

1 **Microbial DNA on the move: sequencing based detection and analysis**
2 **of transduced DNA in pure cultures and microbial communities**
3

4 Manuel Kleiner¹, Brian Bushnell², Kenneth E. Sanderson³, Lora V. Hooper^{4,5}, Breck A. Duerkop⁶
5

6 Affiliations:

7 1: Department of Plant and Microbial Biology, North Carolina State University, Raleigh, NC, USA

8 2: Department of Energy, Joint Genome Institute, Walnut Creek, CA, USA

9 3: Department of Biological Sciences, University of Calgary, Calgary, AB, Canada

10 4: Department of Immunology, University of Texas Southwestern Medical Center, Dallas, TX, USA

11 5: Howard Hughes Medical Institute, University of Texas Southwestern Medical Center, Dallas, TX, USA

12 6: Department of Immunology and Microbiology, University of Colorado School of Medicine, Aurora,
13 CO, USA

14

15

16

17 Correspondence:

18 Manuel Kleiner – manuel_kleiner@ncsu.edu

19 Breck Duerkop – breck.duerkop@cuanschutz.edu

20 **Abstract**

21 Horizontal gene transfer (HGT) plays a central role in microbial evolution. Our understanding of the
22 mechanisms, frequency and taxonomic range of HGT in polymicrobial environments is limited, as we
23 currently rely on historical HGT events inferred from genome sequencing and studies involving cultured
24 microorganisms. We lack approaches to observe ongoing HGT in microbial communities. To address this
25 knowledge gap, we developed a DNA sequencing based “transductomics” approach that detects and
26 characterizes microbial DNA transferred via transduction. We validated our approach using model
27 systems representing a range of transduction modes and show that we can detect numerous classes of
28 transducing DNA. Additionally, we show that we can use this methodology to obtain insights into DNA
29 transduction among all major taxonomic groups of the intestinal microbiome. This work extends the
30 genomic toolkit for the broader study of mobile DNA within microbial communities and could be used to
31 understand how phenotypes spread within microbiomes.

32 **Significance Statement**

33 Microbes can rapidly evolve new capabilities by acquiring genes from other organisms through a process
34 called horizontal gene transfer (HGT). HGT occurs via different routes, one of which is by the transfer of
35 DNA carried by microbe infecting viruses (phages) or virus-like agents. This process is called
36 transduction and has primarily been studied in the lab using pure cultures or indirectly in environmental
37 communities by analyzing signatures in microbial genomes revealing past transduction events. The
38 transductomics approach that we present here, allows for the detection and characterization of genes that
39 are potentially transferred between microbes in complex microbial communities at the time of
40 measurement and thus provides insights into real-time ongoing horizontal gene transfer.

41 Introduction

42 The importance of horizontal gene transfer (HGT) as a driver of rapid evolution and adaptation in
43 microbial communities and host-associated microbiomes has become increasingly recognized(1, 2).
44 Publicly available genomes and metagenomes have revealed pervasive horizontally acquired genes in
45 almost all available genomes. A study of HGT in the human microbiome, for example, showed >10,000
46 recently transferred genes in 2,235 analyzed genomes(3). HGT has been implicated in the spread of
47 antibiotic resistance genes(4), toxin and other virulence genes(5, 6), as well as genes that enable digestion
48 of dietary compounds by microbes in the intestine(7), and metabolic genes that augment microbial
49 metabolism with critical functions in environmental populations(8). Despite its recognized importance,
50 our understanding of the taxonomic range, frequency, and mechanisms of HGT are still limited. Most
51 studies of HGT in microbiomes rely on analysis of microbial genomes(3, 9) and as such these methods
52 attempt to reconstruct historical HGT. What we currently lack are methods that measure ongoing HGT
53 and identify the mechanism of DNA transfer. Here we present a novel method that specifically determines
54 the sequence of DNA that is transferred between cells via one of the major known pathways for DNA
55 transfer – transduction.

56 Currently, there are three major ways that genetic material is known to be exchanged between microbial
57 cells, (1) transformation – uptake of DNA by naturally competent cells, (2) conjugation – exchange of
58 genetic material (e.g. plasmids) using direct contact between donor and recipient cells, and (3)
59 transduction – transfer of genetic material by viruses or virus-like particles (VLPs)(2). Here we focus on
60 transduction only. There are several known types of transduction including classic specialized and
61 generalized transduction, and more recently discovered types, including gene transfer agents (GTAs),
62 lateral transduction and hijacking of bacteriophage (phage) particles by genomic islands(10–12). During
63 specialized transduction DNA adjacent to prophage integration sites in the bacterial genome are co-
64 excised at a low frequency and packaged into phage heads after prophage genome replication. In
65 generalized transduction non-random pieces of the host bacterial genomic DNA or plasmids get packaged
66 at low frequency into phage particles when a lytic phage infects and replicates in a bacterial cell. This non-
67 random packaging is mediated by genomic features that resemble the packaging site (*pac* site) on the
68 phage genome, which is used by the phage particle packaging machinery as the start site phage DNA
69 packaging into the capsid(13). In lateral transduction prophages replicate while still integrated in the host
70 genome and prophage packaging initiates *in situ* ultimately leading to high frequency packaging of host
71 DNA in a unidirectional fashion away from the prophage integration site(12).GTAs are phage-like
72 particles encoded in bacterial genomes that package random pieces of the genomic DNA upon production
73 and can transfer these pieces to other cells(10). In contrast to phages, GTAs do not carry the DNA content
74 sufficient to support their reproduction in the target cells. Lastly, some genomic islands, including

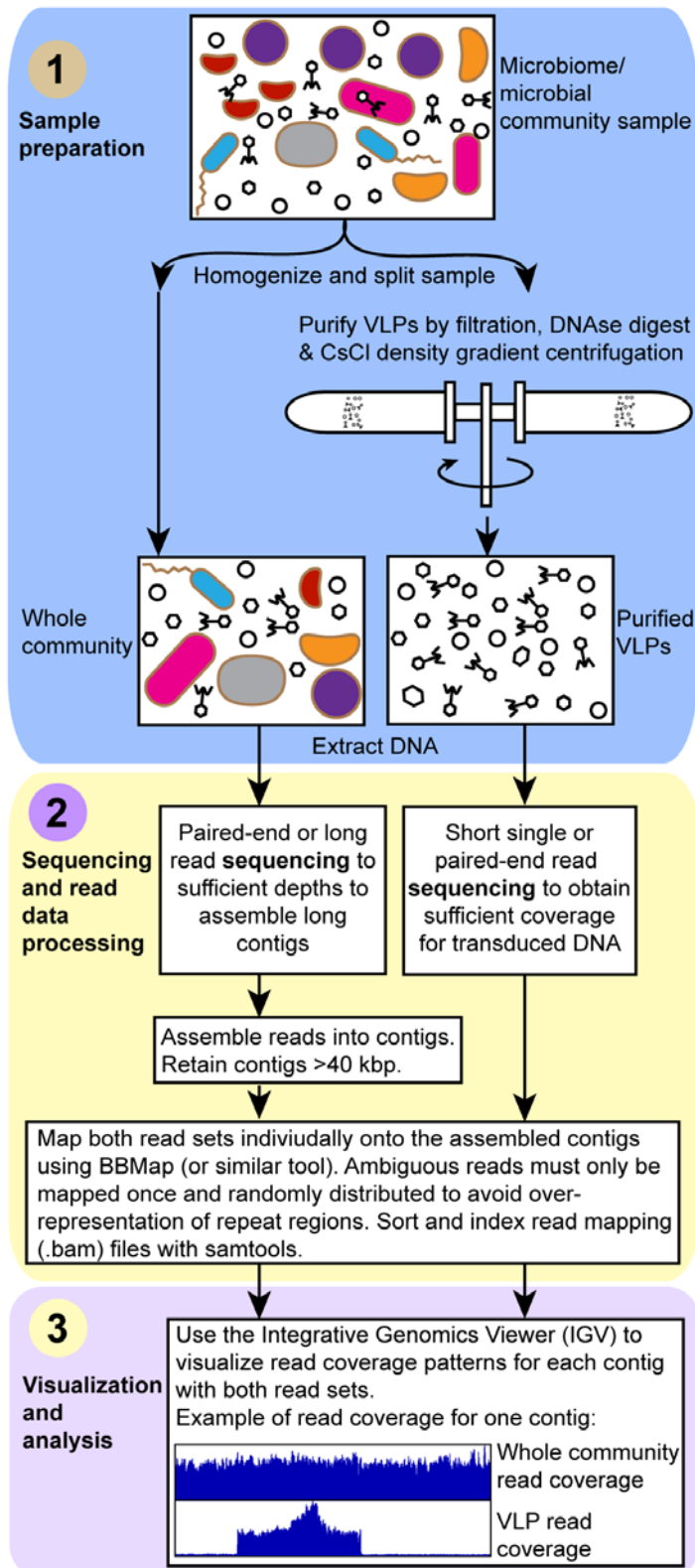
75 pathogenicity islands, can hijack phages capsids in an act of molecular piracy that enables their
76 transduction(14, 15).

77 The unifying characteristic of all types of transduction is that virus or VLPs serve as the vector for transfer
78 of genetic material between cells. Evidence so far indicates that these particles are abundant in most
79 environments and that transduction occurs with a high frequency(16, 17). However, approaches for
80 measuring the abundance of transducing particles and transduction frequencies in microbiome samples are
81 limited. These approaches usually rely on the application of cultured phage to environmental samples(16,
82 17) or sequencing of bacterial 16S rRNA genes from purified VLPs(18). The latter approach can
83 determine which bacterial taxa's DNA is carried in a VLP. However, the approach is limited to marker
84 genes for which conserved PCR primer pairs exist and thus the majority of transduced DNA cannot be
85 detected.

86 Here we describe an unbiased approach, termed “transductomics”, which uses DNA sequencing to
87 identify and characterize DNA originating from microbial cells that is carried in VLPs. This DNA is thus
88 part of the pool of potentially transduced DNA termed the “transductome”. Our approach is based on two
89 observations of transduced DNA in VLPs. First, transduced DNA often represents the genome of hosts
90 that are present in the same sample as the VLPs. Therefore, if DNA from a microbe is found within VLPs
91 purified from the same sample this indicates a potential transduction event. Second, unique regions of the
92 microbial host's genome are unevenly enriched in the VLPs, as most mechanisms of transduction do not
93 lead to random packaging of the host's genome. In recent years, the uneven sequence coverage patterns
94 produced by phages or GTAs carrying microbial host DNA have been used to characterize the genome
95 biology and mechanisms of DNA packaging of specific host-phage/GTA systems(19–22). Our
96 transductomics approach exploits these sequence coverage patterns to identify and characterize transduced
97 DNA in microbial communities. In the past, host DNA carried by VLPs may have been sequenced during
98 metagenomic sequencing of purified VLPs, but without appropriate analysis tools these host derived
99 sequences were classified as host contamination of the VLP sample rather than being recognized as
100 transduced DNA(23).

