

# Identification of type IV conjugative systems that are systematically excluded from metagenomic bins

**Benjamin R. Joris**

The University of Western

**Tyler S. Browne**

The University of Western

**Thomas A. Hamilton**

The University of Western

**David R. Edgell**

The University of Western

**Gregory B. Gloor** (✉ [ggloor@uwo.ca](mailto:ggloor@uwo.ca))

The University of Western

---

## Research Article

### Keywords:

**Posted Date:** March 9th, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-1428512/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 Identification of type IV conjugative systems that are  
2 systematically excluded from metagenomic bins

3 Benjamin R. Joris, Tyler S. Browne, Thomas A. Hamilton, David R. Edgell and Gregory B. Gloor

4 07 March, 2022

5 Author Affiliations: Department of Biochemistry, Schulich School of Medicine, The University of Western  
6 Ontario, London, Ontario, N6A 3K7, Canada

7 Corresponding author: Gregory B. Gloor, ggloor@uwo.ca

## 8 Abstract

### 9 Background

10 Conjugation enables the exchange of genetic elements throughout environments, including the human gut  
11 microbiome. Conjugative elements can carry and transfer clinically relevant metabolic pathways which makes  
12 precise identification of these systems in metagenomic samples clinically important.

### 13 Results

14 Here, we outline two related methods to identify conjugative elements in the human gut microbiome. Con-  
15 jugative systems can effectively be identified from metagenomic assemblies either using a curated sets of  
16 profile hidden Markov models or by searching against large-scale databases, such as UniRef90. Both meth-  
17 ods were successful at identifying type IV conjugative systems with profile hidden Markov models being  
18 faster, but less sensitive than alignment to the UniRef database. Finally, we demonstrate that the ma-  
19 jority of assembled conjugative elements are not included within metagenomic bins, and that only a small  
20 proportion of the binned conjugative systems are included in “high-quality” metagenomic bins.

### 21 Conclusions

22 Analysis of the human gut microbiome by shotgun metagenomic sequencing has revealed numerous con-  
23 nections to human health outcomes. Our findings emphasize the need to identify and analyze conjugative  
24 systems outside of standard metagenomic binning pipelines. We suggest that analysis of type IV conjugative  
25 systems should be added to the current metagenomic analysis approaches as they contain much information  
26 that could explain differences between cohorts.

## 27 Background

28 Bacteria can acquire exogenous DNA through horizontal gene transfer. Conjugation is a common mechanism  
29 of horizontal gene transfer that relies on direct cell-cell contact to unidirectionally transfer DNA from a  
30 bacterial donor to a recipient cell. In bacteria, integrative conjugative elements (ICEs) and conjugative  
31 plasmids are mobilizable through the actions of type IV secretion systems (T4SS). Approximately half of the  
32 known plasmids are mobilizable in *trans* where the conjugative machinery is on a different genetic element  
33 than the transferred element, and the remainder are mobilizable in *cis* because the conjugative machinery is  
34 present on the same genetic element [1]. ICEs encode their own T4SS, and can mobilize other elements [2].  
35 Conjugative elements (CEs) often contain antibiotic resistance genes, but also can harbour useful biosynthetic  
36 and biodegradation genes [3]. Furthermore, conjugative systems can serve as vectors to introduce clustered  
37 regularly interspaced short palindromic repeats (CRISPR) systems, metabolic pathways or novel functions  
38 into the gut microbiota [4–9]. Therefore, characterizing the full complement of conjugative systems in  
39 the human gut could expand the number of useable vectors for these applications. Precise identification of

40 conjugative systems from metagenomic samples could also provide insights to their distribution in populations  
41 and their correlation with antibiotic exposure, age, and health status.

42 For a DNA sequence to be considered mobilizable by conjugation, it must encode an origin of transfer (*oriT*)  
43 sequence that is recognized and nicked by a relaxase protein [1, 10]. Relaxase proteins contain a conserved  
44 histidine triad that coordinates a divalent metal ion, as well as tyrosine residues that catalyze the nicking  
45 reaction at the *oriT* DNA sequence [11, 12]. In addition to a relaxase gene and an *oriT* sequence, a full  
46 complement of type IV secretion system and coupling proteins are required for a sequence to be conjugative.  
47 In the well-studied *Agrobacterium tumefaciens* conjugative system, there are 12 proteins involved in the  
48 transfer of the DNA-relaxase complex from one bacterial cell to another [13, 14]. Homologs of the VirB4  
49 ATPase that are essential for assembly of the conjugative system and DNA transfer are generally similar to  
50 the phylogeny of the bacteria harbouring them [15] and thus are useful for classifying conjugative systems  
51 [16]. The synteny of conjugative transfer genes is also highly conserved among conjugative systems [14].  
52 Both the synteny and presence of highly-conserved genes involved in conjugation facilitates the classification  
53 of genetic elements as potentially conjugative if the sequences are annotated as belonging to the components  
54 of the T4SS [17] (Figure 1).

