

# MICROBIAL SIMILARITY BETWEEN STUDENTS IN A COMMON DORMITORY ENVIRONMENT REVEALS THE FORENSIC POTENTIAL OF INDIVIDUAL MICROBIAL SIGNATURES

Miles Richardson<sup>1,2</sup>, Neil Gottel<sup>3</sup>, Jack A Gilbert<sup>3</sup>, Simon Lax<sup>4</sup>

Correspondence to: M.R. ([miles.richardson@columbia.edu](mailto:miles.richardson@columbia.edu)) & S.L. ([simonlax@mit.edu](mailto:simonlax@mit.edu))

<sup>1</sup>Department of Systems Biology, Columbia University, New York, NY

<sup>2</sup>Integrated Program in Cellular, Molecular, and Biomedical Studies, Columbia University, New York, NY

<sup>3</sup>Department of Pediatrics, University of California San Diego, La Jolla, CA

<sup>4</sup>Physics of Living Systems, Department of Physics, Massachusetts Institute of Technology, Cambridge, MA

## Abstract

1           The microbiota of the built environment is an amalgamation of both human and  
2           environmental sources. While human sources have been examined within single-family households or  
3           in public environments, it is unclear what effect a large number of cohabitating people have on the  
4           microbial communities of their shared environment. We sampled the public and private spaces of a  
5           college dormitory, disentangling individual microbial signatures and their impact on the microbiota of  
6           common spaces. We compared multiple methods for marker gene sequence clustering, and found that  
7           Minimum Entropy Decomposition (MED) was best able to distinguish between the microbial  
8           signatures of different individuals, and was able to uncover more discriminative taxa across all  
9           taxonomic groups. Further, weighted UniFrac- and random forest-based graph analyses uncovered  
10          two distinct spheres of hand or shoe associated samples. For hand-associated samples, connection  
11          between cliques was enriched for hands, implicating them as a primary means of transmission. By  
12          contrast, shoe-associated samples were found to be freely interacting, with individual shoes more  
13          connected to each other than to the floors they interact with. Individual interactions were highly  
14          dynamic, with groups of samples originating from individuals clustering freely with other individuals,

15 while all floor and shoe samples consistently clustered together.

## Importance

16 Humans leave behind a microbial trail, regardless of intention. This may allow for the identification  
17 of individuals based on the ‘microbial signatures’ they shed in built environments. In a shared living  
18 environment, these trails intersect, and through interaction with common surfaces may become  
19 homogenized, potentially confounding our ability to link individuals to their associated microbiota.  
20 We sought to understand the factors that influence the mixing of individual signatures, and how best  
21 to process sequencing data to best tease apart these signatures.

## Introduction

22 Numerous recent studies have uncovered the extent to which humans influence the microbial  
23 ecology of the spaces they occupy through microbial exchange between skin and the built  
24 environment. Most of these studies have focused on home-associated microbial communities(1–3),  
25 with home size, number of occupants, and building materials differentiated between sampling  
26 locations. Each of those confounding factors may have significant impacts on microbial community  
27 structure, and they are difficult to disentangle. Other studies have focused instead on the microbial  
28 ecology of public spaces, such as classrooms and hospital entrance halls(4–8). Although they have  
29 been able demonstrate that most of the taxa colonizing those spaces are skin-associated, they are  
30 unable to link individual human microbial signatures to their data.

31 Microbial flow in the built environment is a keen topic of interest. Cohabitation of multiple  
32 individuals has been shown to influence the microbiota of common spaces, and of the constituents  
33 themselves(1, 7, 9). Common areas may also serve as sites of exchange between individuals, with  
34 implications for disease control. Also unclear are how methodological differences in sequence

35 clustering impact the ability of these studies to link individuals to their surroundings through microbial  
36 similarity.

37 Dorm buildings, which have a standardized architectural design, common building materials  
38 and furnishings between rooms, and even a common ventilation system, represent an intriguing model  
39 system in which to characterize the direct effects of an individual's skin microbiota on their  
40 surroundings, and to further elucidate the forensic potential of skin microbial signatures. Individuals  
41 shed around 30 million bacterial cells per hour(10), and thus leave behind a “microbial fingerprint”.  
42 In one sense, dorm rooms represent a number of identical replicates that can be used to uncover  
43 general patterns of human microbial exchange with the built environment. In a different sense, they  
44 are a “metacommunity” in which it is possible to record a network of interaction by logging visits  
45 between rooms and the use of common spaces. The divide between private rooms and common  
46 spaces such as hallways, lounges, and restrooms further enables us to tease apart individual microbial  
47 signatures in shared spaces.

48 To explore the divide between public and private, we sampled 37 participants from the  
49 University of Chicago's eight floor South Campus residence hall, with four timepoints over 3 months.  
50 Participants were drawn from one “house” in the dormitory, which serves a subset of the dormitory  
51 floor plan with shared common space and bathrooms. From participants, we swabbed both skin sites,  
52 such as hands, and personal effects, such as bed sheets. Additionally, common surfaces such as tables  
53 and bathrooms were also sampled. Together, this collection of surfaces encompasses the divide  
54 between private and public space in the dormitory.

55 To determine how to optimize the inference of individual microbial signatures, we employed  
56 three sequence processing methods to find the most discriminative in characterizing individuals. It has  
57 been observed that in many built environment studies, a large fraction of reads were attributed to a

58 small number of OTUs(1, 9). These OTUs come from a small selection of skin-associated taxonomic  
59 groups, including corynebacteria, staphylococci, pseudomonads, and streptococci. As much of the  
60 differentiation between individuals occurs within a small number of taxonomic groups, it is unclear  
61 how to optimize sequence clustering for forensic inference as OTU clustering may lump together  
62 similar sequences by design. UPARSE(11) is an greedy OTU clustering algorithm for sequence  
63 processing, relying on a greedy clustering method that uses highly abundant sequences as seeds for  
64 clustering. Sequence based methods do not employ OTU clustering and provide a higher resolution  
65 for sequence differentiation. DADA2(12) is a reference-free sequence based algorithm that partitions  
66 reads based on an error model generated from the dataset. Minimum Entropy Decomposition(13)  
67 (MED) is an unsupervised version of oligotyping(14), a method that derives sequences by looking for  
68 regions of high Shannon entropy among 16S regions, and decomposing them into constituent  
69 oligotypes. Oligotyping has been used to explore variation in host associated bacteria, such as in *Blautia*  
70 found in sewage systems(15).

## Results

### Clustering Methodology Impacts the Success of Forensic Inference

71 Each of the sequence processing methods produced a different picture of the microbial  
72 diversity of the dormitory. UPARSE recovered the largest number of distinct sequences (6011) along  
73 with the greatest number of phyla. MED recovered fewer sequences (3353), and fewer phyla (9), but  
74 recovered more members within each phylum (**Supplementary Table 1**). DADA2 recovered nearly  
75 the same phylum level diversity as UPARSE (23 vs 25), but fewer sequences (4307). MED also had a  
76 significantly smaller phylogenetic distance between taxa (Wilcoxon Rank-sum test,  $p < 2.2e-16$ ) than  
77 both DADA2 and UPARSE (**Figure 1a**), indicating that MED recovered much more closely related  
78 sequences.

79           Since we were most interested in classifying individuals, we compared each method using a  
80 random forest trained on surfaces that closely associate with the hands of only one individual, in order  
81 to test their forensic inference. There is a major divide between floor and hand-associated samples  
82 (**Figure 2**). Floor associated samples, including shoes and floors, inhabit a different space compared  
83 to hand associated samples, and this division significantly structures these communities (ANOSIM on  
84 Bray-Curtis Distance  $R=0.2821$ ,  $P=0.001$ ). Thus, to predict which individual's hands a surface had  
85 interacted with, bed sheets, desks, and door handles of the participant rooms are most useful.

86           These models were implemented using an random forest(39), which allows for the  
87 interrogation of similarity between samples. The model was then tested on hand samples from the  
88 same individuals, with the resulting accuracy summarized in **Table 1**. The standardized method of  
89 interpreting the success of classifiers is the error ration, which quantifies how well the random forest  
90 does at predicting the correct individual relative to the success expected by chance(40). An error ratio  
91 above two is commonly used as a significance threshold, and a higher ratio indicates better  
92 performance. All methods performed significantly better than random, but MED clearly  
93 outperformed UPARSE and DADA2 in our dataset. **Supplementary Figure 2** presents the  
94 confusion matrix generated by MED. Samples that fall on the diagonal are correctly classified by the  
95 random forest model. Most (79.57%) fall on the diagonal of the plot. However, for certain individuals,  
96 their hand samples are misclassified in every instance.

97           Interestingly, the largest source of classification error was the presence of roommates in the  
98 room. In fact, the classification error of an individual was linearly related to the number of roommates  
99 that individual had ( $R$ -squared 0.3143,  $P < 0.0001$ ), with classification error increasing by 18  
100 percentage points for each roommate. The relationship is shown in **Supplementary Figure 3**. The  
101 random forest model attempts to use differences in taxa abundance between individuals to classify

102 individuals. If two individuals interact and exchange bacteria, differences in abundance decrease,  
103 which in turn increases model error. Roommates had a significantly smaller weighted UniFrac distance  
104 between them than individuals residing in different rooms. (Wilcoxon Rank-sum test,  $W = 409660000$ ,  
105  $p < 2.2 \times 10^{-16}$ )

### Classification of Individuals is Driven by Specific Taxa

106 The random forest model is able to rank individual sequences or OTUs by their importance  
107 in successful classification. As seen in **Figure 1**, there are differences in the distribution and average  
108 importance score across phyla, and all of these are significant to .05 by Kruskal-Wallis test.  
109 Furthermore, MED has a significantly higher importance score in all phyla that overlap between all  
110 three methods except for Cyanobacteria, Fusobacteria, and Deinococcus-Thermus (Wilcoxon Rank-  
111 Sum Test, FDR  $p < .05$ ) (**Supplementary Figure 4**).

