

1
2 **CRISPR arrays as high-resolution markers to track microbial transmission during**
3 **influenza infection**
4

5 Lingdi Zhang¹, Jahan Rahman¹, Lauren Lashua¹, Aubree Gordon², Angel Balmaseda^{3,4},
6 Guillermina Kuan^{3,5}, Richard Bonneau¹, Elodie Ghedin^{1,6*}

7
8 1. Center for Genomics and Systems Biology, Department of Biology, New York University, New
9 York, NY 10003

10 2. Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, MI
11 48109

12 3. Sustainable Sciences Institute, Managua, Nicaragua

13 4. Laboratorio Nacional de Virología, Centro Nacional de Diagnóstico y Referencia, Ministry of
14 Health, Managua, Nicaragua.

15 5. Centro de Salud Sócrates Flores Vivas, Ministry of Health, Managua, Nicaragua.

16 6. Systems Genomics Section, Laboratory of Parasitic Diseases, National Institutes of Health,
17 NIH, Bethesda, MD 20894

18
19
20

21 *Corresponding author: elodie.ghedin@nih.gov

22
23
24
25

26 Keywords: Metagenomics, metatranscriptomics, influenza virus, microbiome

27 **Abstract**

28 **Background**

29 The microbial community present in the respiratory tract can be disrupted during influenza virus
30 infection, leading to functional effects on the microbial ecology of the airways and potentially
31 impacting transmission of bacterial pathogens. Determining the transmission of airway
32 commensals, which can carry antibiotic resistance genes that could in turn be transferred to
33 bacterial pathogens, is of public health interest. Metagenomic-type analyses of the microbiome
34 provide the resolution necessary for microbial tracking and functional assessments in the airways.

35 **Results**

36 We obtained 221 respiratory samples that were collected from 54 individuals at 4 to 5 time points
37 across 10 households, with and without influenza infection, in Managua, Nicaragua. From these
38 samples we generated metagenomic (whole genome shotgun sequencing) and
39 metatranscriptomic (RNA sequencing) datasets to profile microbial taxonomy and gene
40 orthologous groups. Overall, specific bacteria and phages were differentially abundant between
41 influenza positive households and control (no influenza infection) households, with bacterial
42 species like *Moraxella catarrhalis* and bacteriophages like *Chivirus* significantly enriched in the
43 influenza positive households. Some of the bacterial taxa found to be differentially abundant were
44 active at the RNA level with genes involved in bacterial physiology differentially enriched between
45 influenza positive and influenza negative samples, primarily for *Moraxella*. We identified and
46 quantified CRISPR arrays detected in the metagenomic sequence reads and used these as
47 barcodes to track bacteria transmission within and across households. We detected a clear
48 sharing of bacteria commensals and pathobionts, such as *Rothia mulcilaginos*a and *Prevotella*
49 bacteria, within and between households, indicating community transmission of these microbes.
50 Antibiotic resistance genes that mapped to *Rothia* and *Prevotella* were prevalent across our
51 samples. Due to the relatively small number of households in our study we could not determine if
52 there was a correlation between increased bacteria transmission and influenza infection.

53

54 **Conclusion**

55 This study shows that microbial composition and ecological disruption during influenza infection
56 were primarily associated with *Moraxella* in the households sampled. We demonstrated that
57 CRISPR arrays can be used as high-resolution markers to study bacteria transmission between
58 individuals. Although tracking of antibiotic resistance transmission would require higher resolution
59 mapping of antibiotic resistance genes to specific bacterial genomes, we observed that individuals
60 connected by shared bacteria had more similar antibiotic resistance gene profiles than non-
61 connected individuals from the same households.

62

63

64

65

66

67

68

69

70 **Introduction**

71 Influenza infection as a contagious respiratory illness causes significant morbidity and mortality
72 worldwide. Bacterial co-infection during influenza infection, particularly in the elderly and
73 immunocompromised populations, can play an important role in disease progression leading to
74 complications and severe disease outcomes [1]. Infections with respiratory viruses can also
75 disrupt the microbiome of the airways and potentially contribute to disease severity [2]. A number
76 of studies have demonstrated viral disruption of the microbiota in the respiratory tract with
77 changes in relative abundance of bacterial taxa such as *Pseudomonas*, *Corynebacterium*, and
78 *Streptococcus* [3, 4]. The use of antibiotics, often prescribed for influenza patients because of
79 secondary bacterial infections, can disrupt the microbiota and diminish the protective function of
80 the microbiome [5], as well as contribute to the emergence of antibiotic resistant bacterial strains.
81 We and others have shown that the respiratory tract can be a potential reservoir of antibiotic
82 resistance genes in humans. Interestingly, the presence and expression of antibiotic resistance
83 genes detected were often to drugs the subjects had not been taking [6, 7], indicating potential
84 transmission between individuals. The transmission of opportunistic pathogens in the respiratory
85 tract, such as *Streptococcus pneumoniae*, is known to be associated with respiratory tract viral
86 infection and younger age of the infected subject [8, 9]. For bacteria to transmit to a new host, the
87 invading bacteria need to interact with the residing microbes and establish colonization [8, 10],
88 more likely to occur when the microbiome in the new host is disrupted.

89
90 Despite these observations for influenza and other viral pathogens [2-4], the dynamics of the
91 respiratory tract microbiome and the disruption of its overall microbial ecology in human seasonal
92 influenza infection remains to be characterized. For example, very little has been reported on the
93 modulation of phages in the respiratory tract during acute respiratory virus infection [11]. Phages,
94 an important entity in the human microbiome, were found to shape and co-evolve with their
95 bacterial hosts impacting bacterial growth, metabolic activities, virulence, and antibiotic resistance

96 [12]. The compositional shift of phages and the microbiome were found to be associated with
97 human disease, such as inflammatory bowel disease [13]. Interacting with both their bacterial
98 hosts and eukaryotic cells, phages can potentially provide another layer of information to the
99 characterization of the microbiome in response to influenza infection. In this study, we used
100 CRISPR spacers to study the interactions between bacteria and phages. CRISPR functions as
101 the bacterial immune system to defend against virus infection by integrating a 20-70 bp viral
102 spacer into the CRISPR locus when the bacteria are first exposed to the virus. Bacteria that have
103 the integrated sequences are then able to defend themselves against viruses that match those
104 spacer sequences [14]. We hypothesized that this unique history record of a bacteria's encounter
105 with a phage could be used to profile the dynamics of the microbial ecology within the respiratory
106 tract.

107

108

109 A goal of our study was to profile the disruption of the microbiome with influenza infection and to
110 determine, from metagenomic analyses of upper respiratory tract samples, whether we could
111 quantify transmission of commensal bacteria and antibiotic resistance genes. Currently, most
112 studies on bacteria transmission focus on one bacteria species and use single nucleotide
113 polymorphisms (SNPs) in marker genes [15] or whole bacterial genomes [16, 17]. If using
114 metagenomics data, this would require very deep sequencing depth and could only sufficiently
115 profile SNPs from the most abundant bacterial genomes. Instead, we leveraged the unique nature
116 of bacterial CRISPR arrays as markers to track bacteria transmission of different bacteria species.
117 Viral spacers are constantly acquired by the bacteria and integrated at the end of CRISPR arrays,
118 proximal to the leader sequence [14]. Although the spacer sequences that the bacteria acquire
119 from a specific virus are not entirely random, as bias in spacer sequence distribution has been
120 observed [18, 19], the possible number of unique spacer sequences bacteria can acquire from a
121 virus infection is large. Given the dynamics of the CRISPR arrays, the probability that individuals

122 share the exact same sequences due to independent spacer integration events is negligible. We
123 demonstrate that CRISPR arrays can indeed be used to study bacteria transmission with better
124 resolution than SNP-type analyses, especially when the CRISPR array is large.

125

126 **Results**

127 **Study cohort and sample collection**

128 We obtained 221 respiratory samples (pooled nasal and throat swabs) that were collected from
129 54 individuals participating in the Household Influenza Transmission Study (HITS) in Managua,
130 Nicaragua. In total, 10 households with 4-8 members in each household participated in the study,
131 and samples were collected at 4 to 5 time points for each individual, at 2-4 day intervals. Sample
132 collection was independent of influenza infection, thus some of the samples were collected at time
133 points when the individual was not infected or had recovered (**Table S1**). The households were
134 assigned to high, low, or no influenza virus (control) infection groups based on the number of
135 individuals per household who tested positive for influenza. High infection households had all or
136 2/3 of the household members testing positive at some point over the serial sampling (5-8
137 household members), while the low infection households had less than a third of household
138 members testing positive for influenza at any time point (2-3 members). The 'no flu' households
139 represent uninfected controls (**Table S1**). We did not sample all the household members from the
140 low influenza and control households. Influenza infection was diagnosed by rtPCR and the
141 infections were all due to influenza A virus subtype H3N2. Total RNA and DNA were extracted
142 from each sample and processed for RNAseq (metatranscriptomics) and whole genome shotgun
143 (metagenomics), respectively, for an in-depth microbiome analysis of the upper respiratory tract
144 across household members. Of the 221 samples, we obtained 167 metagenomics and 178
145 metatranscriptomics libraries; 135 samples had both metagenomics and metatranscriptomics
146 data and we focused on these samples for a subset of the analyses.

147

148 **Impact of influenza infection on microbial composition in the upper airways**

149 To profile the microbial composition in subjects with and without influenza across the households
150 studied, we did a taxonomic classification on the metagenomics and metatranscriptomics reads
151 post filtering (human reads were removed from both datasets and rRNA reads were removed
152 from the metatranscriptomics dataset). We detected bacteria, phages, and eukaryotic viruses
153 across the samples (Human viruses are shown in **Tables S2 and S3**) by using kraken2, which is
154 based on exact k-mer matches to reference genomes [20]. Of the human viruses detected, beta
155 herpesvirus was the most prevalent across samples. There was no correlation between viruses
156 detected and the influenza infection status of the individuals, except for influenza A virus
157 sequence reads, which were, as expected, enriched in the flu positive samples, validating the
158 quality of the metatranscriptomics datasets (Fisher's Exact test and $FDR \leq 0.05$).

159

160 To analyze the microbial community across the household groups and influenza infection status,
161 we compared the relative abundance of bacteria and viruses for different comparisons; the
162 household subjects and samples are summarized in **Table 1** and **Table 2**. Children were enriched
163 in the flu infection group as they were the index cases and the first household members tested to
164 have influenza infection. Pooling samples from different timepoints for each individual, we
165 identified significant differences in bacteria diversity between household groups at both DNA and
166 RNA levels (PERMANOVA [21] p value=0.0005 and 0.043 for metagenomes and
167 metatranscriptomes, respectively; **Fig. 1a** and **1b**). We applied differential expression analysis
168 (DESeq2 [22]) to identify specific bacterial taxa that drove the differences in the beta diversity.
169 *Dolosigranulum* was significantly enriched in the metagenomes of the control households
170 compared to the low infection households but enriched in the metatranscriptomes of the high flu
171 infection households compared to the low infection households (**Fig. 1c** and **1d**). *Moraxella* was

172 enriched in the high flu infection households for both metagenomes and metatranscriptomes (**Fig.**
 173 **1c** and **1d**).