101 The transductomics approach that we present requires the sequencing of both the complete microbial
102 community sample, and VLPs that are ultra-purified using CsCl density gradient centrifugation from the
103 same sample (Fig. 1). The VLP and complete sample sequencing reads are mapped to long genome
104 contigs assembled from the complete sample metagenome. These contigs represent both microbial and
105 viral genomes. Visualization of the read mapping coverages along the contigs comparing VLP and
106 complete metagenome read coverages reveals patterns that can be associated with host DNA transport via
107 VLPs. We demonstrate this method first using pure culture models of different transducing phages and

108 other transducing particles. This is followed by the application of the approach to a murine intestinal
109 microbiome community.



110

111 **Figure 1: The “transductomics” workflow.** In the sample preparation step the sample is gently
112 homogenized and split into two subsamples. One subsample is directly used for whole community DNA
113 extraction, the other subsample is subjected to ultra-purification of virus-like particles (VLPs) using a
114 combination of filtration, DNase digest and CsCl density gradient centrifugation as previously
115 described(24) followed by DNA extraction from the purified VLPs. Both DNA samples are sequenced to
116 different depths and potentially with different sequencing approaches, although in many cases the same
117 sequencing approach could be applied to both samples. For the whole community DNA sample, the
118 sequencing must focus on ultimately achieving assembly of long metagenomics contigs. For the VLP
119 DNA sample, the sequencing must focus on maximal read coverage, and no assembly is needed for these
120 reads. The whole community sequencing reads are assembled using a suitable assembler. Contigs smaller
121 than 40 kbp are discarded. Both the whole community and VLP sequencing reads are mapped onto the
122 contigs >40 kbp using BBMap(25) ensuring that ambiguously mapped reads are only used once and
123 randomly assigned. To find transduced regions, the contigs read coverage patterns for both whole
124 community and VLP reads are visualized using the Integrative Genomics Viewer(26).

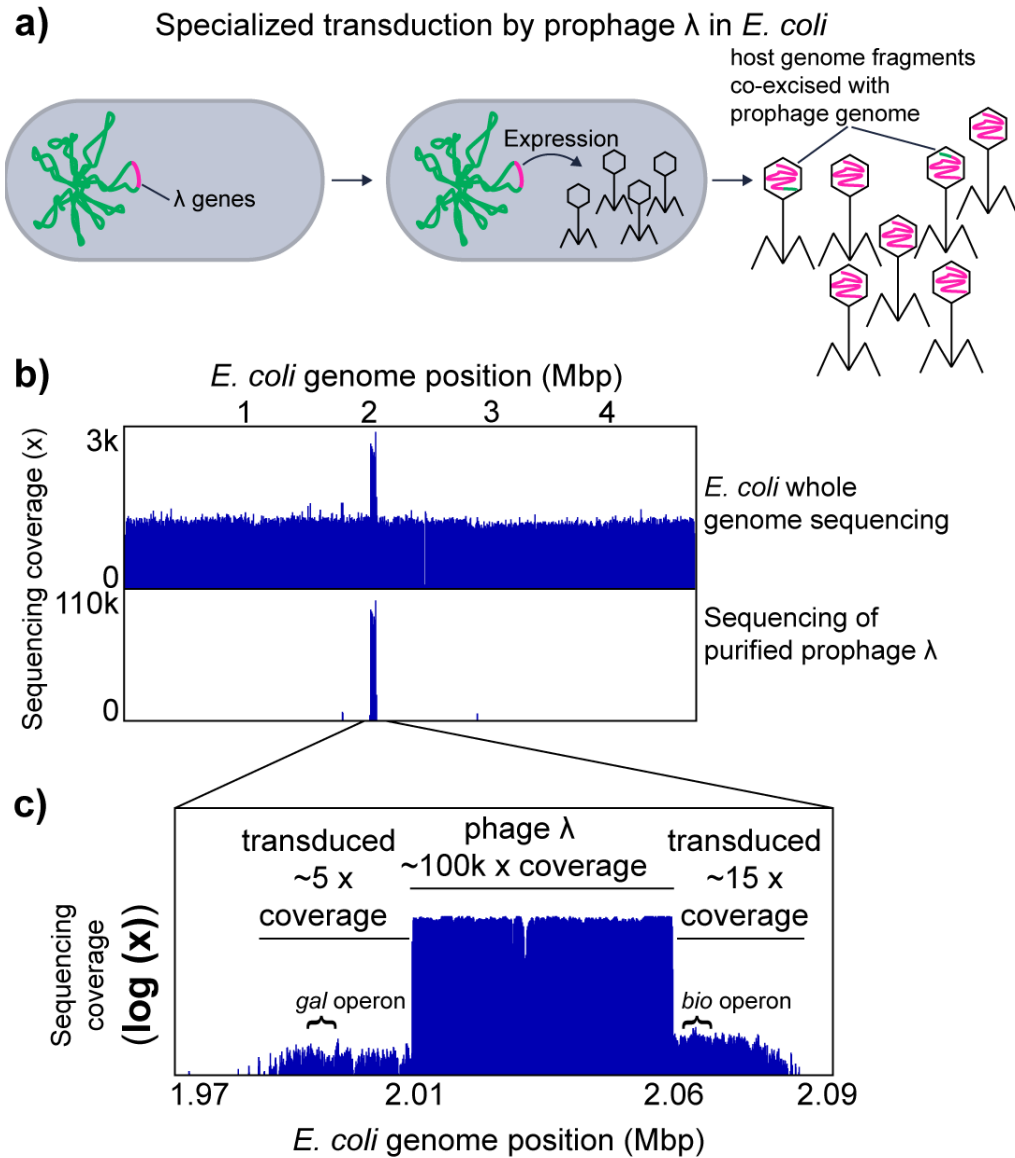
125 **Results and Discussion**

126 **Characterization of sequence coverage patterns associated with different transduction** 127 **modes in model systems**

128 **Specialized transduction by *Escherichia coli* prophage λ (27):** We used the well-studied
129 specialized transducing bacteriophage λ to analyze the sequencing coverage patterns produced by
130 specialized transduction. In specialized transduction a prophage, which is integrated in the chromosome of
131 the bacterial host, packages host genome derived DNA with low frequency due to imprecise excision from
132 the genome upon prophage induction. Prophage λ integrates between the *gal* (galactose metabolism) and
133 *bio* (biotin metabolism) operons in the *E. coli* genome. In rare cases λ excision is imprecise and either the
134 *gal* or the *bio* operon is excised and packaged in the phage particle (Fig. 2a)(27). This packaging of
135 adjacent *E. coli* host derived DNA can lead to the transduction of the *bio* and the *gal* operons.
136 Transduction of recipient cells can be temporary or permanent, depending on if the DNA gets recombined
137 into the chromosome or remains as an extrachromosomal element, which is diluted out in the population
138 during cell divisions.

139 Using the transducomics approach we found that coverage of the *E. coli* genome with sequencing reads
140 derived from purified λ phage particles is almost exclusively restricted to the λ phage integration site and
141 two ~25 kbp regions on the left and right of the λ integration site (Fig. 2b and c). These flanking regions
142 with read coverage represent the regions that are transduced by λ phage as indicated by the presence of the
143 *bio* and the *gal* operons in these flanking regions (Fig. 2c). The coverage of the λ prophage region is
144 roughly 10,000 fold greater than the coverage of the flanking transduced regions indicating that only a
145 small number of phage particles actually carry transduced DNA and thus are specialized transducing
146 particles.

147 Using the *E. coli*-prophage λ model we show that specialized transduction by a prophage produces a
148 unique read coverage pattern. Furthermore, analysis of the read coverage pattern of the transduced DNA
149 region adjacent to the prophage DNA allows determination of both the size and content of the transduced
150 host genome region (~50 kbp in total in case of λ), as well as estimation of the frequency with which
151 transducing particles are produced (1:10,000 in case of λ). The number of transducing particles produced
152 based on our data is roughly 100-fold higher than previously reported values for successful transduction of
153 the *gal* operon by phage λ (1:1,000,000 successful transductions per λ particles)(28), which indicates that
154 only a small fraction of λ carrying host DNA ultimately leads to successful transduction.



155

156 **Figure 2:** Specialized transduction by *E. coli* prophage λ . a) Illustration of specialized transduction. The
 157 prophage λ genome is integrated into the host chromosome. Upon induction of the prophage, the prophage
 158 genome is excised and replicated. The phage structural genes are expressed, phage particles are produced
 159 and the replicated phage genome is packaged into phage heads. Ultimately the phages are released to the
 160 environment by lysis of the host cell. Imprecise excision of the prophage λ genome happens at low
 161 frequency and leads packaging of the host chromosome into phage heads. These parts of the host
 162 chromosome can be transferred to new host genomes in the process called specialized transduction. b)
 163 Genome coverage pattern associated with prophage λ induction and specialized transducing prophage λ .
 164 The upper box shows coverage patterns for whole genome sequencing reads and purified phage particle
 165 reads mapped to the *E. coli* genome. c) In the lower box, an enlargement of the purified phage read
 166 coverage for the prophage λ region is shown (log scale). The positions of the *gal* and *bio* operons, which
 167 are known to be transduced by prophage λ , are indicated(27).

168 **Generalized transduction of the *Salmonella enterica* serovar typhimurium LT2 genome by**
169 **phage P22 and the *E. coli* genome by phage P1:**

170 We used two well-studied generalized transducing bacteriophages P22 and P1 to analyze the sequencing
171 coverage patterns produced during generalized transducing events. In generalized transduction nonspecific
172 host chromosomal DNA is packaged into phage particles during lytic infection and can then be injected
173 into a new host cell (Fig. 3a). The DNA can then recombine into the host chromosome by homologous
174 recombination.

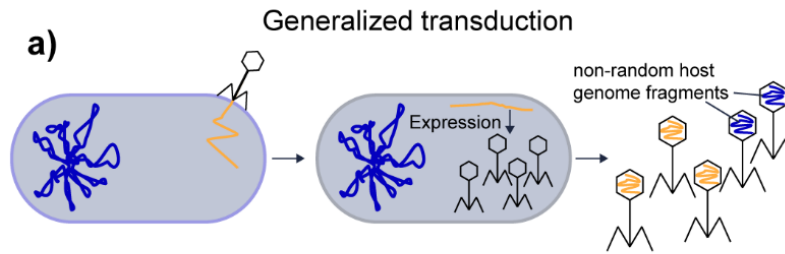
175 98.2% of the sequencing reads from purified P22 particles mapped to the P22 genome leaving 1.8% of
176 reads that map to the *S. enterica* genome. The percentage of P22 particles mapping to the *S. enterica*
177 genome corresponds to the reported percentage of 1.5% transducing P22 particles (i.e. carry host DNA)
178 previously reported(29). The mapped P22 derived reads covered the *S. enterica* genome unevenly, while
179 whole genome sequencing of *S. enterica* yielded even coverage (Fig. 3b). Regions of high or low P22 read
180 coverage corresponded in 23 out of 28 previously reported transduced chromosomal markers(30) (Fig.
181 3b). Only one region at around 4 Mbp, for which high transduction frequencies had been reported, did not
182 show high coverage (Fig. 3b), which might be due to differences in *pac* sites within this region between
183 the *S. enterica* strain used in our study and the strain used in 1982.