112 It has been noted that there are taxa indicative of different sexes.<sup>(41)</sup> To see if there were enriched  
113 taxa between men and women from room samples, we looked for differentially enriched taxa using  
114 DESeq2. The most significantly enriched taxon is *Lactobacillus iners*, an inhabitant of the female  
115 reproductive tract. Certain corynebacteria were also noted to be more abundant in men, as seen in  
116 **Supplementary Figure 5**. Using these enriched taxa, we used the random forest to predict whether  
117 a subject is a man or woman, with an error ratio of about 2.5.

### Metacommunity Structure

118 In addition to classifying individuals, we sought to recapitulate the geographical structure of  
119 the dorm using graphical models. To do this, we constructed a threshold graph of the weighted  
120 UniFrac distance between samples, with a threshold of 0.1. As seen in **Figure 3**, the dorm has two  
121 large subgraphs, along with a number of orphaned graphs. These two groups consist of floor-  
122 associated (shoes and floors) and hand-associated (hand, doorknob, bed, and desk) samples. The

123 orphaned graphs are mostly samples from one individual. As expected, common surfaces in the hand-  
124 associated realm serve as an anchor for their subgraph, connecting a number of different people, while  
125 hallway floors serve the same role for individual shoes. By contrast, orphaned graphs appear to  
126 indicate the stability of an individual's microbial signature over time and a lack of interaction with  
127 other samples.

128 Further, we can also calculate the assortativity of different metadata criteria. Assortativity is a  
129 metric used to quantify how often a node attaches to a similar node, with higher assortativity reflecting  
130 higher connectivity between similar nodes. As seen in Table 2, sex and floor have the highest  
131 assortativity, while timepoint is the least important. This indicates that floor and sex are more  
132 important in generating the graph structure, and implies that the microbial signature of individual  
133 surfaces across the sampling period is stable.

134 While a graph can be constructed using a beta-diversity metric (in our case weighted UniFrac  
135 distance) as above, the distance metric may not be sensitive to the microbial community of an  
136 individual. Since there is information to be gained from aggregating samples into a larger individual  
137 signature, we also constructed a graph using random forest proximity. The proximity values from the  
138 random forest are akin to a distance, and take into account the same signature used to classify  
139 individuals. It is also much sparser, as the random forest is trying to minimize distances between  
140 samples from the same individual, while keeping samples between individuals distinct. The resulting  
141 graph can be seen in **Supplementary Figure 6a and 6b**, where samples are colored by individual and  
142 surface type, respectively.

143 Graph based clustering analysis methods are often used in describing interactions in social  
144 networks. Using the Infomap clustering algorithm(42), which uses flow within a network to generate  
145 groupings, we looked at how bacterial exchange grouped our samples. The Infomap algorithm is also

146 hierarchical(43), allowing for samples to inhabit “Top Modules” which are large scale groupings, and  
147 then submodules that indicate community clustering within top modules. Using this algorithm, we  
148 identified 8 top modules (**Figure 4**), with module 1 encompassing almost all shoe and floor samples.  
149 Of particular interest were how samples grouped over time. Further, among surfaces connecting top  
150 modules, hands were significantly enriched, and other hand-associated surfaces showed enrichment,  
151 including common tables, doors, and bathroom doors. (**Supplementary Figure 7**) Since the dorm  
152 represents a multilayer graph, where each timepoint forms a distinct layer of interaction, we employed  
153 a multilayer implementation of the algorithm(44) to look at the stability of interactions over time. This  
154 is presented in **Figure 5**, where samples are clustered at each timepoint and their membership in  
155 clusters is tracked over time. Shoe and floor samples showed high stability over time, where most  
156 samples co-cluster over time in the same clusters (**Figure 5a**). Similarly, common surfaces had the  
157 same pattern, wherein common floor samples clustered consistently, while common hand samples  
158 could be affiliated with different samples (**Figure 5b**). Individual samples, for example individual 1  
159 and 29 (**Figure 5c**), showed the ability to cluster freely with other samples.

## Discussion

160 The use of human microbial signatures as trace evidence remains a young and inexact science.  
161 In order for this developing field to become a useful forensic tool, methods will need to be optimized  
162 and the myriad factors which influence our microbial interaction with built environments will need to  
163 be disentangled. Here, we compared classification methods to link residents to their rooms and  
164 personal effects in a common dormitory environment. For classifying individuals, Minimum Entropy  
165 Decomposition seems to be the best choice, but it appears that exact sequence variants in general are  
166 better at identifying individuals than OTU-clustering methods. This is unsurprising, as exact sequence  
167 variants maximize our insight into the microbial strains that differ between individuals that can be



168 obscured by higher-level OTU-clustering. At the same time, this contrasts with observations that  
169 MED can produce incorrect sequence variants from mock communities. The advantage of MED  
170 seems to be that it is able to recover more diversity within the main skin-associated taxa from the  
171 phyla Proteobacteria, Fusobacteria, Bacteroidetes, Fusobacteria, and Actinobacteria. It is also able to  
172 recover higher importance scores even at the genus level, indicating that it is able to produce more  
173 individual-specific sequences within common skin taxa, as the importance score only measures the  
174 usefulness in classification between individuals.

175         The high accuracy of classification shows that skin associated samples, in particular bed sheets  
176 and door handles are useful in predicting the individuals inhabiting those rooms. Furthermore,  
177 samples within 2-4 weeks also appear useful in prediction, indicating that the signature is stable over  
178 long time scales. It is confounded by roommates, which is unsurprising, due to microbial exchange  
179 within one room.

180         This study explores the metacommunity structure of the college dormitory, which clusters into  
181 two distinct spheres- hand or shoe associated samples. Within each type, the arrangement between  
182 common and personal samples differs, with personal shoe associated samples freely associating, while  
183 hand associated sample often only associate between individuals when connected by a common  
184 surface. When examining graphs generated by random forest, the shoe associated structure remains  
185 the same, while individual signatures closely cluster with themselves. A multilayer graph reveals that  
186 individual signatures freely intermingle at different timepoints, while shoe and floor samples have large  
187 continuous interaction. This is likely a result of the high exchange between shoes and floors, which  
188 homogenizes the signature of shoe and floor samples.

## Materials and Methods

### Study Design and Sample Collection

189 We collected personal samples from 37 participants in 28 distinct dorm rooms  
190 (**Supplementary Table** ). Samples were collected by swabbing a sterile cotton BD-Swube applicator  
191 against the dry surface of interest. Sampling kits were given to study participants for self-sampling  
192 with instructions. The desk, floor, fitted bed sheet, and interior doorknobs of each participant's room,  
193 along with the dominant hand and shoe of the participant, were sampled at four timepoints. The first  
194 timepoint occurred before occupants left for a scheduled school break (end of a quarter) and then  
195 immediately upon return. The third and fourth timepoints were taken 2 and 4 weeks after spring break.

196 Participants also completed a questionnaire which collected basic information on the subject,  
197 the conditions specific to their dorm room, and who they interacted with in their dorm room during  
198 the sampling period. This questionnaire was completed each time a set of samples was collected.

199 Common surfaces were also sampled similarly. Common surfaces specific to the 5<sup>th</sup> floor  
200 included tables in the dormitory lounge, and the handle of the entry door to the lounge. On each floor  
201 of the dormitory, the door handles of bathrooms, the floors of each hallway, and the elevator buttons  
202 were sampled. Each floor had its own unique combination, and these were swabbed at the same time  
203 as personal surfaces.

### **Sample Processing**

204 DNA was extracted from each sample using a low biomass variation of the MO BIO Powersoil DNA  
205 extraction protocol. 16s rRNA was amplified with the Earth Microbiome 16S Illumina Amplicon  
206 Protocol(16). The V4 region of the 16s rRNA gene was targeted with the 515F-806RB primer pair.  
207 Sequencing was performed using a Illumina Miseq sequencer with the protocol described in Caporaso  
208 et al. 2012(17).

### **Sequence Processing**

209 Each method was processed using the default workflows provided in reference papers.

210 UPARSE

211 Demultiplexed sequences were merged using vsearch(18) with 10,040,708 successful paired end reads  
212 merged together. Sequences were quality filtered with a maximum expected error of 0.5, with  
213 9,057,613 remaining sequences. Sequences were then dereplicated for 1,276,202 unique sequences.  
214 Sequences were then clustered at 97% identity, with 11658 OTUs and 42539 chimeras. Sequences  
215 were then matched to OTUs with 93.28% of sequences matched to OTUs. 6011 OTUs passed  
216 sequence processing. Chloroplast and mitochondrial DNA was removed, and samples were rarefied  
217 to 4000 counts per sample.

218 MED

219 Sequences were processed according to the methods described in Meren et al 2015(13). Demultiplexed  
220 paired-end reads were merged using illumina-utils(19), with Q30 check imposed on sequences, leading  
221 to 10,023,266 successfully merged out of 10,023,266 reads. Gaps between sequences were padded  
222 with blanks, and samples were decomposed using a -M of 100. 1,732,615 outliers were removed by  
223 quality control, and remaining sequences were sorted into 3,748 nodes after refinement. 3,352 passed  
224 quality control. Chloroplast and mitochondrial DNA was removed, and samples were rarefied to 4000  
225 counts per sample.