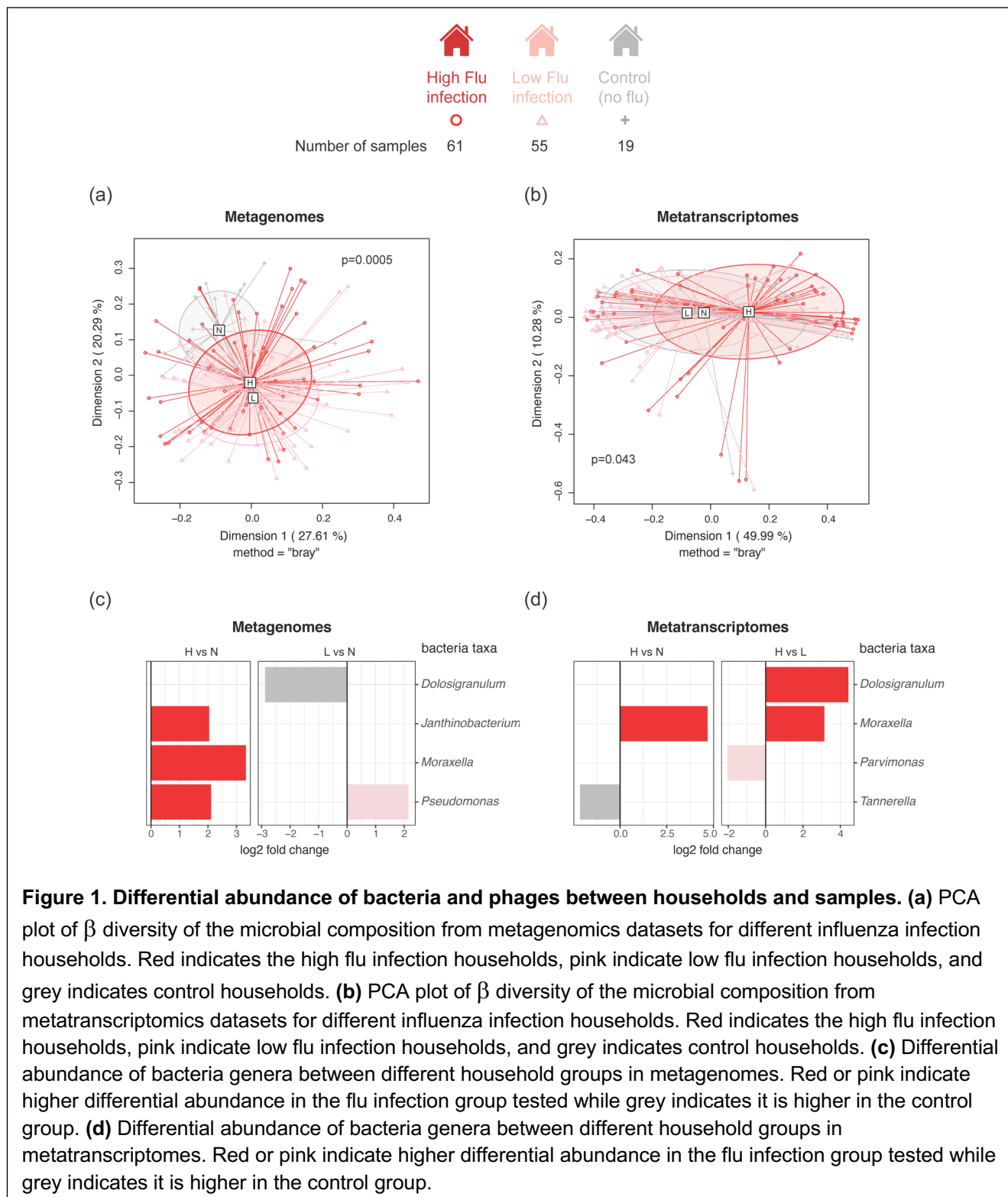
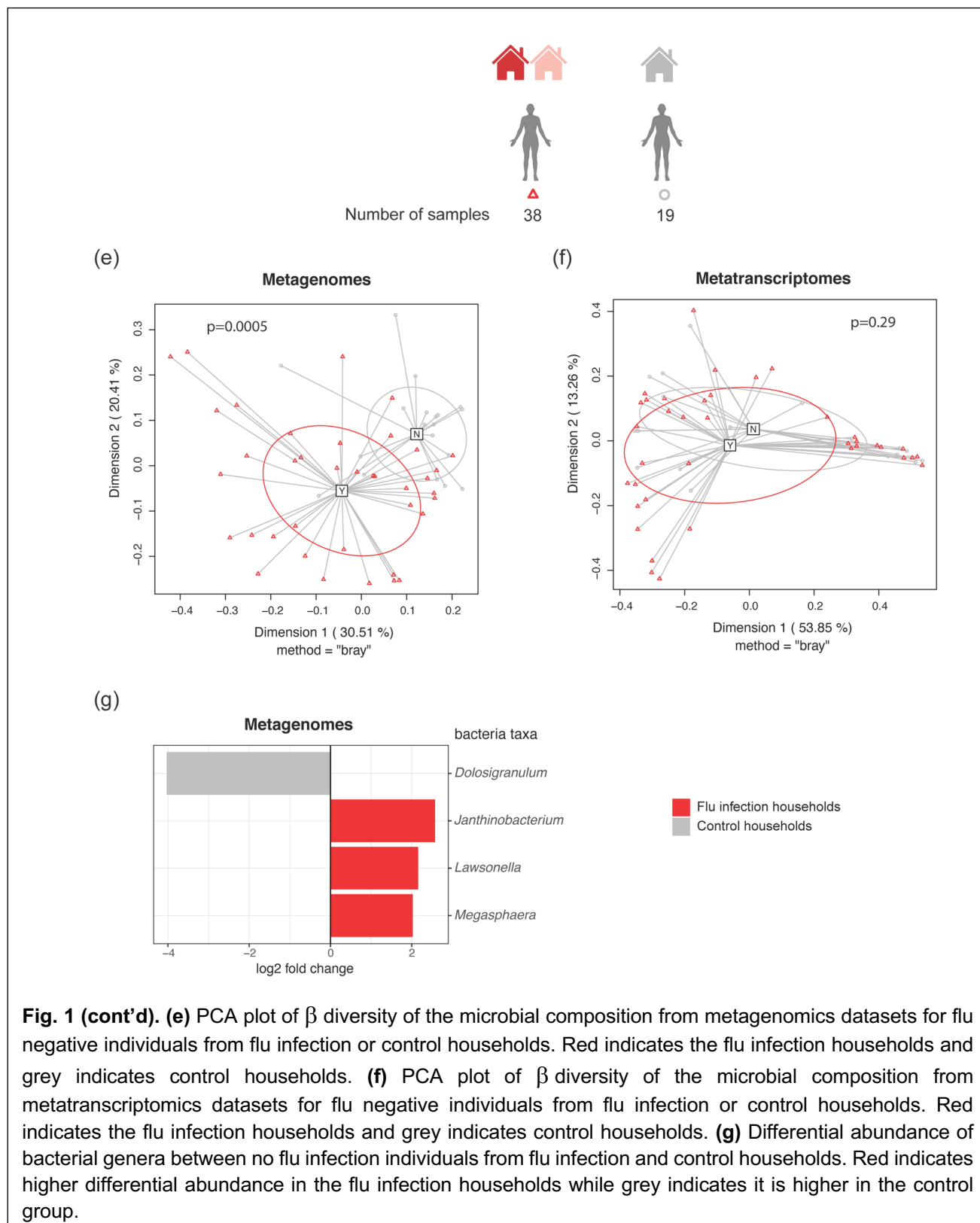


Figure 1. Differential abundance of bacteria and phages between households and samples. (a) PCA plot of β diversity of the microbial composition from metagenomics datasets for different influenza infection households. Red indicates the high flu infection households, pink indicate low flu infection households, and grey indicates control households. **(b)** PCA plot of β diversity of the microbial composition from metatranscriptomics datasets for different influenza infection households. Red indicates the high flu infection households, pink indicate low flu infection households, and grey indicates control households. **(c)** Differential abundance of bacteria genera between different household groups in metagenomes. Red or pink indicate higher differential abundance in the flu infection group tested while grey indicates it is higher in the control group. **(d)** Differential abundance of bacteria genera between different household groups in metatranscriptomes. Red or pink indicate higher differential abundance in the flu infection group tested while grey indicates it is higher in the control group.



175 By comparing flu negative samples from individuals in the influenza infection households with the
176 (flu negative) samples from the control households, we identified significant differences in
177 bacterial composition of the metagenomes (p-value=0.0005) but not of the metatranscriptomes
178 (p-value=0.29) (**Fig. 1e** and **1f**). A few bacteria taxa were differentially enriched between the two
179 groups, including *Dolosigranulum*, which was enriched in the control households (**Fig. 1g**).

180
181 We compared samples from individuals who did not test positive for influenza at any time point
182 from any of the households (including the controls) and compared them to influenza positive
183 samples. The microbial composition was significantly different (metagenomes p-value=0.015) but
184 not the microbial expression profile (metatranscriptomes p-value=0.099, **Fig. 2a**). *Moraxella* was
185 enriched in the flu positive samples (**Fig 2b**). Since *Moraxella* was enriched in the high infection
186 households and the flu infection samples, indicating it correlated with influenza infection, we
187 compared the relative abundance of the *Moraxella* species between flu negative time points
188 (baseline) and flu positive time points from the same individuals. A few of the *Moraxella* spp,
189 especially *Moraxella catarrhalis*, were moderately enriched in the flu positive time points (log2 fold
190 change >1) in both metagenome and metatranscriptome datasets (**Fig. 2c**). We also identified
191 phages, such as Salmonella phage Chi (or *Chivirus*) which was enriched in the high flu infection
192 households compared to the 'no flu' infection households (FDR<=0.05, log2 fold change =2.17)
193 while Pahexavirus was enriched in the low flu infection households compared to the 'no flu'
194 infection households (FDR<=0.05, log2 fold change =2.69)

195
196 To characterize microbial functional changes during influenza infection, we profiled the relative
197 abundance of microbial genes with FMAP, a tool that provides a functional analysis from
198 metagenomic and metatranscriptomic sequencing data [23]. We also identified gene orthology
199 groups (KEGG) and their relative abundance. As previously done, we partitioned the samples into
200 influenza infection timepoints and no infection groups (only samples from individuals that did not