55 Previous work has identified novel CEs in the human and animal gut microbiomes, but the focus was mainly  
56 on ICEs and not on conjugative plasmids [3, 18, 19]. Identifying conjugative plasmids from a short-read  
57 metagenomic assembly is difficult for several reasons. The initial barrier is the difficulty in assembling  
58 circularized plasmids from short-read sequencing data [20]. A second barrier is that the contiguous DNA  
59 sequences (contigs) that compose metagenome-assembled genomes (MAGs) are binned together based on  
60 sequence composition and coverage. Binning of a plasmid with its cognate genome will not happen unless  
61 the contigs that compose the plasmid are maintained in the same copy-number and have the same sequence  
62 composition as the chromosome. These criteria are generally not met because conjugative systems are usually  
63 more AT rich than the cognate chromosome [1] and often do not have a unit copy number. Since nearly 80%  
64 of the non-redundant set of genomes from the human-gut microbiome are from difficult-to-culture species that  
65 are known only from MAGs [21], alternate methods must be employed to assemble and identify conjugative  
66 plasmids from the metagenomic sequencing data. Computational tools have recently been developed to  
67 identify plasmids from metagenomic assemblies [22], but would be less than optimal if applied to already  
68 binned data that systematically excludes plasmids [23]. Methods that identify CEs prior to binning should  
69 be able to capture the full spectrum of ICEs and conjugative plasmids.

70 Here, we show that T4SS conjugative systems can be identified using two distinct methods (Figure 2). First,  
71 we used profile HMMs (pHMMs) to identify conjugative systems directly from metagenomic assemblies of  
72 North American inflammatory bowel disease (IBD) and North American pre-term infant samples. Second, we  
73 searched predicted protein sequences from those same assemblies versus UniRef90 [24] for proteins involved  
74 in conjugation to identify conjugative systems. Both methods were able to find conjugative systems in raw

75 metagenomic assemblies with pHMMs being computationally more efficient but less sensitive. Finally, we  
76 demonstrate that the majority of conjugative systems produced by a metagenomic assembly are not included  
77 in high-quality bins that are used as proxies for bacterial genomes in metagenomic analysis pipelines. Our  
78 findings provide a roadmap to integrate the analysis of conjugative systems alongside the chromosomal  
79 content of bacteria.

## 80 **Methods**

### 81 **Assembly and identification of conjugative systems in North American short-** 82 **read data**

83 Samples belonging to a North American IBD (n=50) [25] and a North American pre-term infant cohort  
84 (n=51) [26] were assembled *de novo* as follows (Supplemental Table 1). Reads from these samples were down-  
85 loaded from the Sequence Read Archive using the SRA toolkit version 2.9.2, deduplicated with `dedupe.sh`  
86 [27], and trimmed with Trimmomatic version 0.36 [28] with options `LEADING:10 TRAILING:10`. Processed  
87 reads were assembled sample-by-sample using SPAdes version 3.14.0, option `--meta` [29].

### 88 **Identification of conjugative systems using Profile hidden Markov models**

89 The resultant assemblies were imported into Anvi'o version 6.0 [30] where the presence of T4SS, T4CP, and  
90 relaxase proteins were predicted using the `anvi-run-hmms` module, which integrates HMMER3 functionality  
91 [31]. Instructions for installation of T4SS pHMMs into Anvi'o can be found in the online code repository  
92 ([https://github.com/bjoris33/humanGutConj\\_Microbiome](https://github.com/bjoris33/humanGutConj_Microbiome)). Contigs that contained pHMM matches for all  
93 three classes of conjugative proteins were extracted and annotated by aligning open reading frames (ORFs)  
94 predicted with Prodigal version 2.6.3 [32] to the UniRef90 database [24]. Taxonomic prediction of the contigs  
95 was conducted with Kaiju version 1.7.2 utilizing the RefSeq non-redundant protein database [33]. MOB-suite  
96 verion 1.4.9.1 was utilized to characterize the incompatibility grouping of the conjugative system, if possible  
97 [34]. PlasFlow version 1.1.0 was used to classify whether the system was chromosomally integrated or a  
98 plasmid [22].