226 DADA2

227 The filtering step of DADA2 was run with no ambiguous base (maxN of 0), maximum expected errors  
228 of 2, quality of truncation of 2. All other commands were run on default settings. Sequences were  
229 merged after performing quality filtering. After merging, 34043 sequences were observed, and 18329

230 sequences were not chimeras. 4307 unique sequences passed final quality filtering. Chloroplast and  
231 mitochondrial DNA was removed, and samples were rarefied to 4000 counts per sample.

### **Taxonomic identification**

232 All sequences were taxonomically identified using the same implementation of RDP(20) implemented  
233 in DADA2 to enable comparison between the sequencing methods. Taxonomy was assigned using  
234 the SILVA(21) training set version 123.

### **Phylogenetic Trees**

235 Sequences were aligned with the R package *MSA*(22), using the Muscle(23, 24) algorithm.  
236 Phylogenetic trees were then generated using the R package *Phangorn*(25). The tree was created by  
237 neighbor joining, and fitted with GTR model.

### **Data Analysis and Visualization**

238 Data cleaning and shaping was performed using R 3.3.2-R3.3.5 and the packages *dplyr*(26)  
239 *reshape2*(27). Visualization and analysis were performed using *phyloseq*(28), *igraph*(29), *ggnetwork*(30),  
240 and *ggplot2*(31). Random forests were generated using *randomForest*(32) and *ranger*(33). Differential  
241 abundance calculation were performed using *DESeq2*(34). Diversity measures were calculated using  
242 *vegan*(35). Inspiration was taken from Callahan et al. 2017(36). Community clustering was performed  
243 using the *InfoMap*(37, 38) and alluvial diagrams generated using the “Map & Alluvial Generator”  
244 (<http://www.mapequation.org/apps/MapGenerator.html>).

### **Acknowledgements**

245 This work was sponsored by National Institutes of Justice award 2015-DN-BX-K430.  
246 Sophia Weaver for indispensable assistance in recruiting individuals.  
247 Study IRB Number: IRB15-0373. Approved by BSD IRB Committee A, The University of Chicago  
248 Biological Sciences Division/University of Chicago Medical Center.

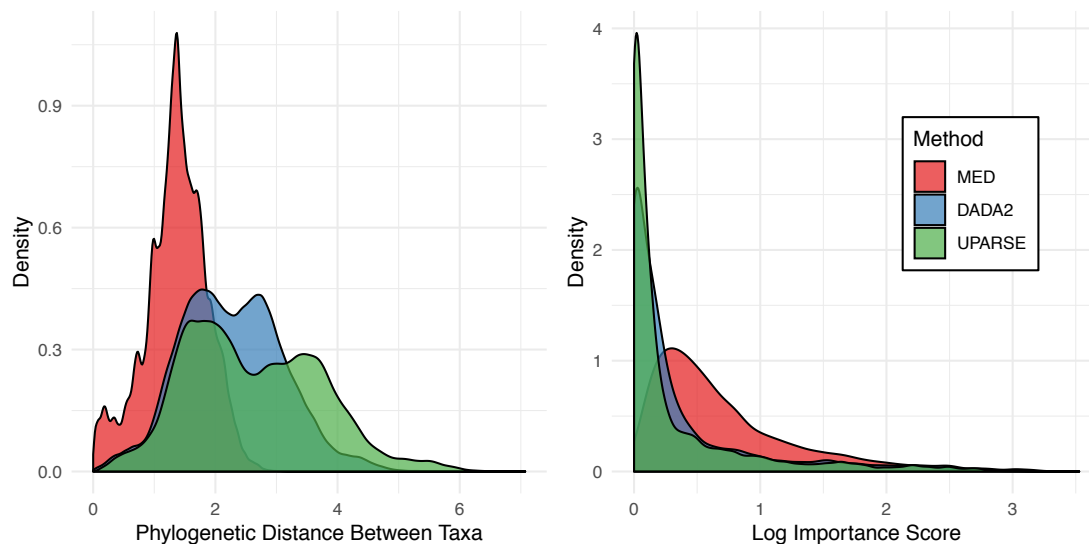
## References

1. Lax S, Smith DP, Hampton-Marcell J, Owens SM, Handley KM, Scott NM, Gibbons SM, Larsen P, Shogan BD, Weiss S, Metcalf JL, Ursell LK, Vázquez-Baeza Y, Treuren WV, Hasan NA, Gibson MK, Colwell R, Dantas G, Knight R, Gilbert JA. 2014. Longitudinal analysis of microbial interaction between humans and the indoor environment. *Science* 345:1048–1052.
2. Dunn RR, Fierer N, Henley JB, Leff JW, Menninger HL. 2013. Home Life: Factors Structuring the Bacterial Diversity Found within and between Homes. *PLOS ONE* 8:e64133.
3. Flores GE, Bates ST, Caporaso JG, Lauber CL, Leff JW, Knight R, Fierer N. 2013. Diversity, distribution and sources of bacteria in residential kitchens. *Environ Microbiol* 15:588–596.
4. Meadow JF, Altrichter AE, Kembel SW, Kline J, Mhuireach G, Moriyama M, Northcutt D, O'Connor TK, Womack AM, Brown GZ, Green JL, Bohannon BJM. 2014. Indoor airborne bacterial communities are influenced by ventilation, occupancy, and outdoor air source. *Indoor Air* 24:41–48.
5. Poza M, Gayoso C, Gómez MJ, Rumbo-Feal S, Tomás M, Aranda J, Fernández A, Bou G. 2012. Exploring Bacterial Diversity in Hospital Environments by GS-FLX Titanium Pyrosequencing. *PLOS ONE* 7:e44105.
6. Kembel SW, Jones E, Kline J, Northcutt D, Stenson J, Womack AM, Bohannon BJ, Brown GZ, Green JL. 2012. Architectural design influences the diversity and structure of the built environment microbiome. *ISME J* 6:1469–1479.
7. Wood M, Gibbons SM, Lax S, Eshoo-Anton TW, Owens SM, Kennedy S, Gilbert JA, Hampton-Marcell JT. 2015. Athletic equipment microbiota are shaped by interactions with human skin. *Microbiome* 3:25.
8. Lax S, Hampton-Marcell JT, Gibbons SM, Colares GB, Smith D, Eisen JA, Gilbert JA. 2015. Forensic analysis of the microbiome of phones and shoes. *Microbiome* 3:21.
9. Lax S, Sangwan N, Smith D, Larsen P, Handley KM, Richardson M, Guyton K, Krezalek M, Shogan BD, Defazio J, Flemming I, Shakhsheer B, Weber S, Landon E, Garcia-Houchins S, Siegel J, Alverdy J, Knight R, Stephens B, Gilbert JA. 2017. Bacterial colonization and succession in a newly opened hospital, Colonization and Succession of Hospital-Associated Microbiota. *Sci Transl Med* 9, 9.
10. Qian J, Hospodsky D, Yamamoto N, Nazaroff WW, Peccia J. 2012. Size-resolved emission rates of airborne bacteria and fungi in an occupied classroom. *Indoor Air* 22:339–351.
11. Edgar RC. 2013. UPARSE: highly accurate OTU sequences from microbial amplicon reads. *Nat Methods* 10:996–998.
12. DADA2: High-resolution sample inference from Illumina amplicon data : Nature Methods : Nature Research.
13. Eren AM, Morrison HG, Lescault PJ, Reveillaud J, Vineis JH, Sogin ML. 2015. Minimum entropy decomposition: Unsupervised oligotyping for sensitive partitioning of high-throughput marker gene sequences. *ISME J* 9:968–979.
14. Eren AM, Maignien L, Sul WJ, Murphy LG, Grim SL, Morrison HG, Sogin ML. 2013. Oligotyping: differentiating between closely related microbial taxa using 16S rRNA gene data. *Methods Ecol Evol* 4:1111–1119.
15. Eren AM, Sogin ML, Morrison HG, Vineis JH, Fisher JC, Newton RJ, McLellan SL. 2015. A single genus in the gut microbiome reflects host preference and specificity. *ISME J* 9:90–100.
16. 16S Illumina Amplicon Protocol : Earth Microbiome Project.