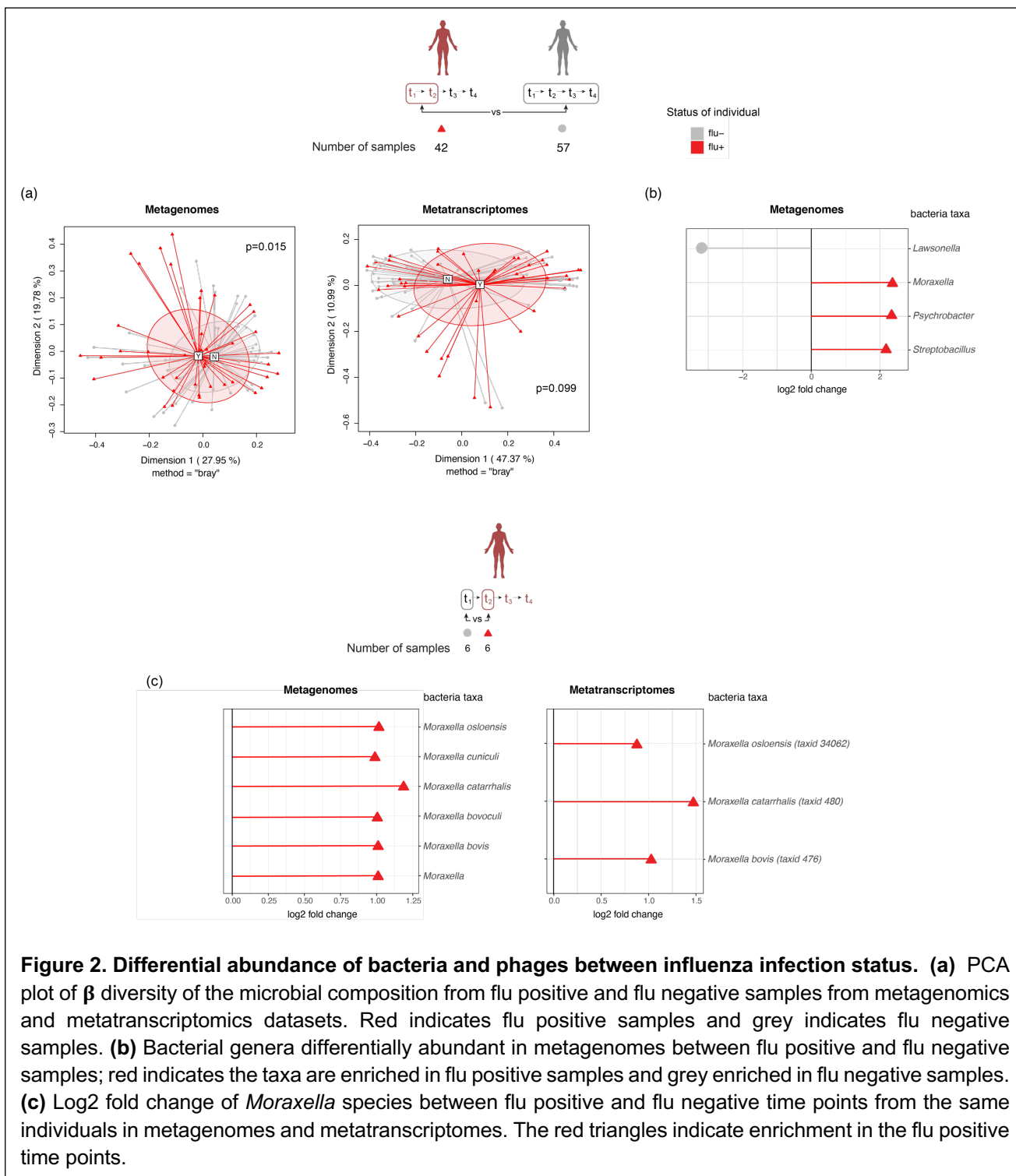


Figure 2. Differential abundance of bacteria and phages between influenza infection status. (a) PCA plot of β diversity of the microbial composition from flu positive and flu negative samples from metagenomics and metatranscriptomics datasets. Red indicates flu positive samples and grey indicates flu negative samples. **(b)** Bacterial genera differentially abundant in metagenomes between flu positive and flu negative samples; red indicates the taxa are enriched in flu positive samples and grey enriched in flu negative samples. **(c)** Log₂ fold change of *Moraxella* species between flu positive and flu negative time points from the same individuals in metagenomes and metatranscriptomes. The red triangles indicate enrichment in the flu positive time points.

201 test as flu positive were included in the 'no infection' group). In the metagenomics data we
 202 identified *i* number of microbial genes differentially abundant between flu positive and flu negative
 203 samples, including genes that allow bacteria to adapt to the environment, such as the two-

204 component system where histidine kinase detect stimuli from the environment and trigger
 205 downstream regulatory responses [24] (**Fig. 3a**). We did not identify genes differentially
 206 expressed in the metatranscriptomics data. To identify which specific bacterial taxa contributed
 207 to the differentially abundant genes between flu positive and flu negative samples, we extracted

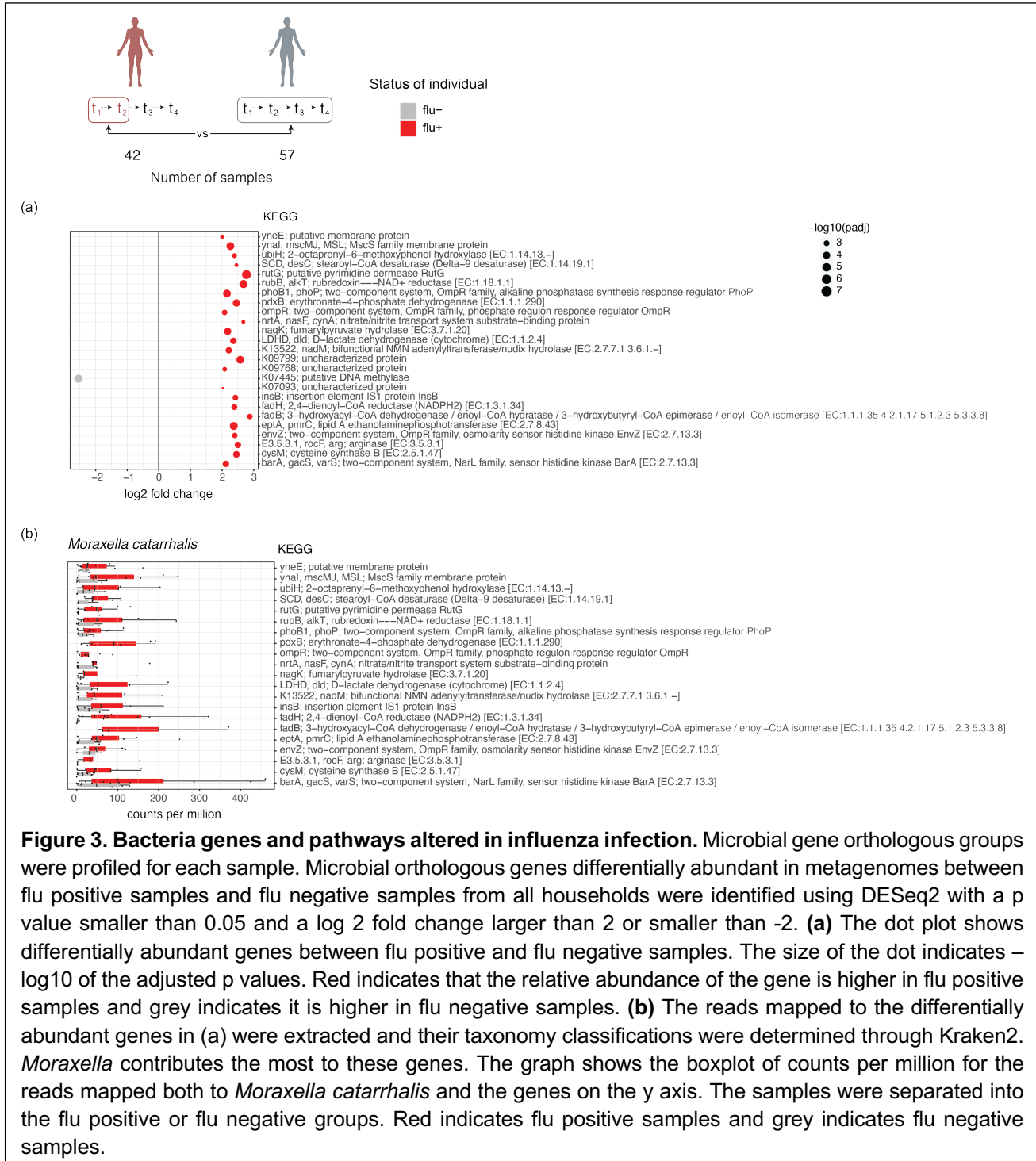


Figure 3. Bacteria genes and pathways altered in influenza infection. Microbial gene orthologous groups were profiled for each sample. Microbial orthologous genes differentially abundant in metagenomes between flu positive samples and flu negative samples from all households were identified using DESeq2 with a p value smaller than 0.05 and a log₂ fold change larger than 2 or smaller than -2. **(a)** The dot plot shows differentially abundant genes between flu positive and flu negative samples. The size of the dot indicates -log₁₀ of the adjusted p values. Red indicates that the relative abundance of the gene is higher in flu positive samples and grey indicates it is higher in flu negative samples. **(b)** The reads mapped to the differentially abundant genes in (a) were extracted and their taxonomy classifications were determined through Kraken2. *Moraxella* contributes the most to these genes. The graph shows the boxplot of counts per million for the reads mapped both to *Moraxella catarrhalis* and the genes on the y axis. The samples were separated into the flu positive or flu negative groups. Red indicates flu positive samples and grey indicates flu negative samples.

208 the reads mapped to these genes and assigned taxonomic classification using Kraken2.
209 *Moraxella* was the dominant genus contributing to the differentially abundant genes (**Fig. S1**),
210 with *Moraxella catarrhalis* largely contributing to these genes in the flu positive samples (**Fig. 3b**).