184 The coverage of P22 derived reads showed a distinct pattern of peaks that rise vertically on one side and
185 decline slowly over several 100 kbp increments on the other side. We speculate that the vertical edge of
186 the peak corresponds to the location of the *pac* site at which the packaging of DNA into phage heads is
187 initiated and that the slope of the peak indicates the range of processivity of the headful packaging
188 mechanism (i.e. how many headfuls are packaged into particles before the packaging apparatus dissociates
189 from the chromosome). This speculation is based on several facts: (1) the size of host DNA carried by
190 transducing particles corresponds to the size of the P22 genome (~44 kbp)(31); (2) the P22 genome is
191 replicated by rolling circle replication, which produces long concatemers of P22 DNA. A specific
192 sequence on the phage DNA (*pac* site) initiates the packaging of these concatemers into phage heads using
193 a headful mechanism(31); (3) the packaging of phage DNA continues sequentially along the P22 genome
194 concatemer with a decreasing probability for each next headful to be encapsulated in a phage particle(30);
195 (4) there are five to six sequences on the *S. enterica* genome that are similar to the *pac* site, which leads to
196 packaging of *Salmonella* DNA into P22 particles upon P22 infection, albeit with much lower frequency as
197 compared to P22 DNA(30).

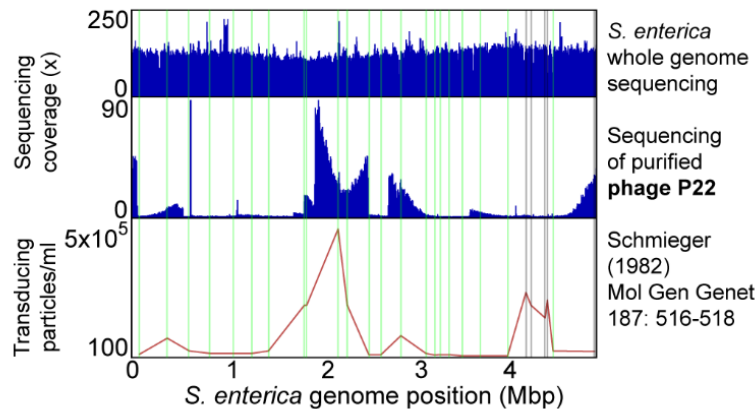
198 For *E. coli* phage P1, the majority of sequencing reads from purified P1 particles mapped to the P1
199 genome and only 4.5% of the reads mapped to the *E. coli* genome. The percentage of transducing P1
200 phages was previously reported to be 6%(32). We also observed that the P1 derived reads mapping to the

201 *E. coli* genome covered the genome unevenly. However, the pattern was less pronounced as compared to
202 P22 and *S. enterica* (Fig. 3c). This low unevenness in sequencing read coverage corresponds to previous
203 data on transduction frequencies of chromosomal markers, which found a maximum transduction
204 frequency across the *E. coli* genome of 10 fold(33).

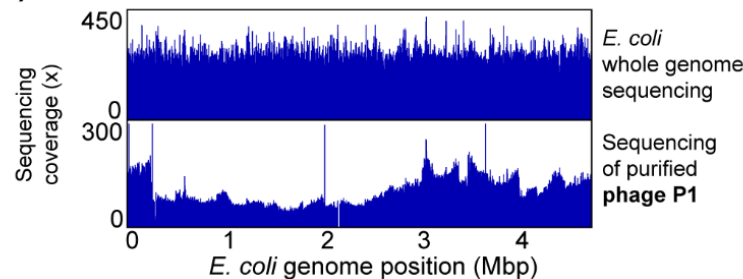
205 Sequencing host DNA carried in generalized transducing phages reveals uneven read coverage patterns
206 along the host genome indicative of transduction. These patterns vary in magnitude depending on the
207 transducing phage and they can only be observed if read coverage is analyzed along long stretches of the
208 host genome covering multiples of the length of the DNA carried by the transducing phage e.g. in case of
209 P22 44 kbp. Additionally, the patterns also provide an indication of the frequency with which different
210 regions of the host genome are transduced, as well as the locations of the *pac* sites.



b) Phage P22 in *Salmonella enterica* sv. Typhimurium LT2



c) Phage P1 in *Escherichia coli*



211

212 **Figure 3: Generalized transduction by *S. enterica* phage P22 and *E. coli* phage P1.** a) Illustration of
213 generalized transduction. Upon phage infection, the phage genome is replicated in the host cell by rolling
214 circle replication resulting in genome concatamers and phage particles are produced. The phage genome is
215 packaged into the phage head by a so called head-full packaging mechanism, which relies on the
216 recognition of a packaging (*pac*) site. The bacterial host chromosomes contain sites that resemble the *pac*
217 site and thus lead to packaging of non-random pieces of the host chromosome into phage heads. The
218 packaging happens in a processive fashion i.e. after one phage head has been filled the packaging
219 machinery continues to fill the next phage head with the remaining DNA molecule. The likelihood that the
220 packaging machinery dissociates from the molecule increases the further away from the *pac* site it gets,
221 thus leading to a decreased packaging efficiency over distance. b) *Salmonella enterica* genome coverage
222 pattern associated with generalized transduction by phage P22. Whole genome sequencing reads and
223 purified phage particle reads were mapped to the *S. enterica* genome. In the lower part transduction
224 frequencies for 28 chromosomal markers along the chromosome are shown as determined by Schmieger
225 (1982)(30). Vertical lines indicate the positions of the chromosomal markers in green where the
226 transduction frequency matches the read coverage, in grey where read coverage does not correspond to
227 reported transduction frequency. c) *Escherichia coli* genome coverage pattern associated with generalized
228 transduction by *E. coli* phage P1.

229 **Hijacking of helper prophage by a phage-related chromosomal island in *Enterococcus*** 230 ***faecalis* and specialized transduction by prophages**

231 Certain chromosomal islands, including pathogenicity islands and integrative plasmids, are mobilized
232 using helper phages(14, 34). This is a form of molecular piracy in which structural proteins of the helper
233 phage are hijacked by the chromosomal island and used as a vehicle for the transfer of the island to other
234 cells. We used *E. faecalis* VE14089, which is a natural resident of the human intestine and causes
235 opportunistic infections, to study the sequencing coverage patterns produced by chromosomal island
236 transfer by way of a helper phage. *E. faecalis* VE14089 is host to a chromosomal island (EfCIV583) that
237 uses structural proteins from a helper phage (pp1) for transfer(15, 34) (Fig. 4a). *E. faecalis* VE14089
238 possesses five additional prophage-like elements (pp2 to pp6). Some of these prophages contribute to *E.*
239 *faecalis* pathogenicity and confer an advantage during competition with other *E. faecalis* strains in the
240 intestine(15, 35).

241 Read coverage differs widely between the different prophage like elements with the coverage of
242 EfCIV583 exceeding the coverage of all other elements by almost an order of magnitude (Fig. 4b). This
243 finding is in line with previous results that showed that EfCIV583 DNA is more abundant than all other
244 prophage DNA in purified VLPs from *E. faecalis* V583(35), an isogenic strain of VE14089. Interestingly
245 pp2, pp4 and pp6 did not yield coverage peaks in the VLP fraction, which confirms previous observations
246 that these prophage elements are not excised under the conditions that we used for prophage
247 induction(15).

248 For pp1, pp5 and EfCIV583 we see patterns (coverage slopes visible in the log-scale coverage plot) that
249 indicate that not only the chromosomal island is transduced but also regions adjacent to these prophages

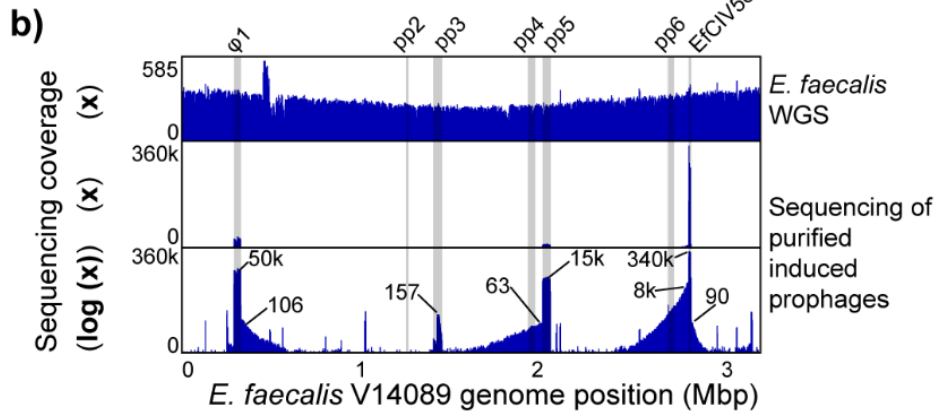
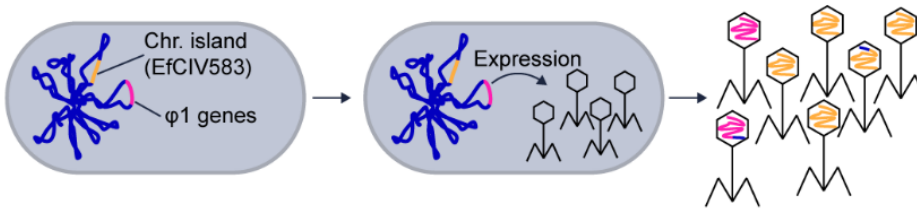
250 and the chromosomal island. Based on the maximal coverage of the transduced regions versus the
251 coverage of the prophage regions (Fig. 4b) we estimate the maximal frequencies of transduction to be
252 1:500 for pp1, 1:240 for pp5, and 1:43 for the left side of EfCIV583 and 1:3780 for the right side of
253 EfCIV583. These relatively high transduction frequencies and the fact that transduced regions span
254 several hundred kbp facing unidirectional from the integration site of the prophages and EfCIV583
255 suggest a lateral transduction mechanism as described by Chen et al.(12).