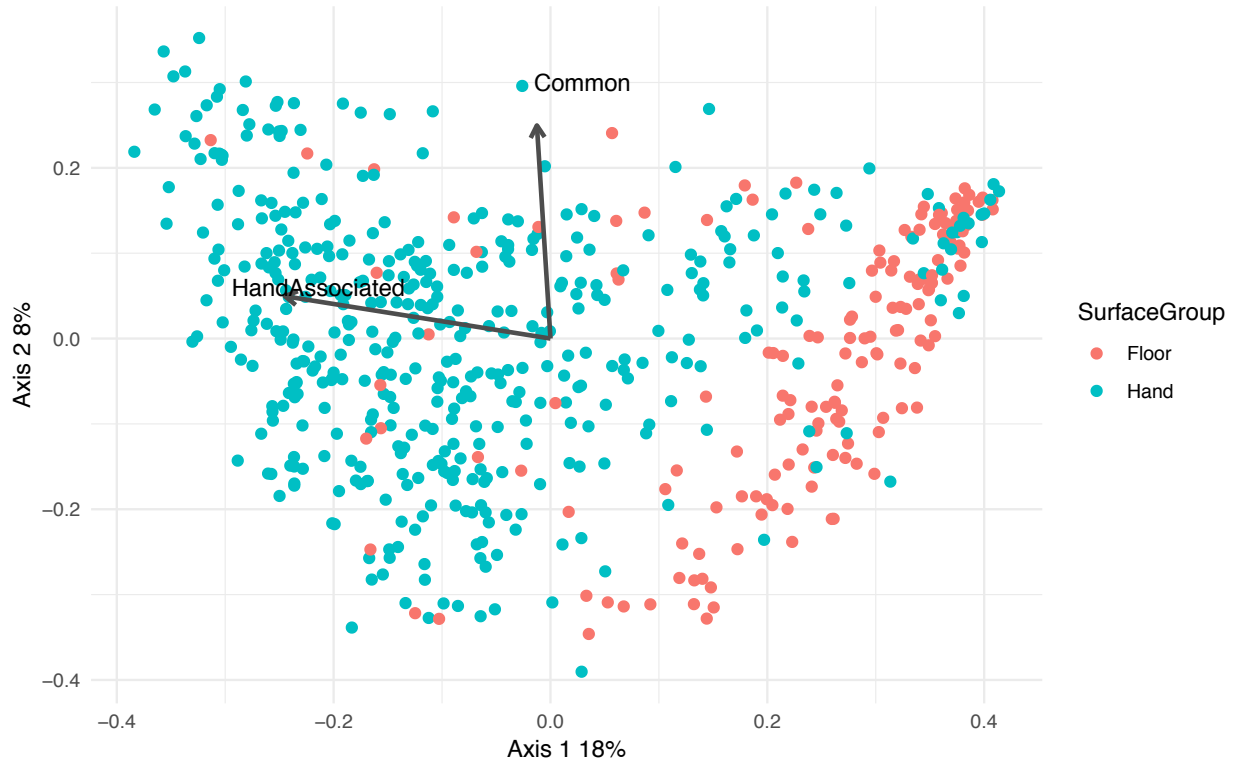
17. Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, Owens SM, Betley J, Fraser L, Bauer M, Gormley N, Gilbert JA, Smith G, Knight R. 2012. Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. *ISME J* 6:1621–1624.
18. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. 2016. VSEARCH: a versatile open source tool for metagenomics. *PeerJ* 4:e2584.
19. Eren AM, Vineis JH, Morrison HG, Sogin ML. 2013. A Filtering Method to Generate High Quality Short Reads Using Illumina Paired-End Technology. *PLOS ONE* 8:e66643.
20. Wang Q, Garrity GM, Tiedje JM, Cole JR. 2007. Naïve Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy. *Appl Env Microbiol* 73:5261–5267.
21. Yilmaz P, Parfrey LW, Yarza P, Gerken J, Pruesse E, Quast C, Schweer T, Peplies J, Ludwig W, Glöckner FO. 2014. The SILVA and “All-species Living Tree Project (LTP)” taxonomic frameworks. *Nucleic Acids Res* 42:D643–D648.
22. Bodenhofer U, Bonatesta E, Horejš-Kainrath C, Hochreiter S. 2015. msa: an R package for multiple sequence alignment. *Bioinformatics* 31:3997–3999.
23. Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* 32:1792–1797.
24. Edgar RC. 2004. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics* 5:113.
25. Schliep KP. 2011. phangorn: phylogenetic analysis in R. *Bioinformatics* 27:592–593.
26. Wickham H, François R, Henry L, Müller K. 2018. dplyr: A Grammar of Data Manipulation.
27. Wickham H. 2007. Reshaping Data with the reshape Package. *J Stat Softw* 21:1–20.
28. McMurdie PJ, Holmes S. 2013. phyloseq: An R Package for Reproducible Interactive Analysis and Graphics of Microbiome Census Data. *PLOS ONE* 8:e61217.
29. Csardi G, Nepusz T. The igraph software package for complex network research. *InterJournal Complex Syst* 1695:9.
30. Briatte F. 2016. ggnetwork: Geometries to Plot Networks with “ggplot2.”
31. Wickham H. 2011. ggplot2. *Wiley Interdiscip Rev Comput Stat* 3:180–185.
32. Liaw A, Wiener M. 2002. Classification and Regression by randomForest 2:6.
33. ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R | Wright | *Journal of Statistical Software*.
34. Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* 15:550.
35. Dixon P, Palmer MW. 2003. VEGAN, a package of R functions for community ecology. *J Veg Sci* 14:927–930.
36. Callahan BJ, Sankaran K, Fukuyama JA, McMurdie PJ, Holmes SP. 2016. Bioconductor Workflow for Microbiome Data Analysis: from raw reads to community analyses. *F1000Research* 5:1492.
37. Rosvall M, Axelsson D, Bergstrom CT. 2009. The map equation. *Eur Phys J Spec Top* 178:13–23.

38. Bohlin L, Edler D, Lancichinetti A, Rosvall M. 2014. Community Detection and Visualization of Networks with the Map Equation Framework, p. 3–34. *In* Ding, Y, Rousseau, R, Wolfram, D (eds.), *Measuring Scholarly Impact: Methods and Practice*. Springer International Publishing, Cham.
39. Shi T, Horvath S. 2006. Unsupervised Learning With Random Forest Predictors. *J Comput Graph Stat* 15:118–138.
40. Knights D, Costello EK, Knight R. 2011. Supervised classification of human microbiota. *FEMS Microbiol Rev* 35:343–359.
41. Luongo JC, Barberán A, Hacker-Cary R, Morgan EE, Miller SL, Fierer N. 2017. Microbial analyses of airborne dust collected from dormitory rooms predict the sex of occupants. *Indoor Air* 27:338–344.
42. Rosvall M, Bergstrom CT. 2008. Maps of random walks on complex networks reveal community structure. *Proc Natl Acad Sci* 105:1118–1123.
43. Multilevel Compression of Random Walks on Networks Reveals Hierarchical Organization in Large Integrated Systems.
44. Aslak U, Rosvall M, Lehmann S. 2018. Constrained information flows in temporal networks reveal intermittent communities. *Phys Rev E* 97:062312.

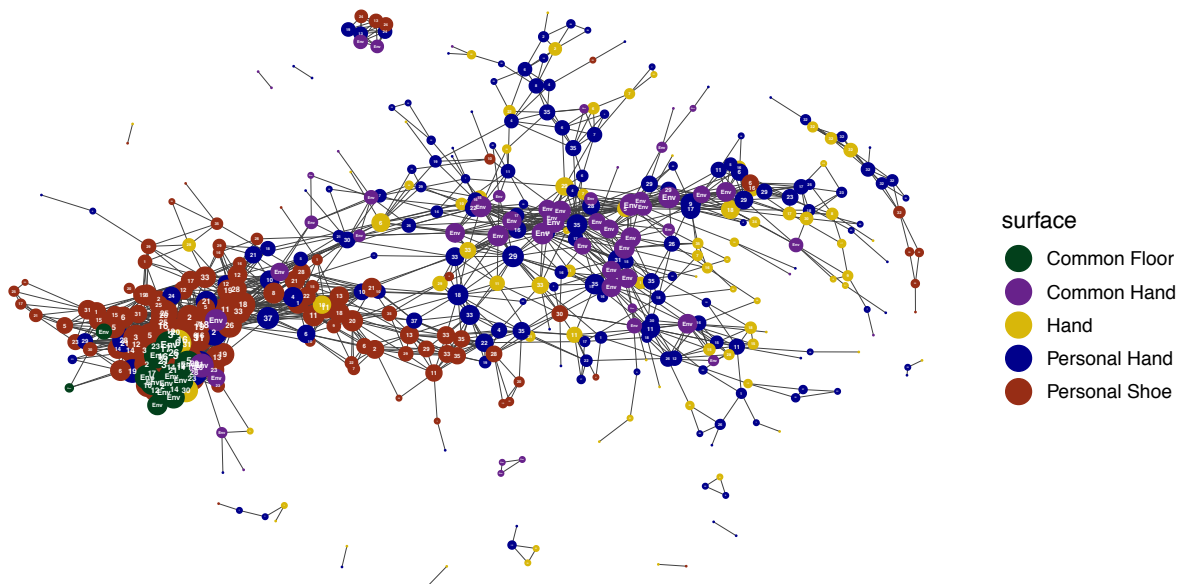
## Figures & Tables



249 **Figure 1:** (a) The distribution of phylogenetic distance between all taxa in each sequence processing  
250 method. MED recovers more highly related taxa than DADA2 or UPARSE. (b) The distribution of  
251 importance scores over all taxa, grouped by sequence processing method. The y-axis is log-  
252 transformed to aid visualization.