211

212 **Shared CRISPR spacers to identify transmission events**

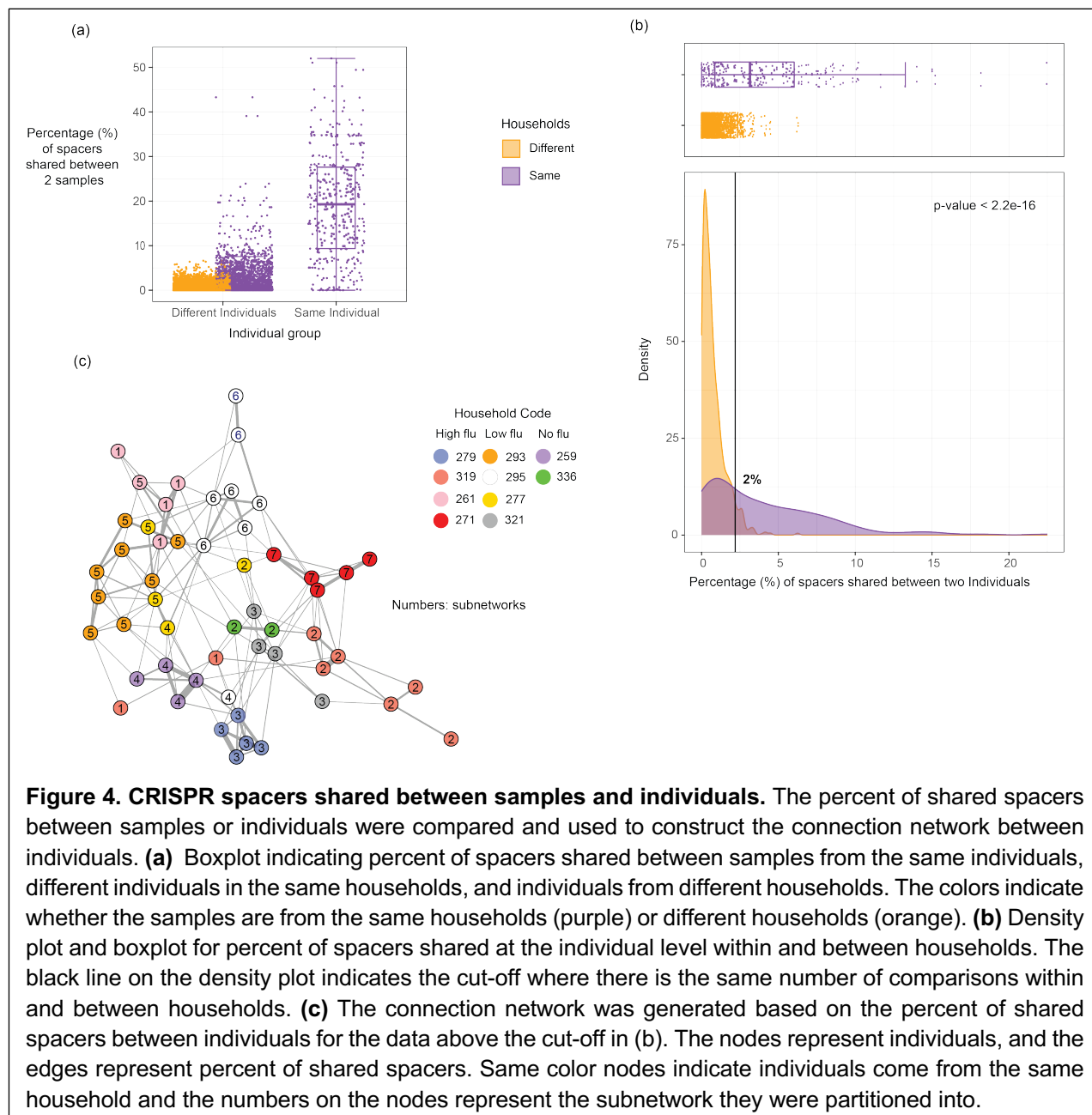
213 We used the metagenomics datasets to study phage-bacteria interactions by focusing on the
214 CRISPR array-integrated spacers. As phages play an important role in shaping the bacterial
215 population and could affect host immunity, we identified phages differentially abundant between
216 influenza positive households and control households. We profiled phage-bacteria interactions to
217 further investigate the microbial ecological changes associated with influenza infection and to
218 identify bacterial host species for the enriched phages found. The spacers in the CRISPR arrays,
219 originating from viruses and integrated into the bacterial genomes, were used to link the bacteria
220 and virus. The spacer sequences were mapped back to the viral and bacterial contigs, leading to
221 the identification of several bacteria and phages that were linked by the shared spacers. Phages
222 were connected to bacterial hosts both in and outside their designated host range (**Fig. S2**). Some
223 interactions between bacteria and phages were shared between individuals and others were
224 unique to an individual (**Fig. S2**). A few bacterial commensals and pathobionts in the respiratory
225 tract, such as *Prevotella* and *Rothia mucilliginosa*, were found to be infected by many phages. No
226 specific phage-bacteria interactions appeared to be associated with influenza infection
227 (determined by Fisher's Exact test). This indicates either that influenza infection does not disrupt
228 the phage-bacteria interaction dynamics in the respiratory tract or the interactions we profiled
229 happened before influenza infection.

230

231 One important question when considering the disruption of the respiratory microbiome in a
232 respiratory viral infection is whether certain bacteria with pathogenic potential are more likely to

233 be transmitted. Since many commensals and pathobionts are natural members of the respiratory
234 community [25, 26], determining the dynamics of respiratory commensal bacteria transmission
235 within households is challenging. We used the CRISPR spacers and arrays to track bacteria
236 transmission. Bacteria with the same set of spacers in the same order are likely to be related,
237 indicating a potential transmission event if the same CRISPR arrays are found in two individuals.
238 We first identified spacer sequences from the metagenomics reads using MetaCRIST. We
239 pooled the spacers across all the samples and identified spacers shared between samples based
240 on 90% sequence similarity. We then determined the proportion of spacers that were shared
241 between any two samples. Samples from the same individual collected at different time points
242 shared more spacers than samples from different individuals (**Fig. 4a**). Although spacer content
243 is dynamic over time within individuals because of continuous interactions between phages and
244 their bacterial hosts, the spacers were not significantly different over the sampling period. We also
245 found that the proportion of shared spacers was higher when comparing samples from individuals
246 living in the same household and individuals from different households, which is what we would
247 expect with transmission (**Fig. 4a**). To further compare spacers identified from individuals within
248 and across households, we pooled the serial samples for each subject and redid the analysis in
249 subject-to-subject comparisons. Individuals from the same households have a higher proportion
250 of shared spacers than individuals from different households (**Fig. 4b**), indicating more shared
251 bacteria. As the number of subject-to-subject comparisons is not balanced for within and between
252 households, we removed comparisons between individuals with lower than 2% shared spacers
253 (**Fig. 4b**), leading to an equal number of comparisons within and between households, helping us
254 to weigh comparisons between individuals from the same and different households equally and
255 removing noise. A connection network was then generated based on the percent of shared
256 spacers between individuals (**Fig. 4c**) where the nodes are the individuals and the edges are
257 weighted by the value of the percent of shared spacers. We also detected subnetworks within the
258 network using the shared spacer data between individuals (**Fig. 4c**). The correlation between the

259 partition of the nodes to the subnetworks and the household metadata was 0.79, indicating the
260 individuals within the same households are more tightly connected based on their percent of
261 shared spacers.



262

263 To determine which bacteria taxa were shared between individuals, we then investigated the
264 genomic sequence where the shared CRISPR arrays were located. To do so we assembled the

265 metagenomics reads into contigs from samples from the same individual and mapped the spacers
266 back to the contigs to identify CRISPR arrays and the order of the spacers. We used a dynamic
267 programming local alignment method to find the best alignments between any two CRISPR arrays
268 that come from different individuals. We only focused on the CRISPR array alignments with more
269 than 5 spacers. We analyzed the 2-dimensional density distribution of the CRISPR array
270 alignments for their alignment similarity and alignment length and compared the density
271 distributions for the alignments with CRISPR arrays from individuals from the same or different
272 households (**Fig. S4**). Using this approach, we found that density distributions are different
273 between the two groups we compared (**Fig. S4**) and statistically significant by using Kolmogorov-
274 Smirnov test (p -value < 0.05) (**Fig. 5a**). We identified a region on the density plot (the upright solid
275 purple region) that was enriched with alignments from the same households, and the alignments
276 in that region have higher alignment similarity and alignment length (**Fig. 5a**), giving us better

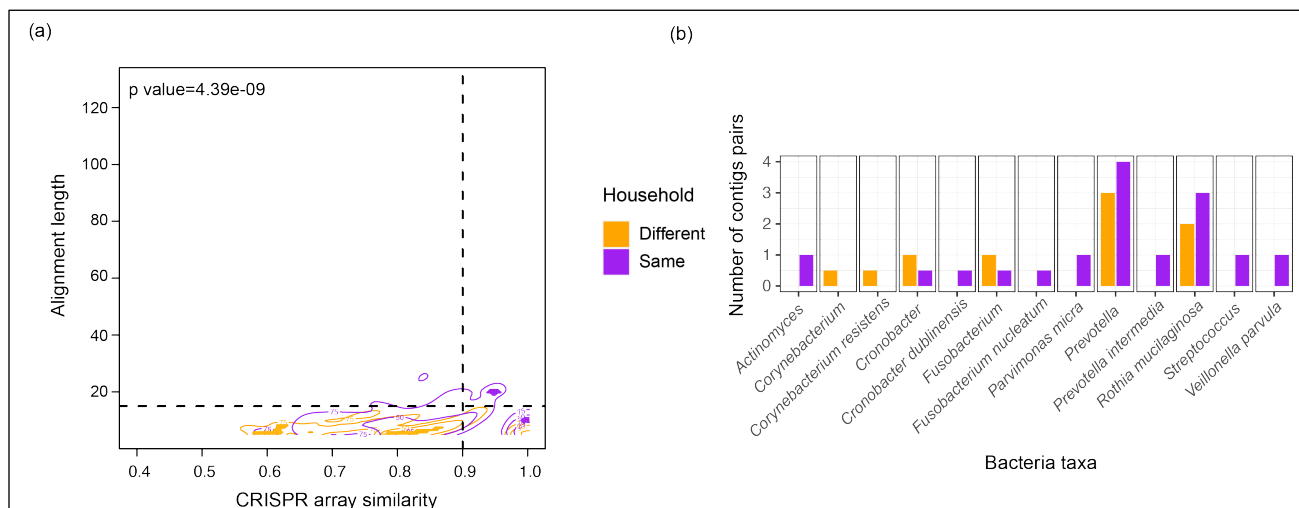


Figure 5. CRISPR array alignment distribution within and between households and bacteria taxa shared between individuals. (a) Contour plots for the 2D density distribution of CRISPR array alignments on alignment similarity and alignment length. The density distribution for the CRISPR array alignments with CRISPR arrays from individuals from the same households or different households were colored in purple and orange, respectively. The numbers on the contour plot indicate the regions have 25%, 50% and 75% of the data. The solid color regions on the plot indicate the density of the data in these regions were significantly different between the two groups. The purple region has higher data density in the same household group and the orange region has higher data density in the different household group. **(b)** The contigs contain the CRISPR arrays from the alignments with a similarity greater than 0.9 and more than 15 spacers mapped to the NCBI nt database. The graph shows the barplot of the number of individual pairs that share the bacteria with the color indicating they are from the same households or different households.

277 confidence that the CRISPR arrays in the alignments in this region are likely to be transmitted
278 between individuals. We filtered the alignments with a similarity greater than 0.9 and an alignment
279 length greater than 15. We extracted the contigs containing the CRISPR arrays from these
280 alignments and aligned the contigs to the NCBI nt database to get taxonomic assignments. We
281 identified contigs assigned to bacterial commensals and pathobionts that were found to be shared
282 within and across households, such as *Prevotella*, *Fusobacterium* and *Rothia* (**Fig. 5b**). Although
283 rare, some bacteria were shared only within households, such as *Streptococcus* and *Veillonella*
284 *parvula*.