### 99 **Identification of conjugative systems using protein alignments to the UniRef90** 100 **database**

101 The contigs of the raw metagenomic assemblies had their ORFs predicted using Prodigal version 2.6.3 [32].  
102 The predicted ORFs were then aligned to the UniRef90 database [24] using the `blastp` module of DIAMOND  
103 version 0.9.14 [35]. By using keywords such as “conjugal” or “mobilization”, the protein alignments were  
104 search for contigs that contained annotations for both a relaxase and either a type IV secretion system  
105 protein or a type IV coupling protein. Through manual curation, type IV secretion system proteins and

106 coupling proteins often shared identical or very similar annotation entries in the UniRef90 database, so the  
107 decision was made not to distinguish between the two.

## 108 **Binning of Assemblies**

109 For each assembly, all 101 samples were mapped to the contigs using Bowtie2 [36]. The mapping files  
110 were sorted and indexed with SAMtools [37] and then the assemblies were binned using MetaBAT2 version  
111 2.12.1 [38]. CheckM version 1.1.2 was used to assess the quality of the resultant bins [39]. High-quality  
112 bins were defined using the same cutoffs (>90% completion and <5% redundancy) as Almeida *et al* (2019)  
113 defined. Bins not passing that threshold were classified as “low-quality”. The previously identified contigs  
114 with conjugative systems were classified based on their presence in bins, and the types of bins they were  
115 present in. Results of this classification were visualized using SankeyMATIC (<http://sankeymatic.com/>).

## 116 **Results**

### 117 **Profile hidden Markov models and database alignment successfully identify con-** 118 **jugative elements from metagenomic assemblies**

119 Fifty-one samples from a pre-term infant cohort and 50 from a North American IBD cohort were assembled  
120 sample-by-sample using metaSPAdes [29] to identify T4SS conjugative systems from a full pool of assembled  
121 contigs (i.e. not binned). Two separate methods we employed to search for contigs containing type IV  
122 conjugative proteins. For the pHMM method of identifying conjugative systems, contigs with conjugative  
123 systems were defined by pHMM matches for a relaxase, a type IV coupling protein, and a type IV secretion  
124 system, which offers a fast and precise method to annotate a limited number of protein families. From the  
125 assembly of the pre-term infant cohort 96 of 470500 contigs met the criteria, whereas 268 of 15100646 contigs  
126 from the IBD cohort did.

127 The second method of identifying conjugative systems utilizes the UniRef90 database by aligning the pre-  
128 dicted ORFs to it using DIAMOND [35]. The alignment results are searched using a keyword strategy for  
129 contigs that contain an alignment for a relaxase or mobilization protein and an alignment for a type IV  
130 secretion or type IV coupling protein. From the pre-term infant cohort assemblies 242 of 470500 contigs met  
131 the described criteria, and 4244 of the 15100646 contigs from the IBD cohort met the same criteria.

132 The two outlined methods represent potentially complimentary methods of tackling the same problem–  
133 identifying conjugative systems from a pool of metagenomic-assembled contigs. There is a large-degree of  
134 overlap between the two methods, however alignment to the UniRef90 database appears to be much more  
135 sensitive with only 280 of the 4486 identified conjugative systems also identified using the pHMM method  
136 (Figure 3). While it may be less sensitive, the pHMM method of identifying conjugative systems has a much  
137 smaller computationally footprint as it does not rely on aligning to a large protein database, but rather to

138 a small and specific set of profile hidden Markov models.

## 139 **The majority of conjugative systems identified are omitted from metagenomic** 140 **bins**

141 Metagenomic assemblies from two distinct cohorts were binned using MetaBAT2 [38] to explore how conjuga-  
142 tive systems are distributed within common metagenomic analyses. Of the 364 assembled contigs containing  
143 pHMM matches to all three protein categories, 270 were not included in any metagenomic bins (Figure 4).  
144 For the 94 contigs included in metagenomic bins, 65 of those were found in high-quality bins (>90% comple-  
145 tion and <5% redundancy). This is in stark contrast to the background binning rate of contigs; For contigs  
146 above 5kb in size the binning rate with MetaBAT2 [38] was 70.4% (116112/164843 contigs) and for contigs  
147 above 10kb the binning rate was 79.1% (57214/72300 contigs). Among the 29 contigs included in bins that  
148 do not meet the aforementioned threshold, 8 are within bins that are less than or equal to 1 megabase in  
149 size, potentially suggesting that fragments of a conjugative plasmid may have binned together.