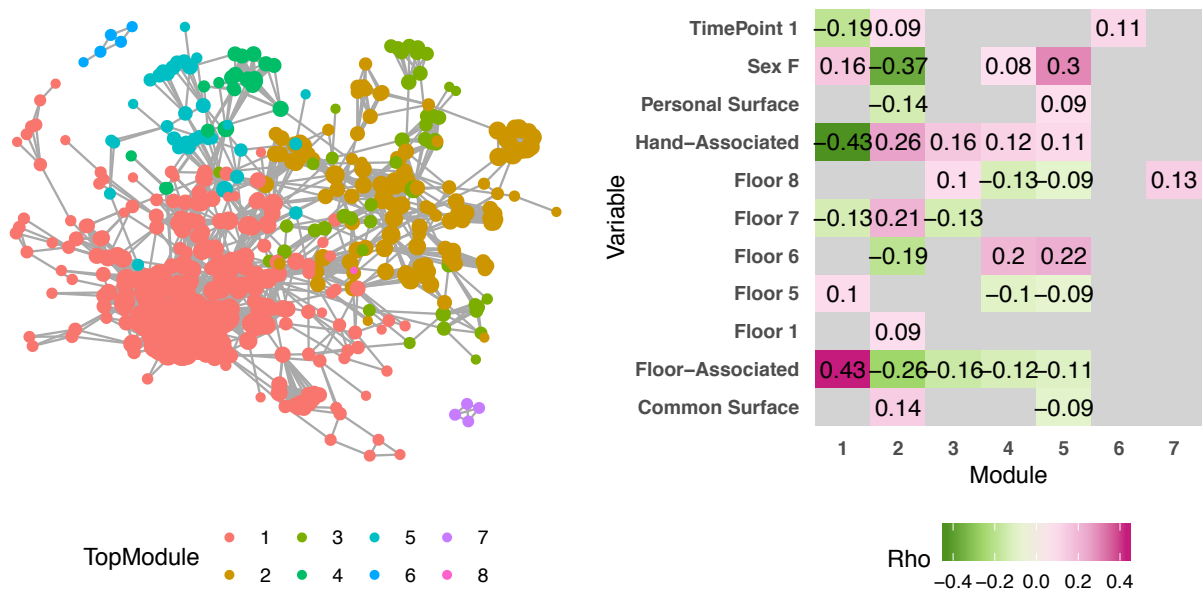
253 **Figure 2:** A principal components (PCoA) plot based upon the Bray-Curtis distance. Significant  
254 environmental vectors are plotted over the data.



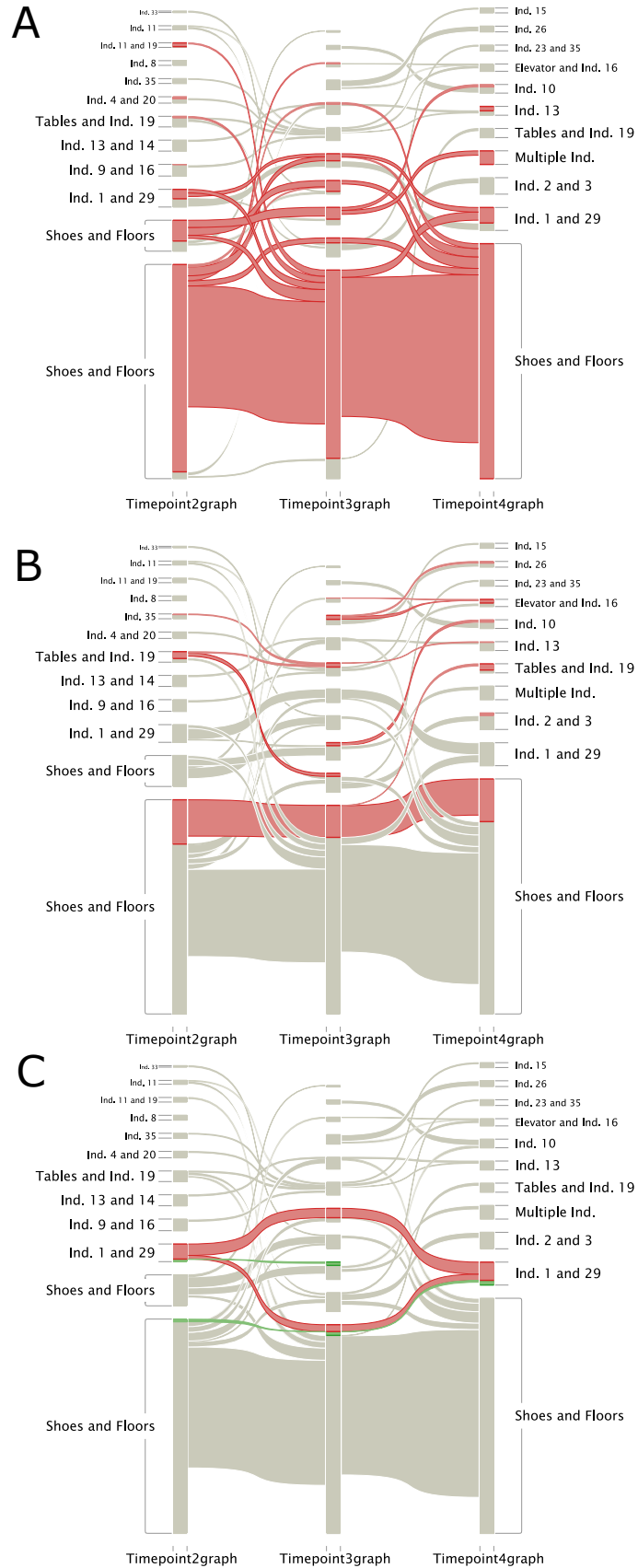
255 **Figure 3:** A Weighted-Unifrac graph of all samples, thresholded to be below .35 distance between  
256 individuals. Samples are colored by individual or environmental ID, while they are shaped by one of



257 the four large sample types, personal vs common and if they are hand or shoe associated. Common  
 258 hand-associated surfaces act as a scaffold, connecting between themselves, along with connecting  
 259 many distinct individuals.



260 **Figure 4: (a)** A graph generated using random forest proximity scores, trained to distinguish  
 261 individuals. It is thresholded by proximity greater than 0.076. It is colored by clique. Module 1 is  
 262 mostly composed of shoe and floor samples, similarly to **Figure 3**. **(b)** Significant Spearman  
 263 correlations ( $p < 0.05$ ) between each module and various metadata categories.



264 **Figure 5:** Alluvial Diagrams depicting the clustering of samples over time. (a) All samples that were  
 265 floor-associated (hallway floors, bedroom floors, and shoes) were colored red. (b) By common surface.  
 266 (c) Individual 1 (red) Individual 29 (green).

Method	UPARSE	DADA2	MED
Accuracy (CV-5)	60.96%	71.06%	79.57%
Error-Ratio	2.49	3.36	4.76

**Table 1:** Random Forest Accuracy and Error-Ratios

Floor	PersonalvCommon	Sex	Surface	Timepoint	SubjectID
0.39408	0.14456	0.39323	0.1985	0.05867	0.3245

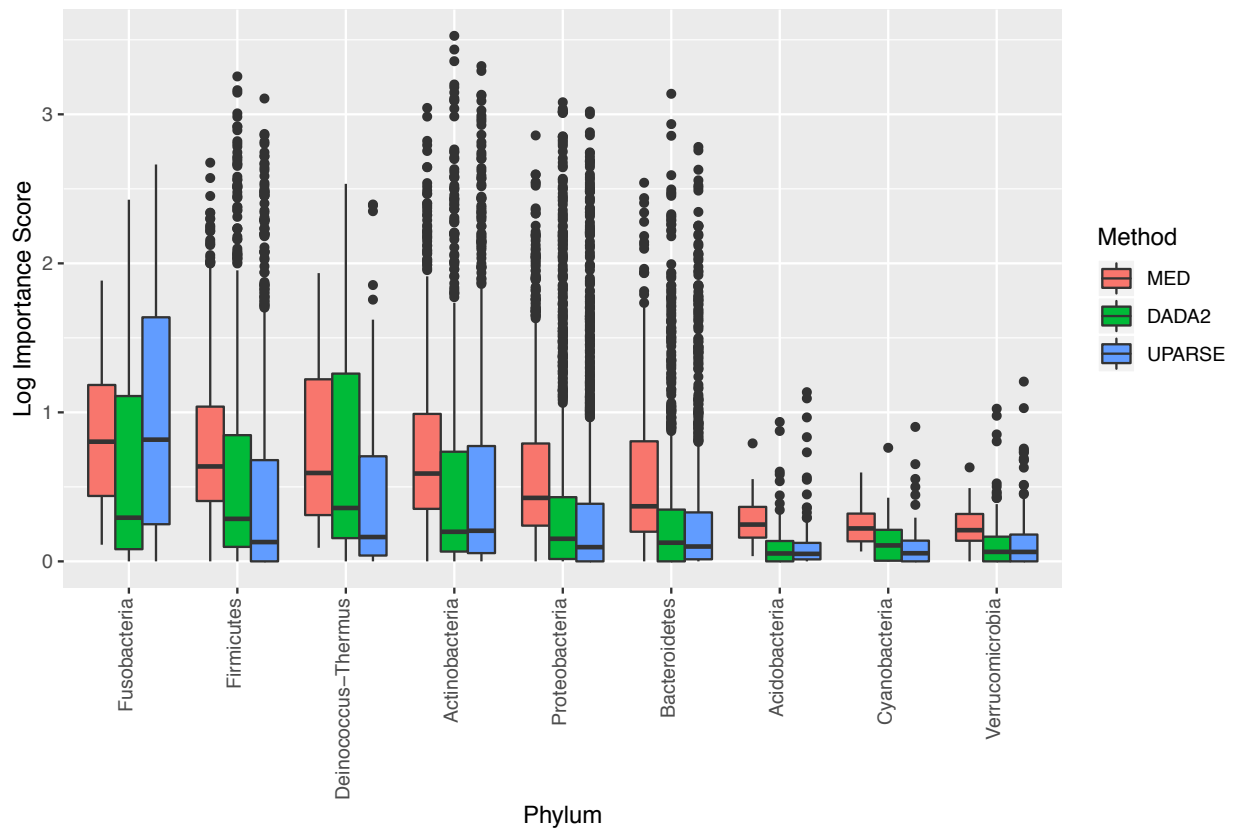
267 **Table 2:** Assortativity of Metadata Factors