256 Our data also revealed that there are several additional regions in the *E. faecalis* VE14089 genome that
257 had an elevated sequencing coverage in the purified phage sample suggesting that these regions encode
258 additional elements that are transported in VLPs. These elements consist of IS-Elements that carry a
259 transposase and surprisingly the three rRNA operons. For the rRNA operons the coverage has a deep
260 valley between the 16S and the 23S rRNA gene suggesting that a specific mechanism for rRNA gene
261 transport is present or that the processed rRNAs were sequenced. We can currently think of three
262 explanations for this intriguing pattern. First, potentially ribosomes are enriched alongside the VLPs in our
263 VLP purification method. However, if this were the case we should have observed similar patterns in VLP
264 fractions of other pure culture organisms, which we did not. Second, intact ribosomes are packaged by
265 VLPs produced in *E. faecalis*. However, this leaves open the question of why the rRNA from these
266 ribosomes was amenable to sequencing by the Illumina method used, which should not enable direct
267 sequencing of RNA. Third, DNA with rRNA genes are packaged with high specificity into one or several
268 types of VLPs from *E. faecalis*.

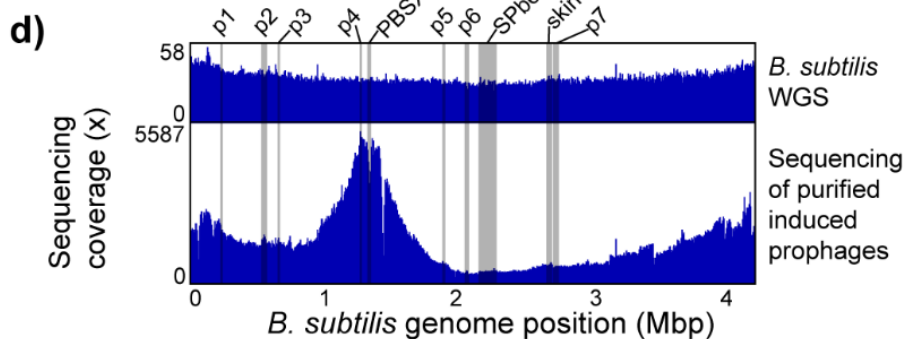
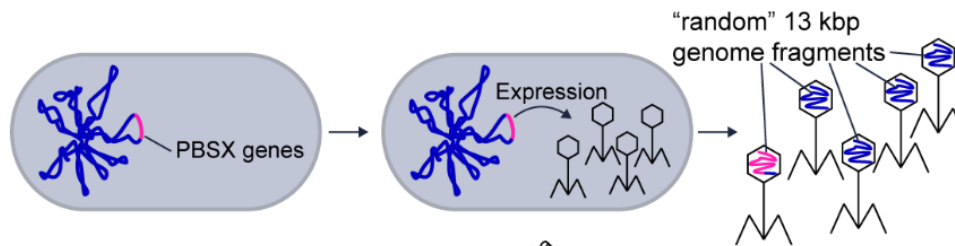
269 Our results show that the transduction of a chromosomal island by a prophage can produce a similar
270 coverage pattern as an induced prophage, indicating that chromosomal islands are easy to detect based on
271 read mapping coverage. However, they can only be distinguished from prophage by annotation of the
272 genes and genomic regions.

273

a) Specialized transduction and transduction of a chromosomal island by prophages in *E. faecalis* V14089



c) Gene transfer agent-like packaging of the *B. subtilis* chromosome by the defective prophage PBSX



274

275 **Figure 4: Other types of transduction.** a) Specialized transduction (see description for prophage λ) and
 276 transduction of a chromosomal island by prophages in *E. faecalis* V14089. The chromosome of *E. faecalis*
 277 contains multiple prophages including φ1 and the chromosomal island EfCIV583. Upon induction φ1 and
 278 EfCIV583 are excised from the chromosome and replicated. EfCIV583 hijacks the structural proteins of
 279 φ1 when they are produced and large number of phage particles that carry the EfCIV583 genome are

280 produced. b) *E. faecalis* V14089 genome coverage patterns associated with prophage induction and
281 EfCIV583 transduction. Whole genome sequencing (WGS) reads and purified VLPs were mapped to the
282 *E. faecalis* genome. The lowest part of the box shows VLP read coverage on a log scale. The small
283 numbers in this plot give x fold coverage for specific genome positions corresponding to prophages or the
284 chromosomal island EfCIV583 and the surrounding areas that are likely transduced. The positions of
285 known prophage-like elements and EfCIV583 in the *E. faecalis* genome are highlighted by grey bars. c)
286 Gene transfer agent-like packaging of the *B. subtilis* chromosome by the defective prophage PBSX. The
287 *B. subtilis* chromosome contains a variety of prophages and prophage-like elements including the
288 defective prophage PBSX(36). Upon expression of the PBSX genes phage-like particles are produced,
289 which contain random 13 kbp pieces of the host chromosome(37). d) *B. subtilis* genome coverage patterns
290 associated with prophage induction. Whole genome sequencing (WGS) reads and purified prophage
291 particle reads were mapped to the *B. subtilis* genome. The positions of known prophages and prophage-
292 like elements in the *B. subtilis* genome (36) are highlighted by grey bars.

293

294 **Transport of *Bacillus subtilis* genome by gene transfer agent-like element PBSX**

295 We used induced the gene transfer agent (GTA)-like element PBSX from the *B. subtilis* ATCC 6051
296 genome to study the sequencing coverage pattern produced by the supposedly randomized incorporation
297 of fragments from the whole genome into GTA type VLPs. PBSX is a defective prophage that randomly
298 packages 13 Kbp DNA fragments of the *B. subtilis* genome in a GTA-like fashion (Fig. 4c)(10, 37, 38). In
299 contrast to other GTAs it does not transfer the packaged DNA between cells but rather acts similar to a
300 bacteriocin against *B. subtilis* cells that do not carry the PBSX gene cluster(10).

301 DNA sequencing reads derived from purified PBSX particles covered the *B. subtilis* genome unevenly
302 with a maximum 30 fold difference between the lowest and highest covered regions (Fig. 4d). Reads from
303 whole genome sequencing of *B. subtilis* covered the genome evenly slightly increasing toward the origin
304 of replication, as expected(39). The genomic region containing PBSX had a lower read coverage in VLP
305 particle derived reads as compared to neighboring genomic regions. This is consistent with results from a
306 previous study where it was found that a genetic marker integrated in the PBSX region was less frequently
307 packaged into particles as compared to a marker in a neighboring region(40). Interestingly, the genomic
308 region containing the prophage SPbeta, which gets excised upon mitomycin C treatment(41), did not
309 show any higher or lower coverage in the VLP particle derived sequencing reads as compared to
310 neighboring genomic regions (Fig. 4d).

311 Our results show that packaging of host DNA by the GTA-like PBSX element of *B. subtilis* produces a
312 distinct and non-random sequencing coverage pattern that bears similarities to the read coverage pattern
313 produced by the generalized transducing phage P1 (Fig. 3c).

314 **Detection limits of the approach**

315 The patterns for different transduction modes have distinct characteristics that will impact sensitivity of
316 detection and the false positive rate. For prophage induction and specialized transduction pattern detection
317 there are three potential challenges: (1) the length of the genome sequence fragment (contig) used for read
318 mapping needs to be sufficiently long to encompass both the prophage genome, as well as a portion of the
319 host genome; (2) potential assembly artifacts (chimeric contigs consisting of multiple source genomes)
320 can lead to highly uneven read coverage that could look similar to an induced prophage pattern. In the
321 case of our approach this is mitigated by the fact that we map whole metagenome and VLP reads to the
322 same contigs and thus we expect to obtain even read coverage for the whole metagenome read mapping,
323 which is indicative of correct assembly; and (3) if read coverage is too low patterns will not be sufficiently
324 distinct. It can be expected that frequency of specialized transduction is specific to specific host species
325 and prophages. Nevertheless, we tested the lower limit of read coverage levels needed for detection of
326 prophage induction and specialized/lateral transduction by down sampling read numbers for the VLP
327 reads from *E. coli* prophage λ and *E. faecalis* prophages (Figs. 2 and 4b) to achieve coverage levels
328 similar to what we observed for our mouse case study below, which ranged from several tenfold to several
329 thousand fold. For *E. coli* prophage λ we found that at ~6000x maximum read coverage (5% of total
330 reads) the specialized transduction pattern was still weakly visible, but disappeared at lower coverages,
331 while the induction of the prophage itself was still identifiable at read coverages of 20x (0.01% of total
332 reads) and less (Fig. S1a). For the *E. faecalis* prophages specialized/lateral transduction patterns were still
333 visible at ~500x coverage for the pp1 region and at ~150x for the pp5 region. Prophage induction was
334 detectable at coverages well below 40x (Fig. S1b). These results indicate that specialized and lateral
335 transduction, as well as prophage induction, can sufficiently be detected with read coverages obtained in
336 shotgun metagenomic sequencing of VLPs.

337 For generalized transduction and GTA mediated DNA transfer pattern detection the two main challenges
338 are; (1) potential generation of similar patterns by contamination of the ultra-purified VLPs with DNA
339 from microbial cells, which can for example be addressed by comparing contig rank abundances between
340 whole metagenome and VLP read coverage (see below in case study); and (2) difficulty to recognize the
341 pattern on short contigs, because sloping can extend across 100s of kbp. To test if generalized transduction
342 or GTA-like patterns can be detected on shorter contigs we used the P22, P1 and PBSX data to simulate
343 how contig length impacts pattern visibility. For this we looked at the coverage patterns of 200 kbp long
344 stretches in the genome (Fig. S2). We found that detectability of generalized and GTA-like patterns in 200
345 kbp sequence stretches depended on where the 200 kbp stretch was located within the overall read
346 coverage pattern. In some cases distinguishable coverage sloping was observed (e.g. #2 in Fig. S2a and #2
347 in S2c) in other cases coverage looked even or irregular (e.g. #4 in S2b and #4 in S2c). These results

348 indicate that generalized transduction and GTA mediated DNA transfer can be detected from contig
349 lengths produced using short read shotgun metagenomics of microbiome samples, however, some DNA
350 transfer events are likely missed if the longer contigs do not cover regions that show the characteristic
351 coverage sloping associated with these transfer events.