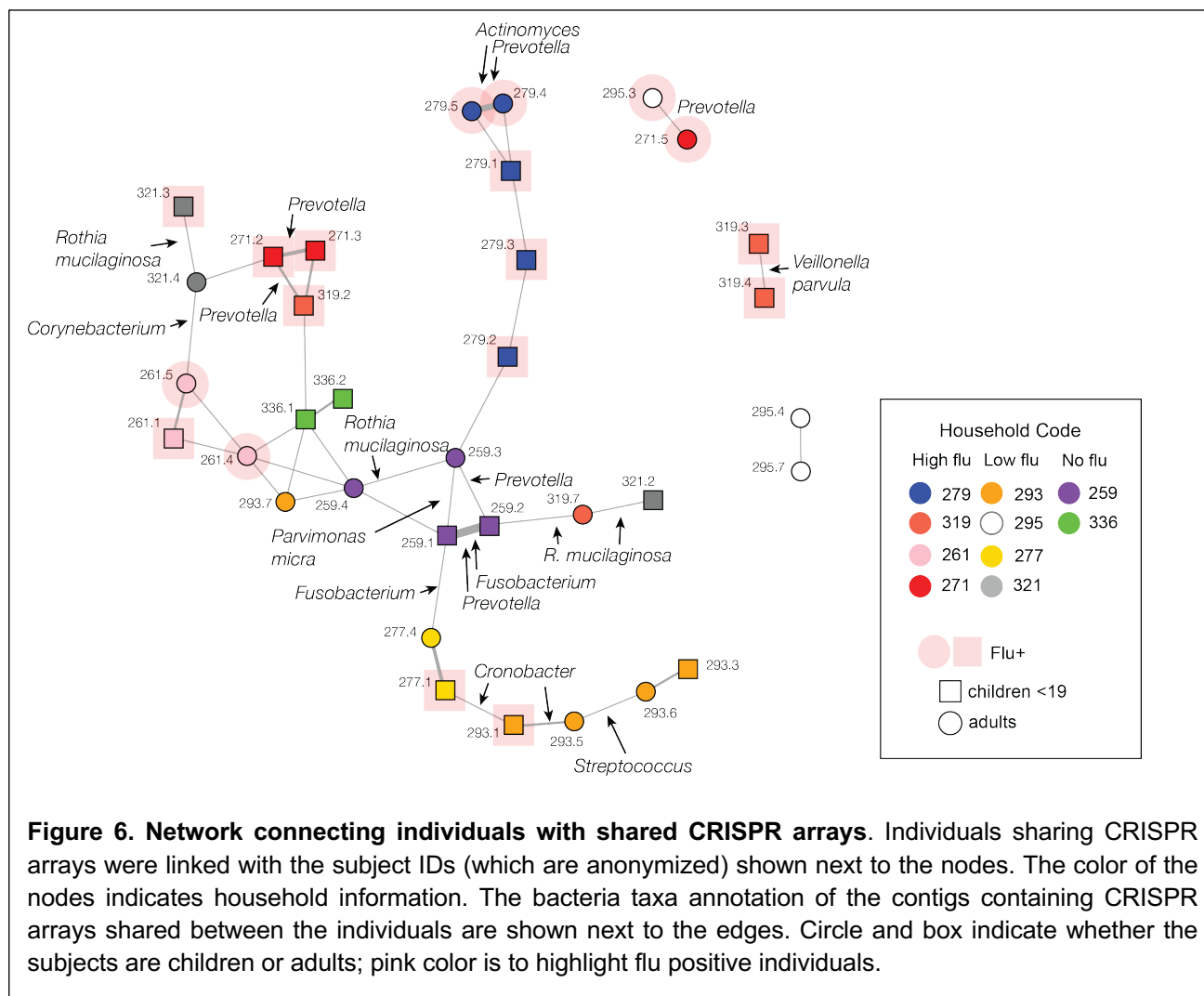
285

286 **Correlation between transmitted bacteria and shared antibiotic resistance genes**

287 Using the spacer profile information, we first tested if there was overrepresentation of shared
288 bacteria within flu infection households as compared to the control households. To do this we
289 removed individual pairs that shared less than 6% of the spacers (**Fig. S5a**), which eliminated all
290 pairs across households and helped focus on within household transmission. We analyzed pairs
291 of individuals within the same households who shared >6% of their spacers (**Fig. S5b**) and
292 compared across the flu infection groups the proportion of individuals in each household who
293 passed this filter. Since we only had 10 households and there were variations in each flu infection
294 group for the data we compared (**Fig. S5c**), we were not able to determine whether there was a
295 correlation between bacteria transmission and flu infection levels.

296

297 We constructed a network to analyze the individuals that shared CRISPR arrays (**Fig. 6**). We
298 looked at antibiotic use history (**Fig. S6**), influenza infection status, and age (we used 18 years
299 old as the cut off for adults). We found that individuals connected in the network could be flu
300 positive or flu negative, with no over representation in either group (**Fig. 6**). To detect the presence
301 of antibiotic resistance genes, we assembled contigs across individuals (different time points from
302 the same individual were pooled together) for both metagenomics and metatranscriptomics



303 datasets. Antibiotic resistance genes were identified by aligning the predicted open reading
 304 frames on the assembled contigs to the CARD database [27]. The sequences predicted as
 305 antibiotic resistance genes were also filtered such that the reference genes in the database were
 306 at least 80% covered over their sequence length. By mapping the reads from the samples back
 307 to the antibiotic resistance genes detected in each subject, we quantified the antibiotic resistance
 308 genes by relative abundance. We aggregated the relative abundance values to the level of gene
 309 families that summarized the gene variants. We applied DESeq2 on the samples for both
 310 metagenomics and metatranscriptomics datasets to identify antibiotic resistance genes that were
 311 enriched or differentially expressed during flu infection. TEM and CfxA4 beta lactamase, as well

312 as sulfonamide resistant genes (sul) were differentially abundant in the metagenomics datasets
313 (FDR p-value < 0.05) (**Fig. S7**) but with a small log₂ fold change, indicating a marginal yet
314 statistically-significant difference. We did not detect any genes differentially expressed in the
315 metatranscriptomics dataset. This indicates that influenza infection does not appear to be
316 associated with increased presence of antibiotic resistance.

317
318 We used the antibiotic resistance gene profiles to calculate the fraction of individuals with a
319 specific antibiotic resistance gene in each household, compared across the high influenza
320 infection, low influenza infection and the control households for both metagenomics and
321 metatranscriptomics datasets. We did not find enrichment of any genes in any group, indicating
322 there was no differential patterns in the presence of antibiotic resistance genes among
323 households (**Fig. S8**). Some antibiotic resistance genes were annotated with bacteria origins
324 using the contigs to which the genes belonged (**Table S4**). We compared the antibiotic resistance
325 gene profiles between individuals from the same households and grouped individuals on whether
326 they were connected by shared bacteria or not. We observed that individuals connected by shared
327 bacteria had more similar antibiotic resistance gene profiles than with other individuals from the
328 same households, demonstrated by a smaller dissimilarity index (the two distributions were tested
329 to be different with a p-value of 0.003) (**Fig. S9**). This implies that there might be transmission of
330 antibiotic resistance genes together with the transmission of bacteria.

331 332 ***Discussion***

333 The respiratory tract microbiome, because of its function in health [25], should play an important
334 role during respiratory tract infections. Here, we generated metagenomic and metatranscriptomic
335 datasets from nasal and throat swabs to profile bacteria taxa and investigate ecological and
336 functional aspects of the microbiome in the upper airways, as well as establish whether
337 transmission of bacteria could be tracked using the metagenomics data.

338

339 We observe that influenza positive households were significantly different in microbial
340 composition from the control (flu negative) households with a few bacteria and phages
341 differentially abundant between the groups. We observed enrichment of *Moraxella* in the high flu
342 infection households and flu positive samples. *Moraxella* has previously been found to be
343 enriched in influenza infection samples [3], supporting our observation. *Moraxella Catarrhalis* was
344 also enriched in the flu infection time points compared to the baseline in the same individuals.
345 *Moraxella catarrhalis* is a common respiratory tract bacteria usually found in children that can
346 potentially cause infections and lead to pneumonia [28]. However, one interesting observation is
347 that the enrichment of certain bacteria that were previously thought to be associated with influenza
348 infection, such as *Dolosigranulum* [3], were found in our study to be enriched in the control
349 households as well, indicating that in some instances there is a household effect for the bacteria
350 present. When comparing between influenza positive and negative time points from the same
351 individuals, we found that the microbial diversity was not significantly different. A recent study
352 found divergence in the respiratory tract microbiome in ferrets over 14 days post influenza
353 infection, where the microbiome at the initial infection time points and day 14 were more like the
354 microbiome in the uninfected ferrets [4]. In our study, although influenza viruses were not
355 detectable after a few days, the microbiome remained similar to the infected time points, although
356 it is possible the microbiome would have gotten back to baseline at later time points.

357

358 We identified microbial genes that were differentially enriched in flu positive and flu negative
359 samples that are involved in bacteria physiology. This includes the two-component systems that
360 are ubiquitous in bacteria and found to regulate virulence and antibiotic resistance in bacteria [29].
361 Thus, the disruption of the microbiome by influenza infection could potentially affect the balance
362 between bacteria competition and bacteria-host interactions. *Moraxella* was found to contribute
363 the most to these differentially abundant genes, indicating not only its relative abundance is

364 disrupted during influenza infection but also effects on its functional potential. We also
365 investigated whether antibiotic resistance genes were associated with influenza infection,
366 however, we did not find a strong relationship between the relative abundance or expression of
367 antibiotic resistance genes and influenza infection as ARGs can be detected as often in the control
368 group.

369

370 The use of CRISPR arrays to track bacteria movement was previously done in a study
371 investigating environmental bacteria strains on different continents [30]. However, our study is the
372 first to use it to track bacteria transmission using clinical samples and in short-range transmissions
373 within and between households. By using the spacers and CRISPR arrays we have shown how
374 we could identify more transmitted bacteria within than across households, as would be expected.
375 However, the number of shared CRISPR arrays between households also indicates that this could
376 be an approach to potentially track transmission of bacteria on a larger scale. The individuals
377 connected with shared bacteria include both children and adults, infected with influenza viruses
378 and not infected. Children within households can also drive bacteria sharing as they may have
379 closer contact with other household members. However, we do not have bacteria isolates or
380 longer time points before influenza infection to validate bacteria transmission and to determine
381 whether this happened during influenza infection.