Method	UPARSE	DADA2	MED
OTUs/Sequences	6011	4307	3352
Phyla	25	23	9
Average Phylogenetic Distance	2.62	2.27	1.36

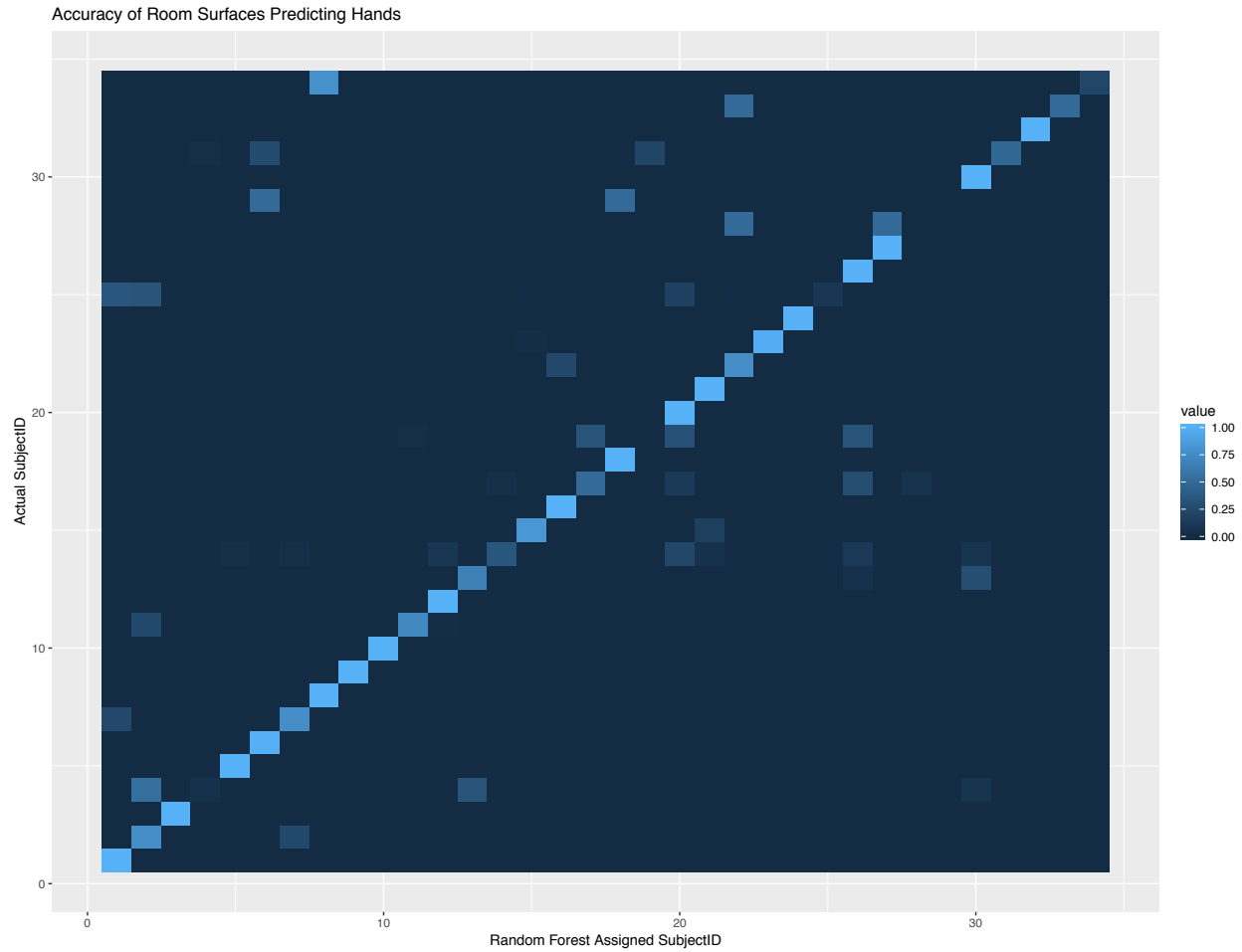
268 **Supplementary Table 1:** The OTU abundance and phylum-level diversity of each method. (fill this  
 269 out with every taxonomic level.)

Floor	Number of Participants	Common Surfaces
5	11	Male Bathroom Door Handle, Hallway Floor, Entry Door, Common Table, Elevator Buttons
6	5	Female Bathroom Door Handle, Hallway Floor, Elevator Button
7	9	Male Bathroom Door Handle, Hallway Floor, Elevator Button
8	12	Female Bathroom Door Handle, Hallway Floor, Elevator Button

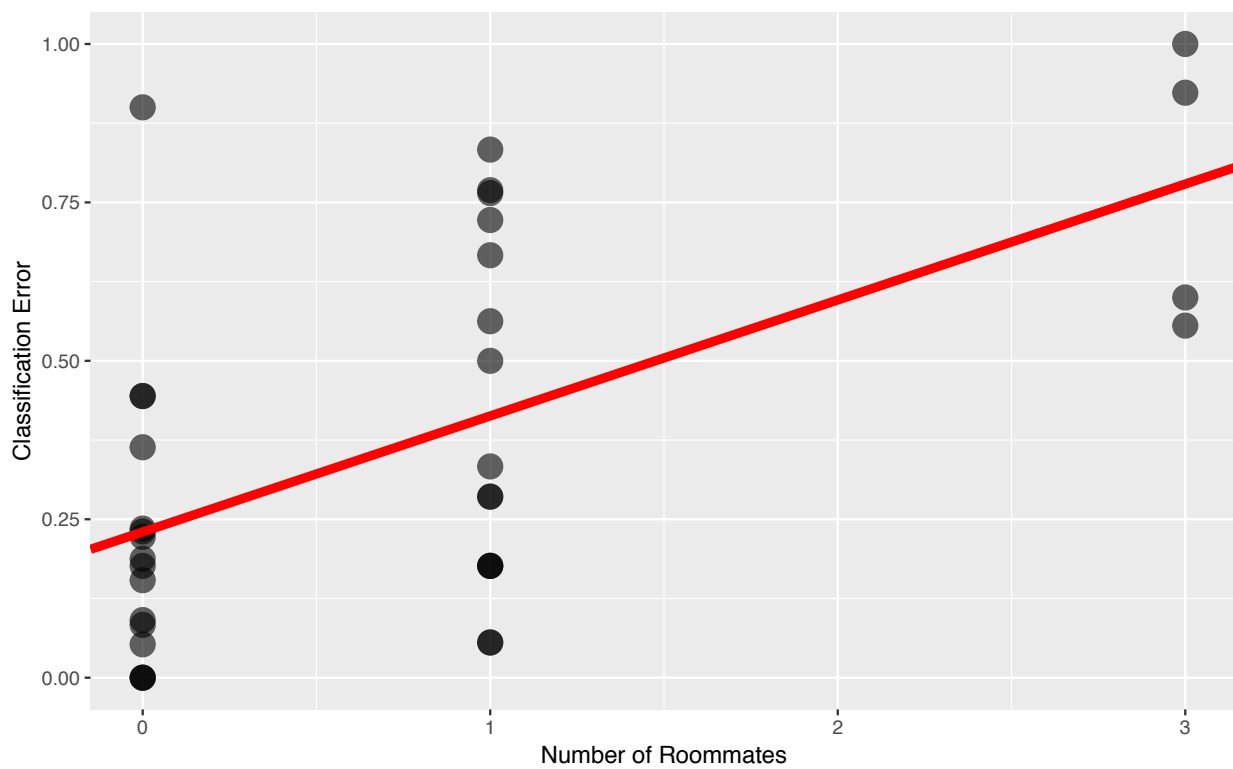
270 **Supplementary Table 2:** Summary of the number of participants on each floor of the dormitory  
 271 and which common surfaces were sampled on each floor.