352 **Case study: High occurrence of transduction in the intestinal microbiome**

353 We next assessed the power and application of our transductions approach for detecting transduced
354 DNA in VLPs from complex microbiomes. We sequenced the whole metagenome (~390 mio reads) and
355 VLPs (~360 mio reads) from a fecal sample of one mouse to high coverage. The VLPs were ultra-purified
356 using the multi-step procedure shown in Fig. 1, for which we previously showed that it efficiently
357 removes DNA from microbial cells and the mouse present in fecal samples(24). We were able to assemble
358 2143 contigs >40 kbp from the whole metagenome reads with the largest contig being 813 kbp (ENA
359 accession for assembly: ERZ1273841). We discarded contigs <40 kbp because detection of transduction
360 patterns requires coverage analysis of a sufficiently large genomic region. We mapped the metagenomic
361 and VLP reads to the contigs >40 kbp to obtain the read coverage patterns. For complete metagenome,
362 44% of all reads mapped to the contigs >40 kbp and for the VLPs 10% of all reads indicating that a large
363 portion of DNA carried in VLPs is derived from prophages and microbial hosts. Of the 2143 contigs, 1957
364 showed a “standard” read coverage pattern (Fig. 5a, Suppl. Table S1), i.e. high even coverage of the
365 contigs with metagenomic reads and low even or no coverage with VLP reads, indicating no mobilization
366 of host DNA in VLPs. The remaining 186 contigs (8.6% of all contigs >40 kbp) showed a read coverage
367 pattern that indicates potential mobilization of DNA in VLPs (Fig. 5b-f, Suppl. Table S1).

368 To verify that the multi-step VLP ultra-purification procedure employed for this study efficiently removes
369 DNA from microbial cells, ruling out potential microbial host DNA contamination, we further assessed
370 read coverage patterns for the 186 contigs in comparison to all contigs. We ranked all 2143 contigs by
371 their normalized coverage for both the whole metagenome and the purified VLP samples (i.e. average x
372 fold read coverage / sum of average x fold read coverage for the sample) with the highest normalized
373 coverage being assigned rank 1 (Table S1). The expectation is that contigs from which DNA is carried in
374 VLPs have the same or lower rank for the VLP sample as compared to the whole metagenome sample,
375 while the rank for contigs for which VLP reads are derived from microbial contamination should have a
376 higher rank as compared to the whole metagenome reads, because randomly contaminating DNA would
377 be depleted in the purified VLP sample. We found that 26 out of the 186 contigs with coverage patterns
378 suggesting DNA mobilization had a normalized coverage based rank that was higher for VLP read
379 coverage than for whole metagenome read coverage indicating that these 26 patterns are potentially due to
380 contamination or alternatively due to very low efficiency of mobilization.

381 We classified all contigs taxonomically using CAT(42) (Suppl. Tables S2 and S3). The majority of contigs
382 were classified as Bacteroidetes (all contigs: 805, transduction pattern contigs: 83), Firmicutes (all: 586,
383 transduction pattern: 42), Proteobacteria (all: 89, transduction pattern: 3), or not classified at the phylum
384 level (all: 527, transduction pattern: 34). We found that with a few exceptions the relative abundance of
385 contigs assigned to specific phyla was similar between the set of all contigs >40 kbp and the subset of
386 contigs with transduction patterns. The phyla that differed in relative contig abundance were
387 Proteobacteria with less than half the relative abundance in the contigs with transduction patterns,
388 Verrucomicrobia with 3.5x and *Candidatus* Saccharibacteria with 11.5x the contig abundance in the
389 contigs with transduction patterns. Since members of *Cand.* Saccharibacteria have been shown to be
390 extremely small (200 to 300 nm)(43) it is likely that they share similar properties with bacteriophages in
391 terms of size and density and thus might get enriched in the VLP fraction. In fact, all transduction patterns
392 of *Cand.* Saccharibacteria contigs were classified as “unknown” or “unknown, potentially a small
393 bacterium” prior to knowing the taxonomic identity of the contigs.

394 We classified the type of DNA mobilization/transduction in the 186 contigs with a mobilization pattern
395 based on the visual characteristics of the mobilized region in the VLP read coverage, as well as based on
396 annotated genes within the mobilized region. For example, we classified mobilization patterns as prophage
397 if the characteristic pattern showed high coverage with sharp edges on both sides (compare Fig. 2) and the
398 presence of characteristic phage genes (e.g. capsid proteins) as an additional but not required criterion.

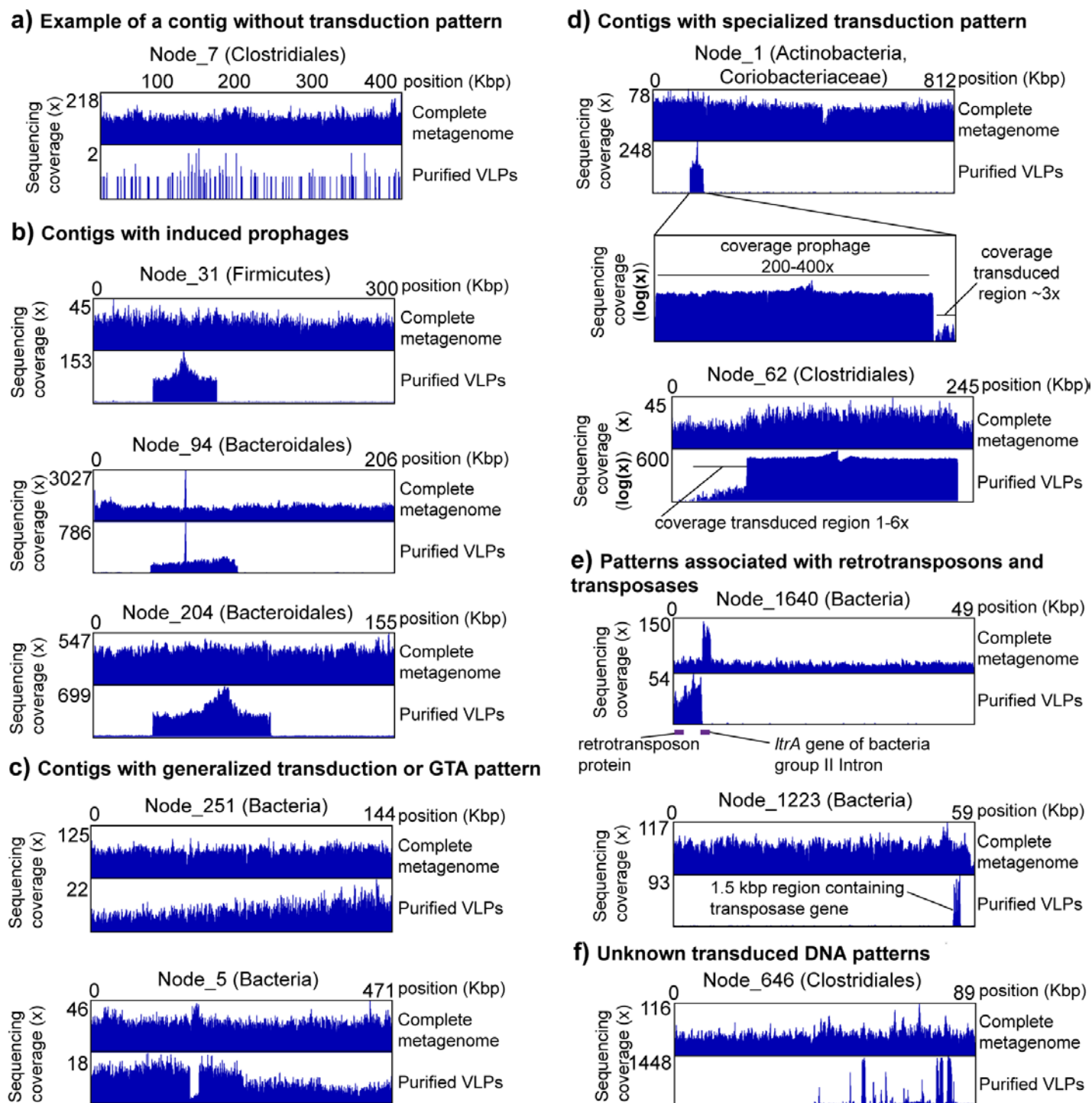
399 We observed 74 contigs that indicated induced prophages. Of these, 12 (16%) prophages showed
400 indications of specialized transduction i.e. read coverage above the base level of the contig in regions
401 adjacent to the prophage (Fig. 5b and d). Additionally, we classified 8 patterns as potential prophages or
402 chromosomal islands, as they showed the same pattern as other prophages, but we were unable to find
403 recognizable phage genes in the annotations.

404 We found patterns of potential generalized transduction or GTA carried DNA in 46 contigs, however,
405 some patterns were observed for shorter contigs and could thus potentially be incorrect classifications
406 (Fig. 5c). One of the contigs (NODE_5, classified as Bacteria) with a generalized transduction or GTA
407 pattern additionally showed a sharp coverage drop in a ~15 kbp region only in the VLP reads (Fig. 5c).
408 This region is flanked by a tRNA gene and carries one gene annotated as a potential virulence factor,
409 internalin used by *Listeria monocytogenes* for host cell entry(44). This region might represent a
410 chromosomal island that was excised from the bacterial chromosome prior to or during production of the
411 unknown VLP and that did not get encapsulated in the VLP. Alternatively, similar strains might be present
412 in the sample, but only some carry the chromosomal island and strains carrying the chromosomal island

413 are less prone to producing the VLPs, e.g., by superinfection resistance provided by the chromosomal
414 island against a generalized transducing phage.

415 We observed 9 patterns that showed strong differences between whole metagenome read coverage and
416 VLP read coverage, but that did not correspond to any of the patterns we analyzed in our proof-of-
417 principle work. However, based on gene annotations we determined that these patterns likely represent
418 retrotransposons or other transposable elements. For example, on contig NODE_1640 (classified as
419 Bacteria by CAT) we observed high coverage with VLP reads on one part of the contig, which carries a
420 gene annotated as a retrotransposon (Fig. 5e). Interestingly, the retrotransposon region is flanked by a *ltrA*
421 gene which is encoded on bacterial group II intron and encodes maturase, an enzyme with reverse
422 transcriptase and endonuclease activity(45). Surprisingly the region containing the *ltrA* gene had above
423 average coverage in the whole metagenome reads, but no coverage in the VLP reads. This suggests that
424 the intron actively reverse splices into expressed RNA with subsequent formation of cDNA(45) leading to
425 increased copy number of this genomic region. As another example, on contig NODE_1223 (classified as
426 Bacteria by CAT) a region containing a transposase gene is strongly overrepresented in the VLP reads
427 suggesting that this region is a transposable element that is packaged into a VLP (Fig. 5e).

428 Finally, we determined that two patterns are likely lytic phages and 47 patterns are classified as
429 “unknown” transduced DNA, as the coverage pattern is uneven indicating transport in VLPs but we could
430 not determine the type of transport. To provide an example, in contig NODE_646 (classified as
431 Clostridiales by CAT) we observed a potential prophage pattern in which we found some of the main
432 phage relevant genes such as major capsid protein, however, within the prophage pattern we observed
433 high coverage spikes for which we currently have no good explanation (Fig. 5f).



