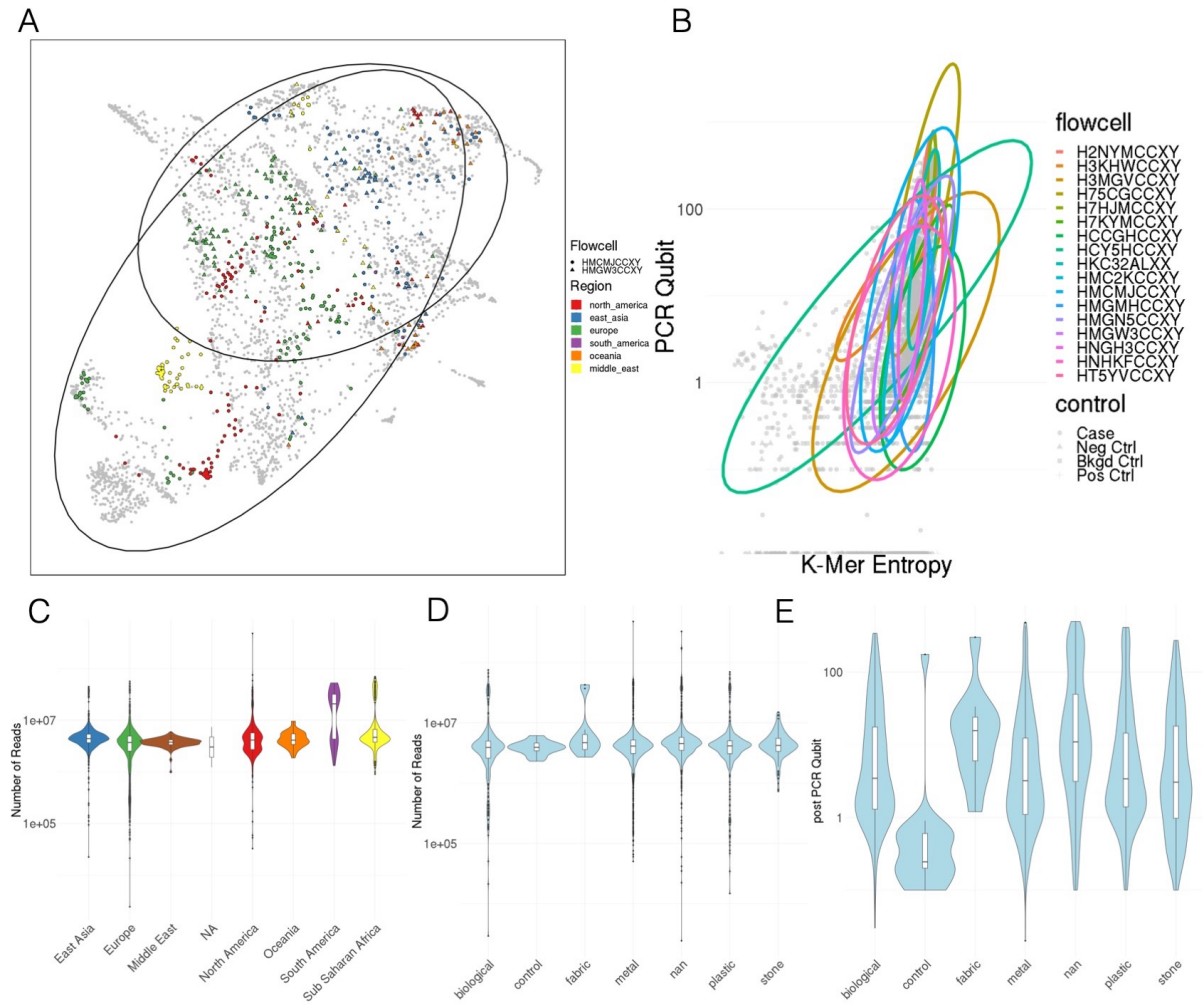


Figure S15: A) UMAP of taxonomic profiles from geographically diverse flowcells B) Flowcells vs quality control metrics C) Number of reads by region D) number of reads by surface material E) PCR Qubit by surface material