382

383 There are a few limitations in this study. First, while the use of CRISPR arrays did allow the
384 identification of shared bacteria between individuals, not all bacteria species have a CRISPR
385 system [31], thus our analysis is restricted to a limited set of bacteria species. Also, we do not
386 have bacteria isolates paired with the metagenomics datasets to validate the CRISPR array
387 analysis. Second, we were limited by the number of households in the study and thus cannot
388 draw any conclusion between bacteria sharing and influenza infection rate. The households with
389 high or low influenza infection only indicate the members in the households were infected with

390 influenza but we could not track specific transmission of influenza viruses among household
391 members. Other type of evidence, such as bacteria isolates or long read sequencing would be
392 needed to accurately link ARGs to bacteria genomes shared between individuals.

393

394 In conclusion, the analysis of the metagenome and metatranscriptome data demonstrate that the
395 microbiome compositional and functional potentials are altered in influenza infection. In both flu
396 positive and control individuals we saw that commensal bacteria and potential pathobionts can
397 be readily transmitted within and across households. This implies that antibiotic resistance genes
398 could also be transmitted. Finally, we demonstrated CRISPR arrays are a powerful tool to track
399 bacteria transmission between individuals, offering a novel approach to leverage metagenomics
400 datasets.

401

402

403 **Ethical Approval and Consent to participate**

404 This study was approved by the institutional review boards at the Nicaraguan Ministry of Health
405 and the University of Michigan (HUM00091392). Informed consent or parental permission was
406 obtained for all participants.

407

408 **Consent for publication**

409 Not applicable.

410

411 **Availability of data and materials**

412 The sequencing datasets supporting the conclusions of this article are available in the Sequence
413 Read Archive (SRA). Metatranscriptome data and the metagenome data are under BioProject
414 PRJNA713420. Scripts generating the data are available on GitHub
415 (https://github.com/GhedinLab/Nicaragua_microbiome_flu_analysis).

416

417 **Competing interests**

418 The authors declare that they have no competing interests.

419

420 **Funding**

421 This work was made possible by the support in part from the NIAID/National Institutes of Health,
422 U01 [AI111598](#) (EG), U01 AI088654 (AG) and National Institutes of Health under contract
423 numbers HHSN272201400031C and HHSN272201400006C (AG). This work was also supported
424 in part by the Division of Intramural Research (DIR) of the NIAID/NIH (EG).

425

426 **Authors' contributions**

427 The project was conceived by EG and RB. LZ performed sample processing, bioinformatics
428 analyses, and writing of the manuscript. JR assisted with the development of bioinformatics
429 pipelines. AG, AB, and GK designed and conducted the Household Influenza Transmission Study.
430 LL assisted with the sample organization and data management. The authors read and approved
431 the final manuscript.

432

433 **Acknowledgements**

434 We thank the NYU Genomics Core at the Center for Genomics and Systems Biology.

435

436

437 **Material and Method**

438

439 **Data collection**

440 Samples were collected from individuals participating in the Household Influenza Transmission
441 Study (HITS) in Managua, Nicaragua between July 2013 and October 2014. The HITS sample
442 cohort included child index cases enrolled in the Nicaraguan Influenza Cohort Study and their
443 family members who developed influenza as well as some influenza negative control households.
444 Respiratory specimens consisted of pooled nasal and throat swabs collected from household
445 members every 2-4 days over a 9-12 day period. Samples were shipped to the Center for
446 Genomics and Systems Biology, New York University, and stored at – 80 °C. The HITS study was
447 approved by the institutional review boards at the Nicaraguan Ministry of Health and the University
448 of Michigan. Informed consent or parental permission was obtained for all participants and
449 children aged 6 years and older provided assent.

450

451 **RNA extraction and library preparation for metatranscriptome sequencing.**

452 Total RNA was isolated from 500uL of each respiratory sample (nasal washes) with the QIAGEN
453 RNeasy Micro Kit (QIAGEN, Hilden, Germany) according to the manufacturer's recommendations
454 and stored at –80°C. No mRNA enrichment or rRNA depletion steps were performed due to the
455 limited biomass of the starting material. NEBNext Ultra II RNA Library Prep kit (New England
456 Biolabs, Ipswich, MA) was used to generate the metatranscriptomics libraries and each library
457 was subjected to 14-17 cycles of PCR and adaptor concentration was diluted 1:100 or 1:200 to
458 maintain sample and adaptor ratio. Libraries were quantified by qPCR using the KAPA Library
459 Quantification Kit (KAPA Biosystems, Wilmington, MA) on a Roche 480 LightCycler (Roche, Basel,
460 Switzerland); their size distributions were measured on a 4200 TapeStation using a D1000
461 ScreenTape (Agilent Technologies, Santa Clara, CA). Libraries were diluted to 4 nM in dilution

462 buffer (10mM Tris, pH 8.5) and combined with equimolar input into 9 sequencing pools (20-25
463 libraries per pool). Paired-end sequencing (2x150 bp) was performed at the Genomics Core
464 Facility (Center for Genomics and Systems Biology, New York University) on the Illumina NextSeq
465 500 instrument according to the manufacturer's instructions (Illumina, Inc., San Diego, CA) with
466 a few libraries sequenced on the Illumina HiSeq 2500 instrument.

467

468 **DNA isolation and library preparation for metagenome sequencing**

469 Genomic DNA was isolated from the remaining volume of each sample with the PowerSoil DNA
470 Isolation Kit (Qiagen) and stored at -80°C . Libraries were generated using Nextera DNA Flex
471 library prep kit (Illumina, Inc., San Diego, CA). Libraries were quantified by qPCR using the KAPA
472 Library Quantification Kit (KAPA Biosystems, Wilmington, MA) on a Roche 480 LightCycler
473 (Roche, Basel, Switzerland); their size distributions were measured on a 4200 TapeStation using
474 a D1000 ScreenTape (Agilent Technologies, Santa Clara, CA). Libraries were diluted to 4 nM in
475 dilution buffer (10mM Tris, pH 8.5) and combined with equimolar input into 9 sequencing pools
476 (20-25 libraries per pool). Paired-end sequencing (2x150 bp) was performed at the Genomics
477 Core Facility (Center for Genomics and Systems Biology, New York University) on the Illumina
478 NextSeq 500 instrument according to the manufacturer's instructions (Illumina, Inc., San Diego,
479 CA) with a few libraries sequenced on the Illumina HiSeq 2500 instrument.

480

481 **Metagenomics and metatranscriptomics data processing**

482 The metagenomics reads were filtered to remove adaptors and low quality reads using
483 Trimmomatic v0.36 [32] followed by DeconSeq2 v1.32.0 [33] to remove human reads. The
484 metatranscriptomics reads were filtered by Trimmomatic, DeconSeq2 and SortMeRNA v2.1 [34]
485 to remove adaptors, human reads and rRNA reads, respectively. The median reads number for

486 metagenomes was 6.8M(IQR=9.8M) and the median reads number for metatranscriptomes was
487 4.9M (IQR=5.6M) post filtering.

488

489 **Bacterial taxonomic assignments and differential abundant analysis**

490 The filtered metagenomics and metatranscriptomics datasets were run through Kraken2 v2 to
491 classify the reads to bacterial and viral taxa. Beta diversity of the metagenomics and
492 metatranscriptomics datasets was determined using Bray Curtis distance and the global diversity
493 between different groups was determined by PERMANOVA. The bacterial and viral taxa
494 differentially abundant between influenza infection and no infection samples, and different
495 household groups were identified by DESeq2 v 1.34.0 with adjusted p values ≤ 0.05 and a log2
496 fold change greater than 2 or smaller than -2.

497

498 **Bacterial gene and taxa contribution**

499 The Functional Mapping and Analysis Pipeline (FMAP [23]) was used to identify the KEGG gene
500 orthologous groups in the metagenomics and metatranscriptomics datasets. Gene families
501 differentially abundant or expressed between influenza positive and influenza negative samples
502 were identified by using DESeq2 with an adjusted p value smaller than 0.05 or a log2 fold change
503 greater than 2 or smaller than -2. The reads mapped to the KEGG gene orthologous groups that
504 are differentially abundant between influenza positive and influenza negative groups were
505 extracted and mapped to bacteria taxa using Kraken2.

506

507 **Bacteria and virus interaction analysis**

508 The metagenomics bacterial reads identified by Kraken2 were pooled across different samples
509 from the same individual and assembled into contigs using metaSPAdes v3.12.0 [35]. The
510 spacers were mapped back to the bacterial contigs with ≤ 1 mismatch and no secondary mapping.
511 All filtered metagenomics reads from the same individuals across time points were assembled

512 into contigs. VirSorter v1.1.0 [36] was used to predict the viral contigs and we focused on the
513 contigs predicted as category 1 and category 2, which are high confidence viral contigs. The
514 spacers were mapped back to the viral contigs using BLAST with parameters modified for short
515 query sequences [37]. The viral and bacterial contigs were searched by BLAST against the NCBI
516 nt database and viral Refseq database, respectively, to get taxonomy assignments. The bacterial
517 contigs were annotated to the level of species or genus if more than 80% of the alignments on
518 that contig were annotated as that species or genus. The viral assignments were filtered by 90%
519 sequence similarity with the best hits (longest alignment length and highest identity). The viral
520 and bacterial contigs linked by at least 5 spacers were identified and presented. Figure S10
521 provides an overview of the analysis steps.

522

523 **CRISPR spacer analysis and network analysis**

524 The spacers were identified from each metagenomics dataset using MetaCRIST [38]. The
525 spacers across all the samples were clustered and spacers with sequence similarity greater than
526 90% using CD-HIT [39] were determined as being the same spacers across samples. The percent
527 of shared spacers between samples was determined as the number of shared spacers between
528 any two samples divided by the average of total spacers in the two samples. Percent of shared
529 spacers was compared between samples from the same individuals, different individuals in the
530 same households and different households. When comparing at the individual level, the spacers
531 from different time points for the same individual were combined to do the analysis. The network
532 connecting individuals based on shared spacers was generated using igraph [40] in R studio and
533 the edge weight was the percent of shared spacers between individuals. Subnetworks were also
534 analyzed using igraph.

535

536 **Identification of bacteria shared between individuals by using CRISPR arrays**

537 The spacers and the order of the spacers were determined by mapping the spacers to the
538 bacterial contigs using bowtie2 v2.2.4 [41]. A dynamic programming Smith-Waterman algorithm
539 was used to find the best alignments between any two CRISPR arrays. We tested a few
540 parameters for gap opening score, match score and mismatch score. Since the results were very
541 similar for alignment length and alignment similarity (**Fig. S3**), we chose the parameters that led
542 to the highest number of alignments between CRISPR arrays (gap score=-1, match score=3,
543 mismatch score=-2). The CRISPR array alignments were analyzed for alignment similarity and
544 alignment length. 2D density distribution was estimated based on the alignment similarity and
545 alignment length using the ks package [42] in R studio. The densities from the same household
546 and different households were compared, and regions on the density plot enriched with data from
547 the same household or different households were identified.