434
435 **Figure 5: Example of transduction patterns in the mouse intestinal microbiome.** The taxonomic
436 classification for each contig specified in parentheses after each contig name is the lowest taxonomic level
437 successfully classified by CAT(42)(Suppl. Table S2). The complete metagenome reads and the purified
438 VLP reads were mapped to the same exact set of contigs assembled from the complete metagenome reads.
439 The read coverage pattern of the complete metagenome reads provides evidence for the correct assembly
440 of the contigs and allows to distinguish potential transduction derived VLP read coverage patterns from
441 VLP coverage patterns due to contamination with microbial DNA. With the exception of panel A) all

442 shown contigs have the same or a lower abundance rank for VLP read coverage as compared to complete
443 metagenome read coverage indicating that their overall read coverage was enriched in the VLP samples.
444 Read coverage due to VLP sample contamination with cellular DNA is expected to result in a higher
445 abundance rank for VLP read coverage, as compared to complete metagenome read coverage.

446 **Conclusions and Outlook**

447 The transductomics approach that we developed should be applicable to a broad range of environments
448 ranging from host-associated microbiomes to soils and aqueous environments. For some environments
449 such as open ocean water samples the approach would need only minor modifications. For example,
450 concentration of VLPs by tangential flow filtration prior to density gradient centrifugation using well
451 established protocols(23). Thus this approach will allow addressing key questions about microbial
452 evolution via HGT in a diversity of microbial communities, including what kind of genes and with what
453 frequency are carried by VLPs. Among these questions, one of the most pressing ones is what the role of
454 transducing particles is in the transfer of antibiotic resistance genes, which is a topic of current debate(46–
455 49). Apart from its application in studying transduction in microbial communities this approach can also
456 be used to increase our understanding of the molecular mechanisms of different transduction mechanisms
457 by careful analysis of read coverage data from pure cultures that shows exact transduction frequencies of
458 each genomic location without tedious analysis of multiple genomic markers. Mechanisms that can be
459 analyzed include, for example, the identification of *pac* sites in generalized transducers, the size range of
460 transduced genomic loci in specialized transducers(20), and the analysis of how random DNA packaging
461 by GTAs really is(19).

462 One of the major surprises for us when analyzing the mouse intestinal transductome data was that around
463 one quarter of the transduction patterns that we identified are unknown. These patterns showed even
464 coverage in the whole metagenome reads and strong uneven coverage in the VLP reads (e.g. Fig. 4e),
465 however, we were unable to associate them clearly with any of transduction modes that we have
466 investigated with pure cultures. We foresee two types of future studies to characterize the nature of the
467 transducing particles that lead to these unknown patterns and to exclude that they are some kind of
468 artifact. First, read coverage patterns of newly discovered modes of transduction have to be analyzed with
469 the “transductomics” approach to correlate the patterns to patterns observed in microbiomes and microbial
470 communities. While we investigated the transduction patterns associated with both major known
471 transduction pathways, as well as more recently discovered transduction pathways, novel modes of
472 transduction are continuously discovered. These novel transduction modes that need to be characterized
473 with our approach include new types of GTAs(10), lateral transduction(12) and DNA transfer in outer
474 membrane vesicles(50, 51). Second, approaches that allow linking specific transduced DNA sequences to
475 the identity of transducing particles in microbial community samples can be developed. We envision, for
476 example, that high resolution filtration and density gradient based separation of individual VLPs will
477 allow linking the transduced DNA (by sequencing) to the identity of the transducing VLPs using
478 proteomics to identify VLP proteins. Using and developing these approaches further will allow us to
479 increase the range of transduction modes that can be detected in microbial communities, as well as

480 potentially reveal currently unknown types of transduction that are not known from pure culture studies
481 yet.

482 We see several pathways for improving the sensitivity, accuracy and throughput of the transductomics
483 approach in the future. Currently, our ability to detect generalized transduction patterns is limited by the
484 fact that detection of these patterns requires long stretches of the microbial host genome to be assembled.
485 Our P22 and P1 data shows these patterns stretch across genomic regions >500 kbp. Additionally, high
486 sequencing coverage is needed for the detection of these patterns. Assembly of long contigs in
487 metagenomes of high diversity communities is currently hampered by the relatively short read lengths of
488 sequencers that allow for high coverage. We expect, that increasing read numbers of long-read sequencing
489 technologies such as PacBio and Oxford Nanopore in the future will allow us to sequence complex
490 microbiomes to sufficient depth for the assembly of long metagenomic contigs. A combination of long-
491 read sequencing for the whole community metagenomes in combination with a short-read, high-coverage
492 approach for the VLP fraction will in the future provide more sensitive and accurate detection of
493 generalized transduction patterns. In addition to improvements in the realm of long-read sequencing we
494 expect the development of computational tools for the automatic or semi-automatic detection of
495 transduction patterns in read coverage data from paired whole metagenome and VLP metagenome
496 sequencing. There is a large number of possible parameters that could be used to train a machine learning
497 algorithm to detect transduction patterns. These parameters include differences between average read
498 coverage and maximal read coverage for VLP reads (Table S1) and the comparison of contig rank
499 abundance based on coverage, which we used to cross check transduction patterns for signatures of
500 microbial DNA contamination. Such computational tools will enable the high-throughput detection of
501 transduction patterns in many samples, which is currently limited by the need for visual inspection of
502 patterns.

503 **Online Methods**

504 ***In vitro* bacteriophage propagation and induction of transducing prophages and other** 505 **elements**

506 ***Lambda***. *E. coli* KL740 was inoculated into 300 ml of LB and grown to an OD₆₀₀ of 0.7 at 28°C with
507 aeration. The culture flask was transferred to a 42°C water bath for 10 minutes and then incubated at 42° C
508 for 30 min with shaking. The temperature was reduced to 28° C and cell lysis was allowed to proceed for 2
509 hrs. The remaining cells and debris were removed by centrifugation at 2750 x g for 10 minutes and the
510 phage containing culture fluid was filtered through a 0.45 µm membrane.

511 ***P22***. The data set used to analyze generalized transduction by *Salmonella* phage P22 was taken from a
512 previous study assessing methods for phage particle purification from intestinal contents(24). For a
513 detailed description of P22 propagation and purification please refer to our previous publication.

514 ***P1***. Lyophilized phage P1 was purchased from ATCC and resuspended in 1 ml of Lennox broth (LB) at
515 room temperature (RT). 200 µl of the phage suspension was added to 10 ml of mid logarithmic phase
516 (OD₆₀₀ ~0.5) *E. coli* ATCC 25922 and incubated for 3 hrs at 37°C with shaking. The bacteria were pelleted
517 at 2750 x g for 10 min and the culture fluid was filtered through a 0.45 µm syringe filter. 100 µl of
518 stationary phase *E. coli* ATCC 25922 was distributed to 15 separate tubes each containing 200 µl of the
519 P1 culture filtrate. The bacteria/phage mixtures were immediately added to molten LB top agar (0.5%
520 agar), poured over LB agar (1.5% agar) plates and incubated overnight (O/N) at 37°C. 2.5 ml of SM-plus
521 buffer (100 mM NaCl, 50 mM Tris·HCl, 8 mM MgSO₄, 5 mM CaCl₂·6H₂O, pH 7.4) was added to the
522 surface each plate and the top agar was scraped off and pooled. The phages were eluted from the top agar
523 by rotation for 1 hour at RT. The top agar suspension was centrifuged at 2750 x g for 10 min, the
524 supernatant was collected and the top agar was washed once with ~30 ml of SM-plus and incubated at RT
525 for an additional 1 hr. Following the wash step centrifugation was repeated and the resulting supernatant
526 was collected. The phage containing supernatants were combined.

527 ***PBSX***. To induce the prophage-like element PBSX from the *B. subtilis* ATCC 6051 genome, a 100 ml
528 culture of *B. subtilis* was grown in LB at 37°C with shaking to an OD₆₀₀ of ~0.5. Mitomycin C was added
529 at a final concentration of 0.5 µg/ml and the culture was incubated at 37°C for 10 minutes. The culture was
530 centrifuged at 2750 x g for 10 min and the pellet was washed with 50 ml of fresh LB and centrifuged a
531 second time. The cell pellet was resuspended in 100 ml of LB and grown for an additional 3 hrs at 37° C
532 with shaking. The cells and debris were removed by centrifugation and the phage containing culture fluid
533 was filtered through a 0.45 µm membrane.

534 ***Enterococcal prophages.*** *E. faecalis* strain VE14089, a derivative of *E. faecalis* V583 that has been cured
535 of its three endogenous plasmids(15), was subcultured to an OD₆₀₀ of 0.025 in 1 L of pre warmed brain
536 heart infusion broth (BHI) and grown statically at 37°C to an OD₆₀₀ of 0.5. To induce excision of
537 integrated prophages, ciprofloxacin was added to the culture at a final concentration of 2 µg/ml and the
538 bacteria were grown for an additional 4 hrs at 37°C. The bacterial cells and debris were centrifuged at
539 2750 x g for 10 min and the culture fluid was filtered through a 0.45µm membrane.

540 **Purification of phage particles from culture fluid**

541 All phage containing culture fluid was treated with 10 U of DNase and 2.5 U of RNase for 1 hr at RT. 1 M
542 solid NaCl and 10 % wt/vol polyethylene glycol (PEG) 8000 was added and the phages were precipitated
543 O/N on ice at 4°C. The precipitated phages were resuspended in 2 ml of SM-plus and loaded directly onto
544 CsCl step gradients (1.35, 1.5 and 1.7 g/ml fractions) and centrifuged for 16 hrs at 83,000 x g. The phage
545 bands were extracted from the CsCl gradients using a 23-gauge needle and syringe, brought up to 4 ml with
546 SM-plus buffer and loaded onto a 10,000 Da molecular weight cutoff Amicon centrifugal filter (EMD
547 Millipore) to remove excess CsCl. The phages were washed 3 times with ~4 ml of SM-plus and then stored
548 at 4°C.

549 **Isolation of phage and host bacterial DNA from pure cultures**

550 Following CsCl purification of phages and phage-like elements, DNA was isolated by adding 0.5 % SDS,
551 20 mM EDTA (pH=8) and 50 µg/ml Proteinase K (New England Biolabs) and incubating at 56°C for 1
552 hour. Samples were cooled to RT and extracted with an equal volume of phenol:chloroform:isoamyl
553 alcohol. The samples were centrifuged at 12,000 x g for 1 min and the aqueous phase containing the DNA
554 was extracted with an equal volume of chloroform. Following centrifugation at ~16,000 x g for 1 min
555 0.3M NaOAc (pH=7) was added followed by an equal volume of 100% isopropanol to precipitate the
556 DNA. The DNA was pelleted at 12,000 x g for 30 min and washed once with 500 µl of 70% ethanol. The
557 samples were decanted and the pellets were air dried for 10 min and resuspended in 100 µl of sterile
558 water.