150 For the conjugative systems established using alignments to the UniRef90 database, there is an even lower  
151 rate of binning (Figure 5). Of the 4486 conjugative systems, only 287 of them were binned—a rate of 6.4%.  
152 Again, a number of the bins that do form are low quality bins below 1mb in size that may be the collection  
153 of contigs that form a conjugative plasmid.

## 154 **Discussion**

155 To produce MAGs, contigs generated by metagenomic assembly are typically binned using a program such  
156 as MetaBAT2 [38]. Conjugative systems are often more AT rich than the parent genomes [1], which would  
157 result in the conjugative system and cognate genome not occurring in the same metagenomic bin because  
158 binning algorithms use GC content as a parameter for clustering. Additionally, plasmids are not necessarily  
159 maintained in a unit copy number within the cell, causing differential sequence coverage in comparison to  
160 the parent genome, which is another factor that leads to plasmids being excluded from MAGs. Therefore to  
161 capture a more complete image of the conjugative systems present in an environment, identification of the  
162 systems must take place before binning.

163 We have outlined two methods for identifying contigs carrying potentially functional type IV conjugative  
164 systems from raw metagenomic assemblies. Using a curated set of pHMMs of the three main classes of type  
165 IV conjugative proteins (relaxases, secretion proteins, and coupling proteins), we were able to classify 364  
166 total contigs as being potentially conjugative. In comparison, the method that utilized predicted protein  
167 alignments to the UniRef90 database found 4486 contigs that met the criteria, which indicates that it may  
168 be the more sensitive method. However, aligning all predicted open reading frames found in a metagenomic  
169 assembly to the full UniRef90 database is a computationally expensive task. Considering that many of

170 the conjugative systems identified by pHMMs are also found by the protein alignment method (280 of 364  
171 contigs), using pHMMs may be appropriate in as a first pass method or in situations where computational  
172 resources are scarce.

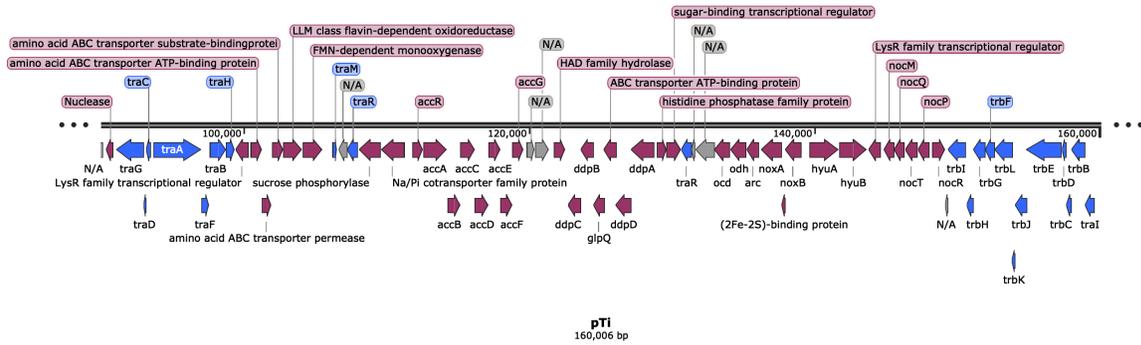
173 The assembled contigs were binned with MetaBAT2 [38] as a way of quantifying the effect of binning, which  
174 revealed that the vast majority of the assembled conjugative systems were not included in metagenomic bins  
175 and therefore would not be included in a MAG database, which confirms recent findings [23]. The binning  
176 rate of contigs carrying type IV conjugative systems identified by pHMMs and alignment to UniRef90 was  
177 considerably lower than the background binning rate of equivalently sized contigs (25.8% and 6.4% compared  
178 to 70.4%, respectively). Many of the binned conjugative systems were not within a bin that would pass the  
179 quality cutoff to be included in a curated MAG genome set as well [21]. Interestingly, eight of the conjugative  
180 systems were binned into low-quality bins that were smaller than <1MB in size, which may suggest that the  
181 fragments of a conjugative plasmid could be binned together, which would increase the completeness of the  
182 conjugative system.