548

549 **Antibiotic resistance gene profiling and household transmission analysis**

550 All filtered metagenomics and metatranscriptomic reads from the same individuals across time
551 points were assembled into contigs using metaSPAdes [35]. Open reading frames from each
552 contig were predicted by using MetaProdigal v2.6.3 [43]. The ORFs were aligned to the CARD
553 v3.1.0 [27] database for ARG annotation. ORFs that covered at least 80% of the ARG reference
554 sequences were identified as present. The contigs carrying the ARGs were mapped to the NCBI
555 nt database; the top hits, and the taxonomic information and definition for the top hits, were
556 extracted. The metagenomic and metatranscriptomic sequence reads were mapped back to the
557 ARGs identified in their respective datasets. The number of reads mapped back to the antibiotic
558 resistance genes were identified and differential abundance and expression analysis were done
559 using DESeq2.

560

561 The beta diversity of the presence/absence of the antibiotic resistance gene profiles were
562 measured between individuals. The dissimilarities (measured by Jaccard similarity index)

563 between same household individuals connected by shared bacteria were compared to the
564 dissimilarities between same household individuals not connected by bacteria. The taxa origins
565 of the antibiotic resistance genes were investigated by searching by BLAST the contigs that carry
566 antibiotic resistance genes against the NCBI nt database; the top hits were recorded.

567

568

569

570

571 **Table 1. Summary table for individuals across the households**

	Flu infection			
	High Infection N=4	Low Infection N=4	No infection N=2	p values
Number of Households				
Number of individuals	23	23	6	
Age Median (sd)	8 (6.5)	13 (9.5)	7 (8.3)	0.0981
Gender				1
Female (%)	16 (70)	17 (70)	4 (70)	
Male (%)	7 (30)	6 (30)	2 (30)	

572 The p values were calculated by using ANOVA and Fisher's Exact tests.

573

574

575

576 **Table 2. Summary table for flu infection and no infection samples**

577

	Flu infection		
	Infection N=42	No infection N=93	p values
Age Median (sd)	6 (5.4)	14 (8.2)	4.47E-07
Gender			0.02051
Female (%)	15 (47)	74 (80)	
Male (%)	17 (53)	19 (20)	

578 The p values were calculated using ANOVA and Fisher's Exact tests

579

580 **Figure legends**

581

582 **Figure 1. Differential abundance of bacteria and phages between households and samples.**

583 **(a)** PCA plot of β diversity of the microbial composition from metagenomics datasets for different
584 influenza infection households. Red indicates the high flu infection households, pink indicate low
585 flu infection households, and grey indicates control households. **(b)** PCA plot of β diversity of the
586 microbial composition from metatranscriptomics datasets for different influenza infection
587 households. Red indicates the high flu infection households, pink indicate low flu infection
588 households, and grey indicates control households. **(c)** Differential abundance of bacteria genera
589 between different household groups in metagenomes. Red or pink indicate higher differential
590 abundance in the flu infection group tested while grey indicates it is higher in the control group.
591 **(d)**. Differential abundance of bacteria genera between different household groups in
592 metatranscriptomes. Red or pink indicate higher differential abundance in the flu infection group
593 tested while grey indicates it is higher in the control group. **(e)** PCA plot of β diversity of the
594 microbial composition from metagenomics datasets for flu negative individuals from flu infection
595 or control households. Red indicates the flu infection households and grey indicates control
596 households. **(f)** PCA plot of β diversity of the microbial composition from metatranscriptomics
597 datasets for flu negative individuals from flu infection or control households. Red indicates the flu
598 infection households and grey indicates control households. **(g)** Differential abundance of
599 bacterial genera between no flu infection individuals from flu infection and control households.
600 Red indicates higher differential abundance in the flu infection households while grey indicates it
601 is higher in the control group.

602

603

604 **Figure 2. Differential abundance of bacteria and phages between influenza infection status.**

605 **(a)** PCA plot of β diversity of the microbial composition from flu positive and flu negative samples
606 from metagenomics and metatranscriptomics datasets. Red indicates flu positive samples and
607 grey indicates flu negative samples. **(b)** Bacterial genera differentially abundant in metagenomes
608 between flu positive and flu negative samples; red indicates the taxa are enriched in flu positive
609 samples and grey enriched in flu negative samples. **(c)** Log₂ fold change of *Moraxella* species
610 between flu positive and flu negative time points from the same individuals in metagenomes and
611 metatranscriptomes. The red triangles indicate enrichment in the flu positive time points.

612

613 **Figure 3. Bacteria genes and pathways altered in influenza infection.** Microbial gene

614 orthologous groups were profiled for each sample. Microbial orthologous genes differentially
615 abundant in metagenomes between flu positive samples and flu negative samples from all
616 households were identified using DESeq2 with a p value smaller than 0.05 and a log₂ fold change
617 larger than 2 or smaller than -2. **(a)** The dot plot shows differentially abundant genes between flu
618 positive and flu negative samples. The size of the dot indicates $-\log_{10}$ of the adjusted p values.
619 Red indicates that the relative abundance of the gene is higher in flu positive samples and grey
620 indicates it is higher in flu negative samples. **(b)** The reads mapped to the differentially abundant
621 genes in (a) were extracted and their taxonomy classifications were determined through Kraken2.
622 *Moraxella* contributes the most to these genes. The graph shows the boxplot of counts per million
623 for the reads mapped both to *Moraxella catarrhalis* and the genes on the y axis. The samples
624 were separated into the flu positive or flu negative groups. Red indicates flu positive samples and
625 grey indicates flu negative samples.

626

627 **Figure 4. CRISPR spacers shared between samples and individuals.** The percent of shared

628 spacers between samples or individuals were compared and used to construct the connection
629 network between individuals. **(a)** Boxplot indicating percent of spacers shared between samples

630 from the same individuals, different individuals in the same households, and individuals from
631 different households. The colors indicate whether the samples are from the same households
632 (purple) or different households (orange). **(b)** Density plot and boxplot for percent of spacers
633 shared at the individual level within and between households. The black line on the density plot
634 indicates the cut-off where there is the same number of comparisons within and between
635 households. **(c)** The connection network was generated based on the percent of shared spacers
636 between individuals for the data above the cut-off in (b). The nodes represent individuals, and the
637 edges represent percent of shared spacers. Same color nodes indicate individuals come from the
638 same household and the numbers on the nodes represent the subnetwork they were partitioned
639 into.

640

641 **Figure 5. CRISPR array alignment distribution within and between households and bacteria**

642 **taxa shared between individuals. (a)** Contour plots for the 2D density distribution of CRISPR

643 array alignments on alignment similarity and alignment length. The density distribution for the

644 CRISPR array alignments with CRISPR arrays from individuals from the same households or

645 different households were colored in purple and orange, respectively. The numbers on the contour

646 plot indicate the regions have 25%, 50% and 75% of the data. The solid color regions on the plot

647 indicate the density of the data in these regions were significantly different between the two groups.

648 The purple region has higher data density in the same household group and the orange region

649 has higher data density in the different household group. **(b)** The contigs contain the CRISPR

650 arrays from the alignments with a similarity greater than 0.9 and more than 15 spacers mapped

651 to the NCBI nt database. The graph shows the barplot of the number of individual pairs that share

652 the bacteria with the color indicating they are from the same households or different households.

653

654 **Figure 6. Network connecting individuals with shared CRISPR arrays.** Individuals sharing

655 CRISPR arrays were linked with the subject IDs shown next to the nodes. The color of the nodes

656 indicates household information. The bacteria taxa annotation of the contigs containing CRISPR
657 arrays shared between the individuals are shown next to the edges. Circle and box indicate
658 whether the subjects are children or adults; pink color is to highlight flu positive individuals.
659
660

661 **Supplementary Figures**

662

663 **Figure S1. Bacteria taxa mapped to genes differentially abundant between flu positive and**

664 **flu negative samples.** Dotplot for the bacteria taxa mapped to genes found to be differentially

665 abundant between flu+ and flu- samples with bacteria taxa on the x axis and genes shown on the

666 y axis. For each gene, the color intensity of the dots shows the fraction of reads mapped to the

667 bacteria taxa relative to total reads contributing to that gene on the y axis.

668

669 **Figure S2. Bacteria and phage interaction network.** Bacteria and viruses were linked by

670 spacers mapped back to both bacterial and viral contigs. The graph shows the presence of

671 interactions between specific bacteria and phages. The phages are shown on the x axis and

672 bacteria on the y axis. The color intensity reflects the number of individuals for which the specific

673 bacteria/phage interactions can be detected. The bacteria and phages in the same genera or

674 found to infect bacteria in the same genera are grouped into same panels.

675

676 **Figure S3. Parameters tested for dynamic programming to align CRISPR arrays.** Different

677 parameters were tested in a local alignment method to find the best alignments between any two

678 CRISPR arrays. With the alignment similarity and alignment length distributions, correlations

679 between methods and statistics are shown in the figure.

680

681 **Figure S4. Contour plot with 2D density of the CRISPR array alignments.** The 2D density

682 distribution was estimated for the CRISPR array alignments on the alignment similarity and

683 alignment length. The color intensity indicates the density of the data in the regions. The two

684 panels indicate the CRISPR array alignments with CRISPR arrays from the same or different

685 households. The y axis was adjusted to highlight the region where most of the data points are

686 located, although there were more points over 50 on the y axis.

687

688 **Figure S5. Bacteria transmission and flu infection.** (a) Density and boxplot plot for percent of
689 spacers shared at the individual level within and between households. The red line on the density
690 plot indicates the cut-off where all the “between household” individual pairs were removed. (b)
691 The connection network was generated based on the percent of shared spacers between
692 individuals for the data above the cut-off in (a). The nodes represent individuals and the edges
693 represent percent of shared spacers. Same color nodes represent individuals come from the
694 same household. (c) Dotplot for percent of individuals in each household that were connected.
695 Number of individuals in (b) in each household were divided by total number of individuals in the
696 households and compared across flu infection groups. The x axis indicates household code and
697 the panels show the household from high, low, or no flu infection groups.

698

699 **Figure S6. Antibiotic use history.** The antibiotics taken by the subjects within 2 years prior to
700 influenza infection is shown in red color. White indicates that particular antibiotic was not taken
701 by the subject, and grey indicates antibiotics use information for that individual is not available
702 (NA).

703

704 **Figure S7. Antibiotic resistance genes differentially abundant between flu positive and flu**
705 **negative samples.** Barplot for counts per million of the ARGs identified as differentially abundant
706 between flu positive and flu negative samples. The color indicates the influenza positive (flu+) or
707 influenza negative (flu-) groups. The panel titles indicate the mechanism used by the genes to
708 confer antibiotic resistance. The log₂ fold change of ARGs counts for the identified antibiotic
709 resistance genes are also shown with barplot.

710

711

712

713 **Figure S8. Prevalence of antibiotic resistance genes in households across flu infection**

714 **groups.** Presence/absence of antibiotic resistance genes was determined for each individual.

715 The ratio of individuals that have a specific ARG in each household was calculated by dividing

716 the number of individuals carrying the gene with the total number of individuals in the households.