559 For the isolation of bacterial genomic DNA, we used the Genra Puregene Yeast/Bacteria Kit (Qiagen)
560 according to the manufacturer's instructions.

561 **Isolation and purification of bacteria and VLPs from mouse fecal pellets for metagenomic** 562 **sequencing**

563 The entire colon contents of one male C57BL6/J mouse were added to 1.2 ml of SM-plus buffer and
564 homogenized manually with the handle of a sterile disposable inoculating loop. After homogenization the
565 sample was brought up to 2 ml with SM-plus. One third of the sample volume was added to a fresh tube

566 containing 100 mM EDTA and set aside on ice. This represented the unprocessed whole metagenome
567 sample. The remaining two thirds of the sample volume were used to isolate VLPs.

568 VLPs from the homogenized feces were ultra-purified as described previously(24). Briefly, the sample was
569 centrifuged at 2500 x g for 5 min, the supernatant transferred to a clean tube and centrifuged a second time
570 at 5000 x g to pellet any residual bacteria and debris. The supernatant was transferred to a sterile 1 ml syringe
571 and filtered through a 0.45 µm syringe filter. The clarified supernatant was treated with 100 U of DNase
572 and 15 U of RNase for 1 hr at 37°C. The sample was loaded onto a CsCl step gradient (1.35, 1.5 and 1.7
573 g/ml fractions) and centrifuged for 16 hrs at 83,000 x g. The VLPs residing at the interface of the 1.35 and
574 1.5 g/ml fractions were collected (~2 ml) and the CsCl was removed by centrifugal filtration as described
575 above. The purified VLPs were disrupted by the addition of 50 µg/ml proteinase K and 0.5% sodium dodecyl
576 sulfate (SDS) at 56° C for 1 hr. The samples were cooled to room temperature and total DNA was extracted
577 by the addition of an equal volume of phenol:chloroform:isoamyl alcohol. The organic phase was separated
578 by centrifugation at 12,000 x g for 2 minutes and the aqueous phase was extracted with an equal volume of
579 chloroform. The DNA was precipitated by the addition of 0.3 M NaOAc, pH 7, and an equal volume of
580 isopropanol. The DNA pellet was washed once with ice cold 70% ethanol and resuspended in 100 µl of
581 sterile water. The DNA was further cleaned on a MinElute spin column (Qiagen) and eluted into 12 µl of
582 elution buffer (Qiagen).

583 To purify total metagenomic DNA, unclarified fecal homogenate was treated with 5 mg/ml lysozyme for
584 30 min at 37° C. The sample was transferred to 2 ml Lysing Matrix B tubes (MP Biomedical) and bead beat
585 in a Bullet Blender BBX24B (Next Advance) at top speed for 1 min followed by placing on ice for 2 min.
586 This was repeated a total of 4 times. The samples were centrifuged at 12,000 x g for 1 min and the DNA
587 from the supernatant was extracted, precipitated and purified as described above.

588 **DNA Sequencing**

589 The concentration of purified microbial DNA was determined using a Qubit 3.0 fluorometer (Thermo-
590 Fisher). Prior to library preparation total microbial DNA was sheared using a Covaris S2 sonicator with a
591 duty cycle of 10%, intensity setting of 5.0 and a duration of 2 x 60 sec at 4° C. Sequencing libraries were
592 generated using the KAPA HTP library preparation kit KR0426 – v3.13 (KAPA Biosystems) with
593 Illumina TruSeq ligation adapters. Library quality was determined using a 2100 Bioanalyzer system
594 (Agilent). Libraries were size selected and purified in the range of 300-900 bp fragments and subjected to
595 Illumina deep sequencing. For DNA sequencing of pure phage cultures and the *E. coli* KL740 genome we
596 used an Illumina NextSeq 500 desktop sequencer. Illumina HiSeq 2500 v3 Sequencing System in rapid run
597 mode was used to sequence the metagenomes and viromes from the feces of the C57BL6/J mouse. All
598 sequencing was performed in paired end mode acquiring 150 bp reads. Per fragment end the following

599 number of reads were obtained; C57BL/6 mouse feces - 76 M reads for the whole metagenome and 97 M
600 reads for the virome, 45 M reads for phage P1, 21 M reads for lambda phage and 27 M reads for the *E.*
601 *coli* KL740 genome, 29 M reads for phage PBSX and 28 M reads for the enterococcal prophages. For the
602 C57BL/6 mouse microbiome we sequenced an additional 75 bp single-end read run to increase coverage.
603 We obtained 313 M 75 bp reads for the whole metagenome and 262 M 75 bp reads for the virome. All
604 DNA sequencing was performed by the Eugene McDermott Center for Human Growth and Development
605 Next Generation Sequencing Core Facility at the University of Texas Southwestern.

606 **Assembly of mouse fecal metagenome**

607 Read decontamination and trimming of the mouse metagenome 75 and 150 bp reads were performed using
608 the BMap short read aligner (v. 36.19)(25) as previously described(52). Briefly, for decontamination,
609 raw reads were mapped to the internal Illumina control phiX174 (J02482.1), the mouse (mm10) and
610 human (hg38) reference genomes using the bbsplit algorithm with default settings. The resulting
611 unmapped reads were adapter trimmed and low-quality reads and reads of insufficient length were
612 removed using the bbdduk algorithm with the following parameters: ktrim = lr, k = 20, mink = 4,
613 minlength = 20, qtrim = f. The reads were assembled using SPAdes version 3.6.1(53) with the following
614 parameters: --only-assembler -k 21,33,55,77,99. Assembled contigs <40 kbp were discarded. The
615 assembly resulted in 2143 contigs >40 kbp.

616 **Taxonomic classification and annotation of metagenomic contigs**

617 The 2143 contigs >40 kbp from the assembly of the mouse fecal metagenome were taxonomically
618 classified using the Contig Annotation Tool (CAT)(42) (version 2019-07-19). Genes were predicted and
619 annotated using the automated PROKKA pipeline (v1.11) (54).

620 **Read mapping and read coverage visualization**

621 The whole (meta)genome and purified VLP read sets were mapped onto the corresponding reference
622 genomes of pure culture organisms or the mouse fecal metagenome contigs using BMap(25) with the
623 following parameters: ambiguous=random qtrim=lr minid=0.97. The generated read mapping files (.bam)
624 were sorted and indexed using SAMtools(v. 1.7)(55). Integrative Genomics Viewer (IGV, v. 2.3.67) tools
625 were used to generate tiled data files (.tdf) from the read mapping (.bam) files for data compression and
626 faster access in IGV using the following parameters: count command, zoom levels: 9, using the mean,
627 window size: 25 or 100(26). Read coverage patterns were displayed and visually assessed in IGV using a
628 linear or if needed log scale.

629 **Data availability**

630 All sequencing read data generated for this study is available from the European Nucleotide Archive
631 (ENA) through study PRJEB33536 (<https://www.ebi.ac.uk/ena/data/view/PRJEB33536>). This data
632 includes the reads for the mouse whole metagenomes and VLP fraction, *E. faecalis* VLPs, *B. subtilis*
633 VLPs, *E. coli* phage P1, *E. coli* phage λ , the whole genome sequencing of *E. coli* KL740 (*E. coli* with
634 lambda phage integrated) and the contigs >40 kbp from the mouse whole metagenome assembly (contig
635 file accession number ERZ1273841).

636 In addition to the de novo generated data we used existing genome assemblies and sequencing read sets
637 for individual bacteria including *Bacillus subtilis* subsp. subtilis str. 168 complete genome from NCBI
638 RefSeq (NC_000964.3), *B. subtilis* 168 genome sequencing read set from the ENA (Study: PRJDB1076,
639 Sample: SAMD00008600), *Enterococcus faecalis* V583 complete genome from NCBI RefSeq
640 (NC_004668.1), *E. faecalis* V583 sequencing read set from the ENA (Study: PRJEB13005, Sample:
641 ERS1085927), *Escherichia coli* K12 complete genome from NCBI RefSeq (NC_000913.3), *E. coli*
642 NCM3722 (*E. coli* K12 with Lambda phage integrated) complete genome sequence from NCBI GenBank
643 (CP011495.1), phage P1 complete genome sequence from NCBI RefSeq (NC_005856.1), *E. coli* K12
644 genome sequencing read set from the ENA (Study: PRJNA251794, Sample: SAMN02840692),
645 *Salmonella enterica* subsp. enterica serovar typhimurium str. LT2 complete genome sequence from NCBI
646 RefSeq (NC_003197.1), *S. typhimurium* LT2 genome sequencing read set from ENA (Study:
647 PRJNA203445, Sample: SAMN02367645), and a read set of CsCl density gradient purified P22 phage
648 (Study: PRJEB6941, Sample: SAME2690949)(24).

649 **Author Contributions**

650 M.K. and B.A.D. designed the study. M.K. and B.A.D. performed experiments. M.K. and B.A.D.
651 performed bioinformatic analyses. B.B. developed BBTools and implemented new parameters in BBMap
652 needed for analyses performed in this study. K.S. and L.V.H. provided conceptual input, strains and
653 specialized reagents. M.K. and B.A.D. wrote the paper with input from all authors.

654 **Acknowledgements**

655 We would like to thank Kelly Ruhn for assistance with animals and Vanessa Schmid and Rachel Bruce of
656 the University of Texas Southwestern Medical Center's Next Generation Sequencing Core for assistance
657 with Illumina library construction and sequencing support.

658 This work was supported in part by the USDA National Institute of Food and Agriculture Hatch project
659 1014212 (M.K.), the Foundation for Food and Agriculture Research Grant ID: 593607 (M.K.), the NC
660 State Chancellor's Faculty Excellence Program Cluster on Microbiomes and Complex Microbial
661 Communities (M.K.), National Institutes of Health Grants R01AI141479 (B.A.D.), K01DK102436
662 (B.A.D), and R01DK070855 (L.V.H.), and the Howard Hughes Medical Institute (L.V.H.).