## 183 Conclusions

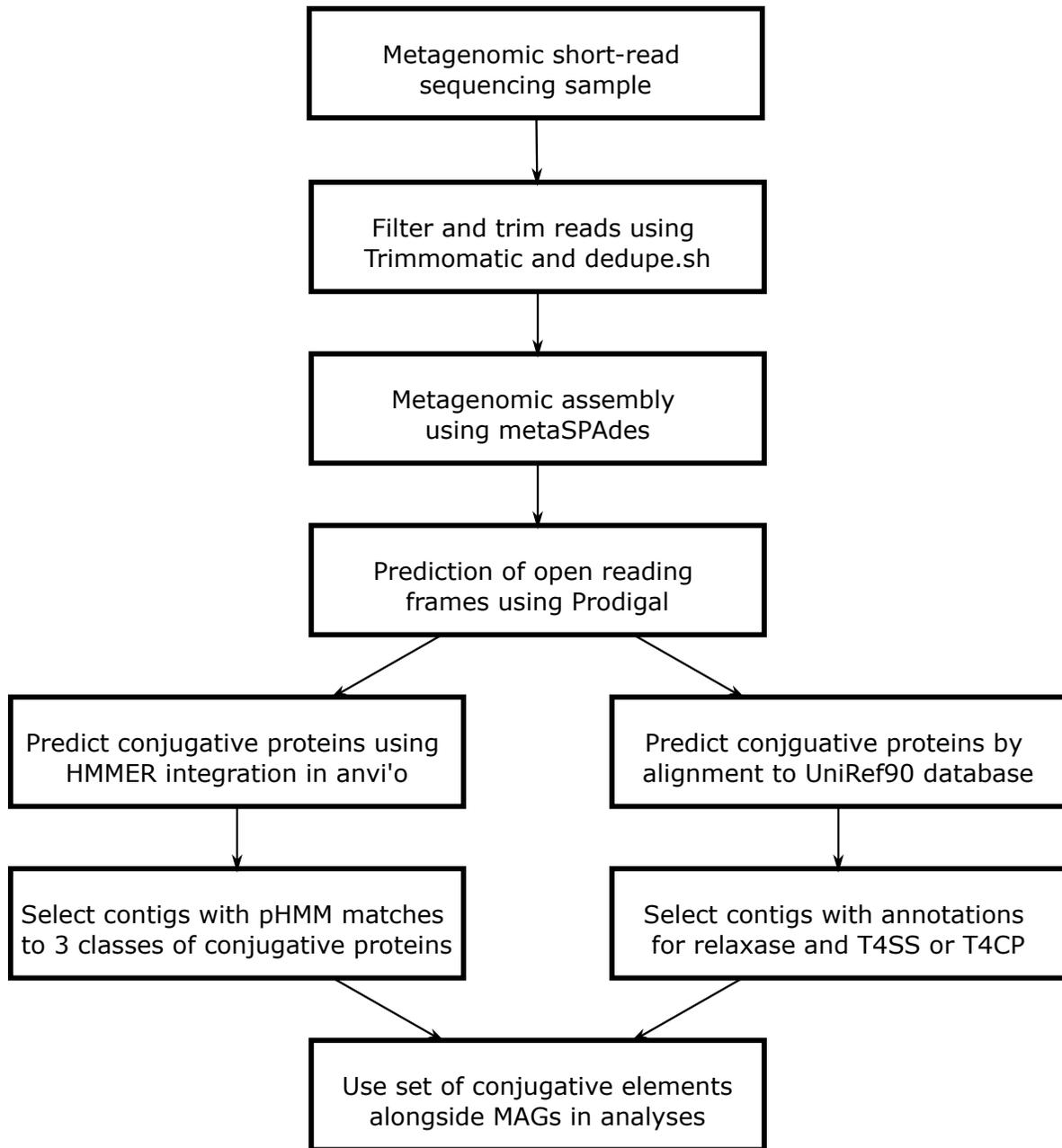
184 Conjugative systems could differ between cohorts and require special consideration to ensure that they are  
185 included in metagenomic analyses. ICEs and plasmids can carry harmful systems, such as antimicrobial re-  
186 sistance, but also can act as vectors for bile salt metabolism and for detoxification modules [3]. These cargo  
187 genes are relevant for research relating to the gut microbiome’s role in pathogenicity as well as metabolism,  
188 digestion, and host effector molecules. Comprehensive identification and quantification of conjugative sys-  
189 tems could allow for association of conjugative systems with different health outcomes. Because assembled  
190 plasmid-based conjugative systems are rarely included in metagenomic bins [23] (Figure 4 and Figure 5),  
191 they need to be identified and analyzed outside of standard binning pipelines. At present, it is not possible  
192 to assemble complete plasmids from short-read metagenomic data [20], so it may helpful to identify bins  
193 containing conjugative systems in an attempt to cluster the fragments of plasmids present in an assembly  
194 together. Identifying type IV conjugative systems using pHMMs or UniRef90 annotations and using tools  
195 such as PlasFlow [22] to identify plasmids out of a full assembly in parallel with standard binning analyses  
196 will enhance research of the associations between the human gut microbiome and human health.