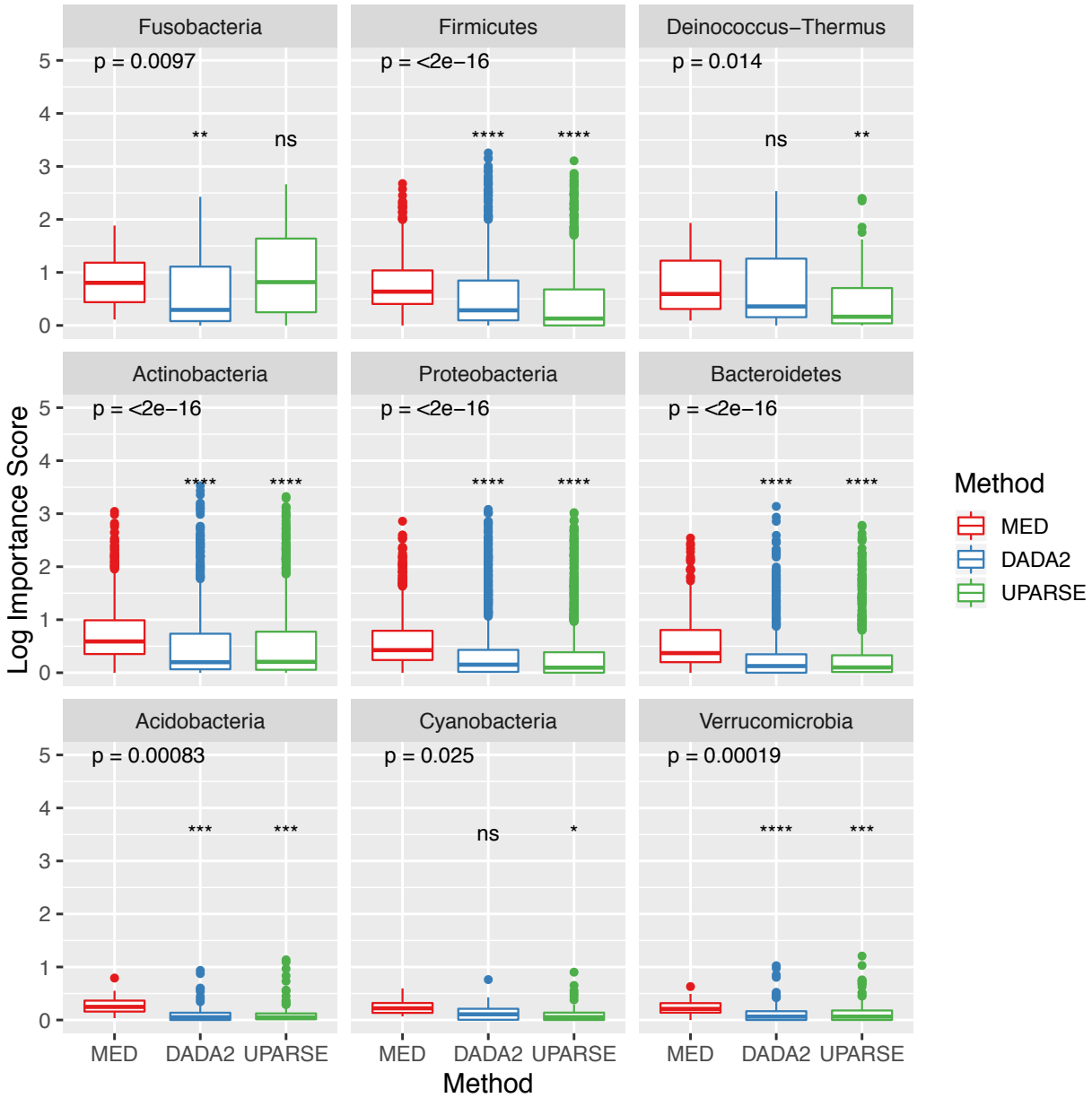
272 **Supplementary Figure 1:** Distribution of importance scores by phylum and sequence processing  
273 method. Importance score is log transformed to aid visualization. All phyla are significantly different  
274 by Kruskal-Wallis test, and MED has a higher average importance score for all phyla except for  
275 Fusobacteria. There are a large number of outliers in abundant taxa, due to the non-normality of  
276 importance scores.



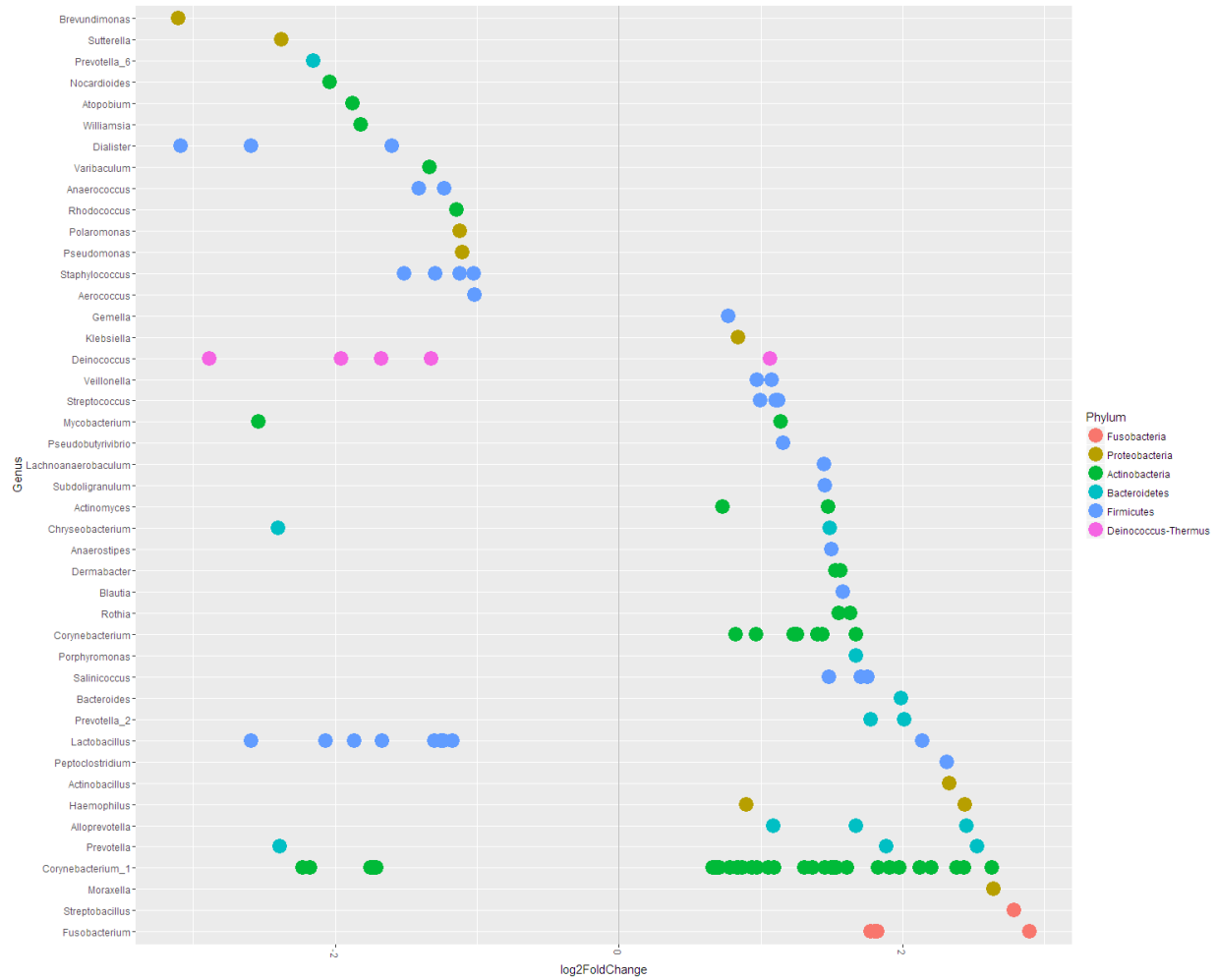
277 **Supplementary Figure 2:** A confusion matrix generated by the results of a random forest. It  
278 compares the actual identity of a sample with the one assigned by the random forest. Accurate  
279 classification appears on the diagonal, and any deviation is a mislabeled sample. Mostly samples fall  
280 on the diagonal, reflecting the 4.76 error ratio.



281 **Supplementary Figure 3:** Classification error plotted against the number of roommates.