663

664 **Competing Interests**

665 The authors declare no competing interests.

666 References

- 667 1. O. Zhaxybayeva, W. F. Doolittle, Lateral gene transfer. *Curr. Biol.* **21**, R242–R246 (2011).
- 668 2. S. M. Soucy, J. Huang, J. P. Gogarten, Horizontal gene transfer: building the web of life. *Nat. Rev.*
669 *Genet.* **16**, 472–482 (2015).
- 670 3. C. S. Smillie, *et al.*, Ecology drives a global network of gene exchange connecting the human
671 microbiome. *Nature* **480**, 241–244 (2011).
- 672 4. A. Oladeinde, *et al.*, Horizontal Gene Transfer and Acquired Antibiotic Resistance in *Salmonella*
673 *enterica* Serovar Heidelberg following In Vitro Incubation in Broiler Ceca. *Appl. Environ.*
674 *Microbiol.* **85** (2019).
- 675 5. S. Borgeaud, L. C. Metzger, T. Scignari, M. Blokesch, The type VI secretion system of *Vibrio*
676 *cholerae* fosters horizontal gene transfer. *Science* **347**, 63–67 (2015).
- 677 6. J. Chen, R. P. Novick, Phage-mediated intergeneric transfer of toxin genes. *Science* **323**, 139–141
678 (2009).
- 679 7. J.-H. Hehemann, *et al.*, Transfer of carbohydrate-active enzymes from marine bacteria to Japanese
680 gut microbiota. *Nature* **464**, 908–912 (2010).
- 681 8. E. F. Mongodin, *et al.*, The genome of *Salinibacter ruber*: Convergence and gene exchange among
682 hyperhalophilic bacteria and archaea. *Proc. Natl. Acad. Sci.* **102**, 18147–18152 (2005).
- 683 9. O. Popa, G. Landan, T. Dagan, Phylogenomic networks reveal limited phylogenetic range of lateral
684 gene transfer by transduction. *ISME J.* **11**, 543–554 (2017).
- 685 10. A. S. Lang, O. Zhaxybayeva, J. T. Beatty, Gene transfer agents: phage-like elements of genetic
686 exchange. *Nat. Rev. Microbiol.* **10**, 472–482 (2012).
- 687 11. Y. N. Chiang, J. R. Penadés, J. Chen, Genetic transduction by phages and chromosomal islands: The
688 new and noncanonical. *PLOS Pathog.* **15**, e1007878 (2019).
- 689 12. J. Chen, *et al.*, Genome hypermobility by lateral transduction. *Science* **362**, 207–212 (2018).
- 690 13. A. Thierauf, G. Perez, and S. Maloy, “Generalized Transduction” in *Bacteriophages: Methods and*
691 *Protocols, Volume 1: Isolation, Characterization, and Interactions*, Methods in Molecular
692 Biology™, M. R. J. Clokie, A. M. Kropinski, Eds. (Humana Press, 2009), pp. 267–286.
- 693 14. G. E. Christie, T. Dokland, Pirates of the Caudovirales. *Virology* **434**, 210–221 (2012).
- 694 15. R. C. Matos, *et al.*, *Enterococcus faecalis* Prophage Dynamics and Contributions to Pathogenic
695 Traits. *PLOS Genet.* **9**, e1003539 (2013).
- 696 16. S. C. Jiang, J. H. Paul, Gene Transfer by Transduction in the Marine Environment. *Appl. Environ.*
697 *Microbiol.* **64**, 2780–2787 (1998).
- 698 17. T. Kenzaka, K. Tani, M. Nasu, High-frequency phage-mediated gene transfer in freshwater
699 environments determined at single-cell level. *ISME J.* **4**, 648–659 (2010).

- 700 18. M. Sander, H. Schmieger, Method for host-Independent detection of generalized transducing
701 bacteriophages in natural habitats. *Appl. Environ. Microbiol.* **67**, 1490–1493 (2001).
- 702 19. J. Tomasch, *et al.*, Packaging of Dinoroseobacter shibae DNA into Gene Transfer Agent Particles Is
703 Not Random. *Genome Biol. Evol.* **10**, 359–369 (2018).
- 704 20. C. Pourcel, C. Midoux, Y. Hauck, G. Vergnaud, L. Latino, Large Preferred Region for Packaging of
705 Bacterial DNA by phiC725A, a Novel Pseudomonas aeruginosa F116-Like Bacteriophage. *PLOS*
706 *ONE* **12**, e0169684 (2017).
- 707 21. J. R. Garneau, F. Depardieu, L.-C. Fortier, D. Bikard, M. Monot, PhageTerm: a tool for fast and
708 accurate determination of phage termini and packaging mechanism using next-generation
709 sequencing data. *Sci. Rep.* **7**, 1–10 (2017).
- 710 22. A. Reyes, M. Wu, N. P. McNulty, F. L. Rohwer, J. I. Gordon, Gnotobiotic mouse model of phage-
711 bacterial host dynamics in the human gut. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 20236–20241 (2013).
- 712 23. R. V. Thurber, M. Haynes, M. Breitbart, L. Wegley, F. Rohwer, Laboratory procedures to generate
713 viral metagenomes. *Nat. Protoc.* **4**, 470–483 (2009).
- 714 24. M. Kleiner, L. V. Hooper, B. A. Duerkop, Evaluation of methods to purify virus-like particles for
715 metagenomic sequencing of intestinal viromes. *BMC Genomics* **16**, 7 (2015).
- 716 25. B. Bushnell, BBMap. *SourceForge* (December 29, 2019).
- 717 26. H. Thorvaldsdóttir, J. T. Robinson, J. P. Mesirov, Integrative Genomics Viewer (IGV): high-
718 performance genomics data visualization and exploration. *Brief. Bioinform.* **14**, 178–192 (2013).
- 719 27. M. E. Gottesman, R. A. Weisberg, Little Lambda, who made thee? *Microbiol. Mol. Biol. Rev.* **68**,
720 796–813 (2004).
- 721 28. M. L. Morse, E. M. Lederberg, J. Lederberg, Transduction in Escherichia Coli K-12. *Genetics* **41**,
722 142–156 (1956).
- 723 29. J. Ebel-Tsipis, D. Botstein, M. S. Fox, Generalized transduction by phage P22 in Salmonella
724 typhimurium: I. Molecular origin of transducing DNA. *J. Mol. Biol.* **71**, 433–448 (1972).
- 725 30. H. Schmieger, Packaging signals for phage P22 on the chromosome of Salmonella typhimurium.
726 *Mol. Gen. Genet. MGG* **187**, 516–518 (1982).
- 727 31. S. Casjens, M. Hayden, Analysis in vivo of bacteriophage P22 headful nuclease. *J. Mol. Biol.* **199**,
728 467–474 (1988).
- 729 32. M. C. Hanks, B. Newman, I. R. Oliver, M. Masters, Packaging of transducing DNA by
730 bacteriophage P1. *Mol. Gen. Genet. MGG* **214**, 523–532 (1988).
- 731 33. M. Masters, “Generalized Transduction” in *Escherichia Coli and Salmonella: Cellular and*
732 *Molecular Biology*, (American Society for Microbiology, 1996), pp. 2421–441.
- 733 34. José R. Penadés, Gail E. Christie, The phage-inducible chromosomal islands: A family of highly
734 evolved molecular parasites. *Annu. Rev. Virol.* **2**, 181–201 (2015).

- 735 35. B. A. Duerkop, C. V. Clements, D. Rollins, J. L. M. Rodrigues, L. V. Hooper, A composite
736 bacteriophage alters colonization by an intestinal commensal bacterium. *Proc. Natl. Acad. Sci.* **109**,
737 17621–17626 (2012).
- 738 36. F. Kunst, *et al.*, The complete genome sequence of the Gram-positive bacterium *Bacillus subtilis*.
739 *Nature* **390**, 249–256 (1997).
- 740 37. H. E. Wood, M. T. Dawson, K. M. Devine, D. J. McConnell, Characterization of PBSX, a defective
741 prophage of *Bacillus subtilis*. *J. Bacteriol.* **172**, 2667 (1990).
- 742 38. C. Canchaya, G. Fournous, S. Chibani-Chennoufi, M.-L. Dillmann, H. Brüssow, Phage as agents of
743 lateral gene transfer. *Curr. Opin. Microbiol.* **6**, 417–424 (2003).
- 744 39. T. Korem, *et al.*, Growth dynamics of gut microbiota in health and disease inferred from single
745 metagenomic samples. *Science* **349**, 1101 (2015).
- 746 40. L. M. Anderson, K. F. Bott, DNA packaging by the *Bacillus subtilis* defective bacteriophage PBSX.
747 *J. Virol.* **54**, 773 (1985).
- 748 41. K. Abe, *et al.*, Developmentally-Regulated Excision of the SP β Prophage Reconstitutes a Gene
749 Required for Spore Envelope Maturation in *Bacillus subtilis*. *PLOS Genet.* **10**, e1004636 (2014).
- 750 42. F. A. B. von Meijenfheldt, K. Arkhipova, D. D. Cambuy, F. H. Coutinho, B. E. Dutilh, Robust
751 taxonomic classification of uncharted microbial sequences and bins with CAT and BAT. *Genome*
752 *Biol.* **20**, 217 (2019).
- 753 43. X. He, *et al.*, Cultivation of a human-associated TM7 phylotype reveals a reduced genome and
754 epibiotic parasitic lifestyle. *Proc. Natl. Acad. Sci.* **112**, 244–249 (2015).
- 755 44. M. Bonazzi, M. Lecuit, P. Cossart, *Listeria monocytogenes* Internalin and E-cadherin: From Bench
756 to Bedside. *Cold Spring Harb. Perspect. Biol.* **1** (2009).
- 757 45. B. Cousineau, S. Lawrence, D. Smith, M. Belfort, Retrotransposition of a bacterial group II intron.
758 *Nature* **404**, 1018–1021 (2000).
- 759 46. J. Haaber, *et al.*, Bacterial viruses enable their host to acquire antibiotic resistance genes from
760 neighbouring cells. *Nat. Commun.* **7**, 1–8 (2016).
- 761 47. F. Enault, *et al.*, Phages rarely encode antibiotic resistance genes: a cautionary tale for virome
762 analyses. *ISME J.* **11**, 237–247 (2017).
- 763 48. S. R. Modi, H. H. Lee, C. S. Spina, J. J. Collins, Antibiotic treatment expands the resistance
764 reservoir and ecological network of the phage metagenome. *Nature* **499**, 219–222 (2013).
- 765 49. W. Calero-Cáceres, M. Ye, J. L. Balcázar, Bacteriophages as Environmental Reservoirs of
766 Antibiotic Resistance. *Trends Microbiol.* **27**, 570–577 (2019).
- 767 50. N. J. Bitto, *et al.*, Bacterial membrane vesicles transport their DNA cargo into host cells. *Sci. Rep.* **7**,
768 1–11 (2017).
- 769 51. S. Fulsundar, *et al.*, Gene Transfer Potential of Outer Membrane Vesicles of *Acinetobacter baylyi*
770 and Effects of Stress on Vesiculation. *Appl. Environ. Microbiol.* **80**, 3469–3483 (2014).

- 771 52. B. A. Duerkop, *et al.*, Murine colitis reveals a disease-associated bacteriophage community. *Nat.*
772 *Microbiol.* **3**, 1023–1031 (2018).
- 773 53. A. Bankevich, *et al.*, SPAdes: a new genome assembly algorithm and its applications to single-cell
774 sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).
- 775 54. T. Seemann, Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).
- 776 55. H. Li, *et al.*, The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079
777 (2009).
- 778