717 The graph shows a dotplot for the fraction of individuals that have ARGs across the flu infection

718 groups for metagenomics (MG) and metatranscriptomics (MT) datasets. The color indicates the

719 data type, red for MG and blue for MT. H, L or N panel titles correspond to high, low, and no flu

720 infection groups, respectively. The color intensity indicates the fraction of individuals that have

721 the ARGs in the household. The red horizontal panel titles indicate the mechanisms used by the

722 ARGs to confer antibiotic resistance.

723

724 **Figure S9. Boxplot for ARG profiles dissimilarity between individuals from same**

725 **households connected by shared bacteria, same households not connected and different**

726 **households.** Each dot represents the dissimilarity in ARG profiles between two individuals. The

727 color indicates the individuals were from the same or different households. The x axis shows

728 whether the individuals were connected by shared bacteria when they are from the same

729 household. The y axis shows the distance in ARG profiles between any two individuals.

730

731 **Figure S10. Overview of CRISPR array analysis steps.** CRISPR arrays were used for two

732 purposes: (1) the spacers that comprise the CRISPR arrays on bacterial contigs and that can be

733 mapped to viral contigs, were used to link bacteria and phages to study phage-bacteria

734 interactions; and (2) the spacers and the CRISPR arrays were used as a barcode to study bacteria

735 transmission within and across households.

736

737 **Supplementary Tables**

738

739 **Table S1. Sample metadata.**

740

741 **Table S2. Human DNA viruses identified from the metagenomics datasets.** The human
742 viruses were identified from Kraken analysis and the presence or absence of the viruses are
743 shown as 1 or 0.

744

745 **Table S3. Human viruses identified from the metatranscriptomics datasets.** The human
746 viruses were identified from Kraken analysis and the presence or absence of the viruses are
747 shown as 1 or 0.

748

749 **Table S4. Antibiotic resistance genes and contig annotations.** The contigs containing the
750 ARGs were annotated with the bacteria genome information.

751

752 References

- 753
754 1. Morris, D.E., D.W. Cleary, and S.C. Clarke, *Secondary Bacterial Infections Associated*
755 *with Influenza Pandemics*. Front Microbiol, 2017. **8**: p. 1041.
- 756 2. Rosas-Salazar, C., et al., *Differences in the Nasopharyngeal Microbiome During Acute*
757 *Respiratory Tract Infection With Human Rhinovirus and Respiratory Syncytial Virus in*
758 *Infancy*. J Infect Dis, 2016. **214**(12): p. 1924-1928.
- 759 3. Ding, T., et al., *Microbial Composition of the Human Nasopharynx Varies According to*
760 *Influenza Virus Type and Vaccination Status*. mBio, 2019. **10**(4).
- 761 4. Kaul, D., et al., *Microbiome disturbance and resilience dynamics of the upper respiratory*
762 *tract during influenza A virus infection*. Nat Commun, 2020. **11**(1): p. 2537.
- 763 5. Hagan, T., et al., *Antibiotics-Driven Gut Microbiome Perturbation Alters Immunity to*
764 *Vaccines in Humans*. Cell, 2019. **178**(6): p. 1313-1328 e13.
- 765 6. Clemente, J.C., et al., *The microbiome of uncontacted Amerindians*. Sci Adv, 2015. **1**(3).
- 766 7. Zhang, L., et al., *Characterization of antibiotic resistance and host-microbiome*
767 *interactions in the human upper respiratory tract during influenza infection*. Microbiome,
768 2020. **8**(1): p. 39.
- 769 8. Weiser, J.N., D.M. Ferreira, and J.C. Paton, *Streptococcus pneumoniae: transmission,*
770 *colonization and invasion*. Nat Rev Microbiol, 2018. **16**(6): p. 355-367.
- 771 9. McCullers, J.A., et al., *Influenza enhances susceptibility to natural acquisition of and*
772 *disease due to Streptococcus pneumoniae in ferrets*. J Infect Dis, 2010. **202**(8): p. 1287-
773 95.
- 774 10. Ribet, D. and P. Cossart, *How bacterial pathogens colonize their hosts and invade*
775 *deeper tissues*. Microbes Infect, 2015. **17**(3): p. 173-83.
- 776 11. Genoyer, E. and C.B. López, *The Impact of Defective Viruses on Infection and Immunity*.
777 Annual Review of Virology, 2019. **6**: p. 547-566.
- 778 12. Chatterjee, A. and B.A. Duerkop, *Beyond Bacteria: Bacteriophage-Eukaryotic Host*
779 *Interactions Reveal Emerging Paradigms of Health and Disease*. Front Microbiol, 2018.
780 **9**: p. 1394.
- 781 13. Norman, J.M., et al., *Disease-specific alterations in the enteric virome in inflammatory*
782 *bowel disease*. Cell, 2015. **160**(3): p. 447-60.
- 783 14. Karginov, F.V. and G.J. Hannon, *The CRISPR system: small RNA-guided defense in*
784 *bacteria and archaea*. Mol Cell, 2010. **37**(1): p. 7-19.
- 785 15. Truong, D.T., et al., *Microbial strain-level population structure and genetic diversity from*
786 *metagenomes*. Genome Res, 2017. **27**(4): p. 626-638.
- 787 16. Nayfach, S., et al., *An integrated metagenomics pipeline for strain profiling reveals novel*
788 *patterns of bacterial transmission and biogeography*. Genome Res, 2016. **26**(11): p.
789 1612-1625.
- 790 17. Costea, P.I., et al., *metaSNV: A tool for metagenomic strain level analysis*. PLoS One,
791 2017. **12**(7): p. e0182392.
- 792 18. Paez-Espino, D., et al., *Strong bias in the bacterial CRISPR elements that confer*
793 *immunity to phage*. Nat Commun, 2013. **4**: p. 1430.
- 794 19. Heler, R., et al., *Spacer Acquisition Rates Determine the Immunological Diversity of the*
795 *Type II CRISPR-Cas Immune Response*. Cell Host Microbe, 2019. **25**(2): p. 242-249 e3.
- 796 20. Wood, D.E., J. Lu, and B. Langmead, *Improved metagenomic analysis with Kraken 2*.
797 Genome Biol, 2019. **20**(1): p. 257.
- 798 21. Anderson, M.J., *Permutational Multivariate Analysis of Variance (PERMANOVA)*, in
799 *Wiley StatsRef: Statistics Reference Online*. 2017. p. 1-15.
- 800 22. Love, M.I., W. Huber, and S. Anders, *Moderated estimation of fold change and*
801 *dispersion for RNA-seq data with DESeq2*. Genome Biol, 2014. **15**(12): p. 550.

- 802 23. Kim, J., et al., *FMAP: Functional Mapping and Analysis Pipeline for metagenomics and*
803 *metatranscriptomics studies*. BMC Bioinformatics, 2016. **17**(1): p. 420.
- 804 24. Mitrophanov, A.Y. and E.A. Groisman, *Signal integration in bacterial two-component*
805 *regulatory systems*. Genes Dev, 2008. **22**(19): p. 2601-11.
- 806 25. Man, W.H., W.A. de Steenhuijsen Piters, and D. Bogaert, *The microbiota of the*
807 *respiratory tract: gatekeeper to respiratory health*. Nat Rev Microbiol, 2017. **15**(5): p.
808 259-270.
- 809 26. de Steenhuijsen Piters, W.A., E.A. Sanders, and D. Bogaert, *The role of the local*
810 *microbial ecosystem in respiratory health and disease*. Philos Trans R Soc Lond B Biol
811 Sci, 2015. **370**(1675).
- 812 27. Alcock, B.P., et al., *CARD 2020: antibiotic resistome surveillance with the*
813 *comprehensive antibiotic resistance database*. Nucleic Acids Res, 2020. **48**(D1): p.
814 D517-D525.
- 815 28. Murphy, T.F. and G.I. Parameswaran, *Moraxella catarrhalis, a human respiratory tract*
816 *pathogen*. Clin Infect Dis, 2009. **49**(1): p. 124-31.
- 817 29. Huffnagle, G.B., R.P. Dickson, and N.W. Lukacs, *The respiratory tract microbiome and*
818 *lung inflammation: a two-way street*. Mucosal Immunol, 2017. **10**(2): p. 299-306.
- 819 30. Lopatina, A., et al., *Natural diversity of CRISPR spacers of Thermus: evidence of local*
820 *spacer acquisition and global spacer exchange*. Philos Trans R Soc Lond B Biol Sci,
821 2019. **374**(1772): p. 20180092.
- 822 31. Burstein, D., et al., *Major bacterial lineages are essentially devoid of CRISPR-Cas viral*
823 *defence systems*. Nat Commun, 2016. **7**: p. 10613.
- 824 32. Bolger, A.M., M. Lohse, and B. Usadel, *Trimmomatic: a flexible trimmer for Illumina*
825 *sequence data*. Bioinformatics, 2014. **30**(15): p. 2114-20.
- 826 33. Schmieder, R. and R. Edwards, *Fast identification and removal of sequence*
827 *contamination from genomic and metagenomic datasets*. PLoS One, 2011. **6**(3): p.
828 e17288.
- 829 34. Kopylova, E., L. Noe, and H. Touzet, *SortMeRNA: fast and accurate filtering of*
830 *ribosomal RNAs in metatranscriptomic data*. Bioinformatics, 2012. **28**(24): p. 3211-7.
- 831 35. Nurk, S., et al., *metaSPAdes: a new versatile metagenomic assembler*. Genome Res,
832 2017. **27**(5): p. 824-834.
- 833 36. Roux, S., et al., *VirSorter: mining viral signal from microbial genomic data*. PeerJ, 2015.
834 **3**: p. e985.
- 835 37. Edwards, R.A., et al., *Computational approaches to predict bacteriophage-host*
836 *relationships*. FEMS Microbiol Rev, 2016. **40**(2): p. 258-72.
- 837 38. Moller, A.G. and C. Liang, *MetaCRAB: reference-guided extraction of CRISPR spacers*
838 *from unassembled metagenomes*. PeerJ, 2017. **5**: p. e3788.
- 839 39. Li, W. and A. Godzik, *Cd-hit: a fast program for clustering and comparing large sets of*
840 *protein or nucleotide sequences*. Bioinformatics, 2006. **22**(13): p. 1658-9.
- 841 40. Ju, W., et al., *iGraph: an incremental data processing system for dynamic graph*.
842 Frontiers of Computer Science, 2016. **10**(3): p. 462-476.
- 843 41. Langmead, B. and S.L. Salzberg, *Fast gapped-read alignment with Bowtie 2*. Nat
844 Methods, 2012. **9**(4): p. 357-9.
- 845 42. Duong, T., *ks: Kernel Density Estimation and Kernel Discriminant Analysis for*
846 *Multivariate Data in R* Journal of Statistical Software, 2007. **21**(7).
- 847 43. Hyatt, D., et al., *Gene and translation initiation site prediction in metagenomic*
848 *sequences*. Bioinformatics, 2012. **28**(17): p. 2223-30.
- 849