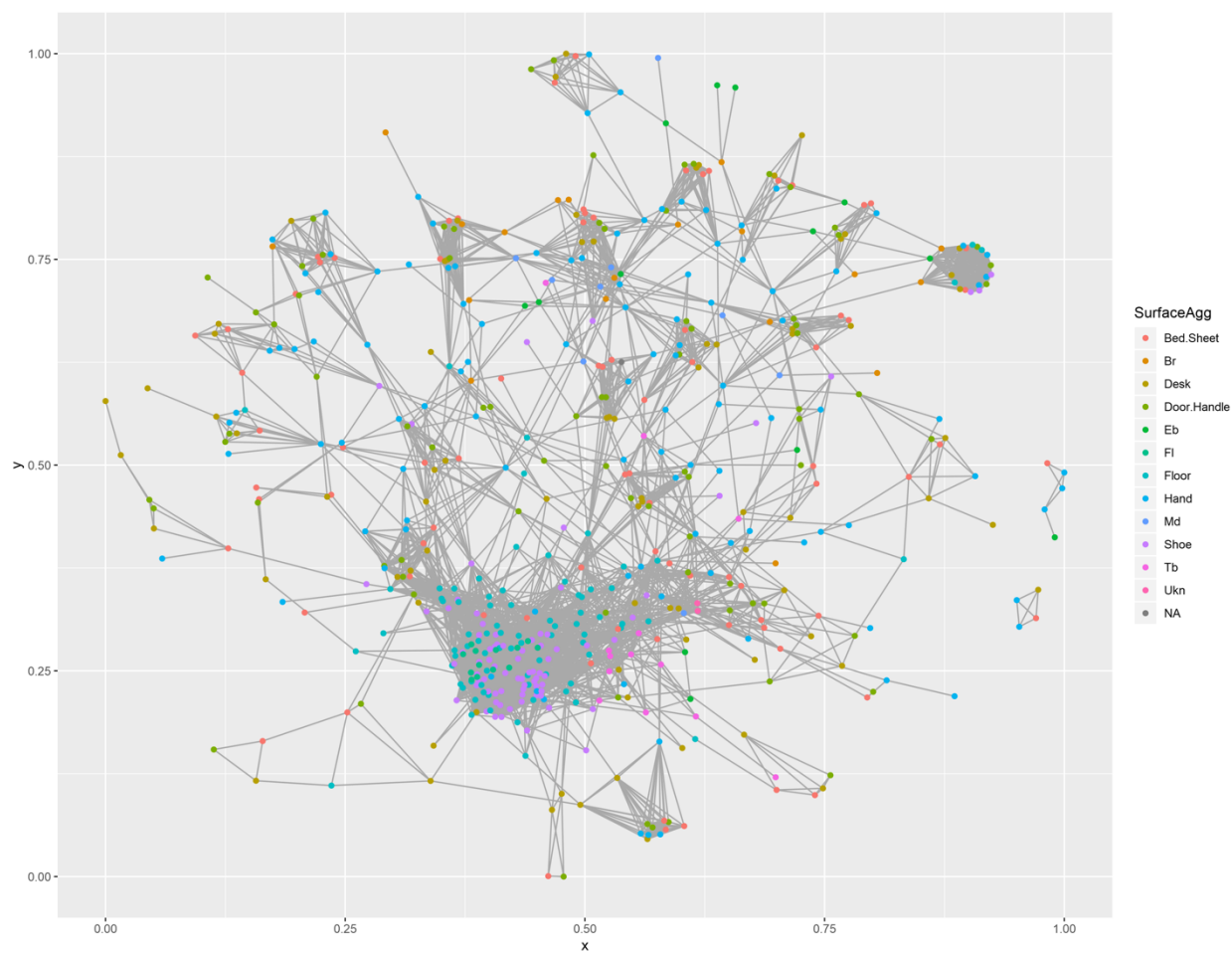


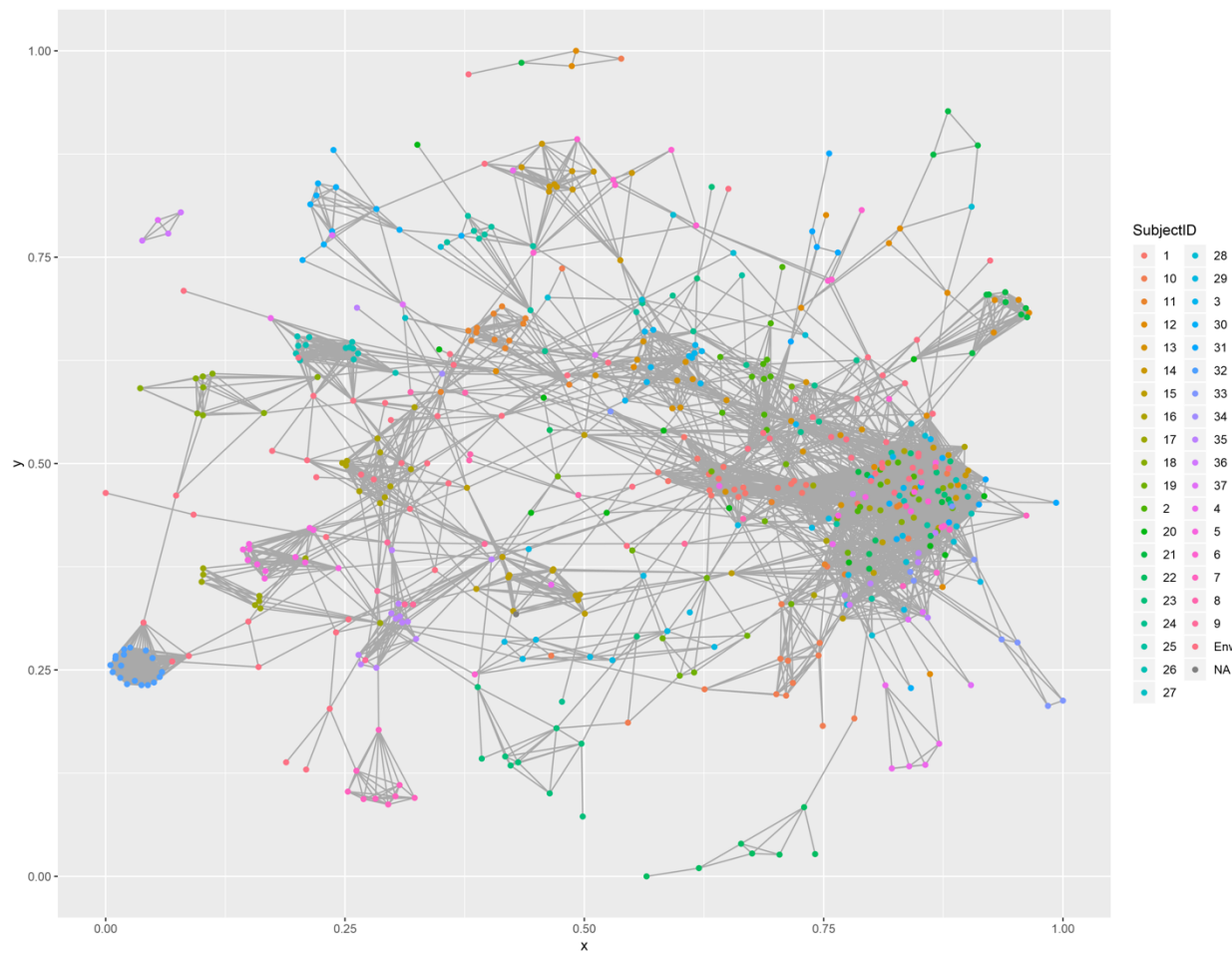
282 **Supplementary Figure 4:** The distribution of importance scores by phylum, with both Kruskal-  
283 Wallis significance between all pairs, and Wilcox-results of means compared to MED importance  
284 scores



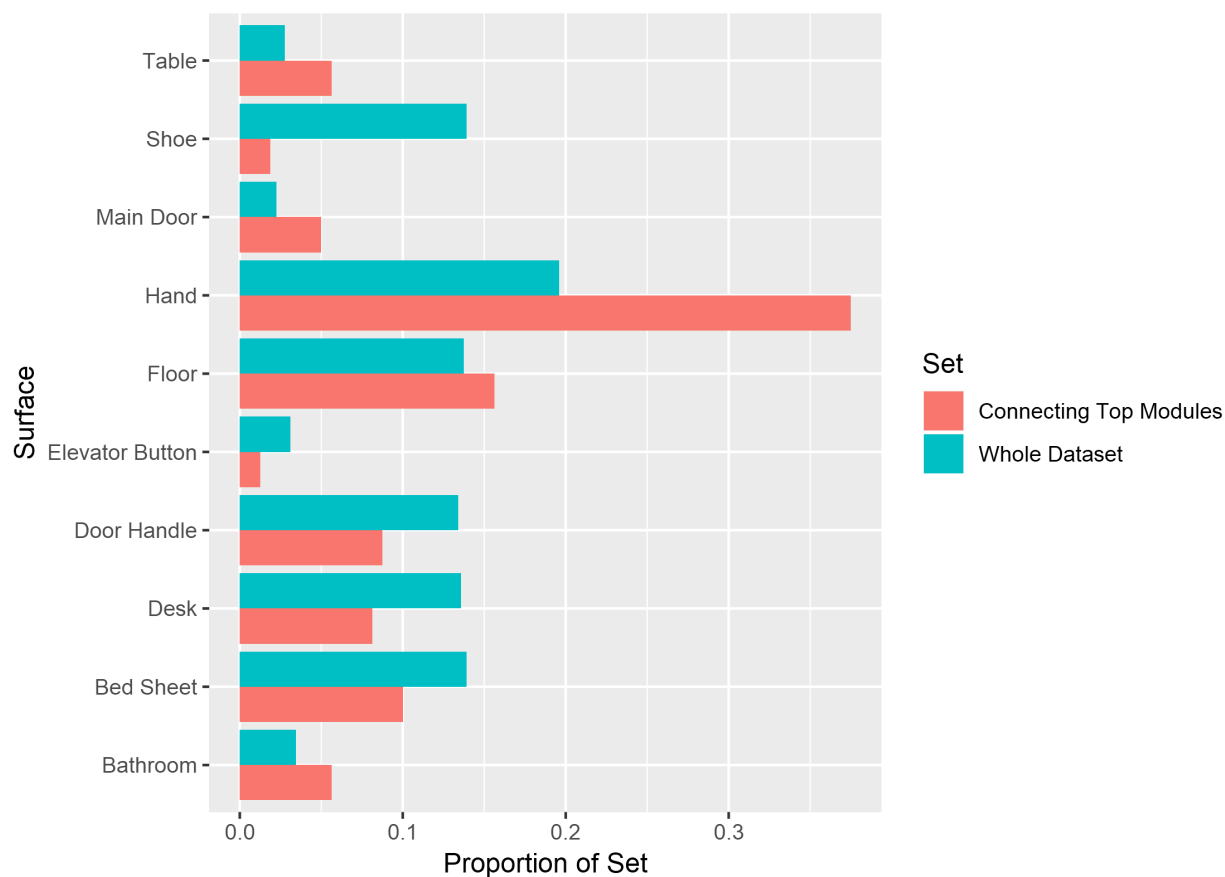
285 **Supplementary Figure 5:** Differential abundant sequences between male and female individuals,  
286 with women at left and men at right.







287 **Supplementary Figure 6:** (a) Random forest proximity graph colored by individual, (b) colored by  
288 surface type.



289 **Supplementary Figure 7:** Bar graph comparing the proportions of samples connecting top  
290 modules compared to those in the dataset at large. Hands show significant enrichment in  
291 connecting modules, indicating that they are the likely source of exchange between modules.