197 In the future, expanding the curated set of pHMMs could increase the sensitivity of the method to detect  
198 conjugative systems and creating a curated set of conjugation proteins from the UniRef90, instead of exhaus-  
199 tively annotating with the full database, could improve the computational efficiency of the other method.  
200 Additionally, improvements in assembly and binning algorithms will continue to improve the recovery of  
201 low relative abundance conjugative elements and improve the completeness and accuracy of the assembled  
202 fragments. For instance, long-read assembly permits the circularization of genomes and plasmids [40, 41]  
203 and the binning of plasmids to their cognate genomes using methylation data [42], which will reduce the

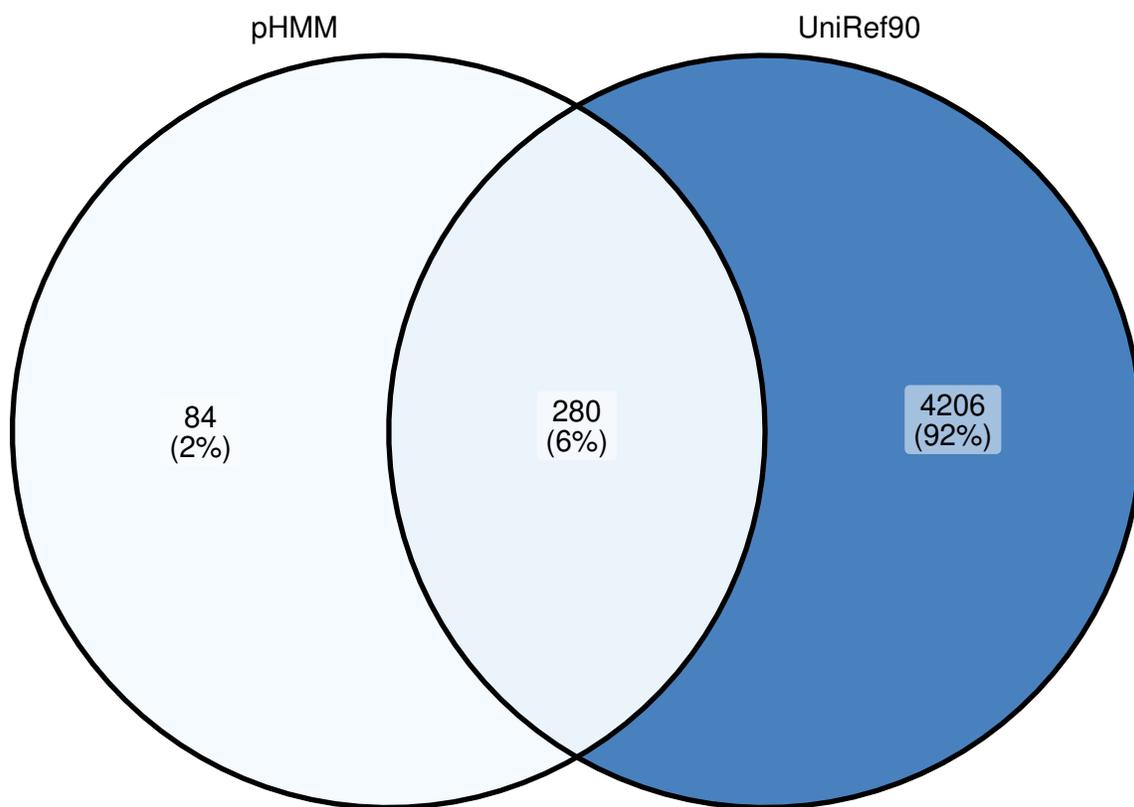
204 ambiguity of the origins of conjugative systems (i.e. whether they are an ICE or independently circularized  
 205 plasmid) and provide a more complete picture of the cargo they carry and the differences between cohorts.



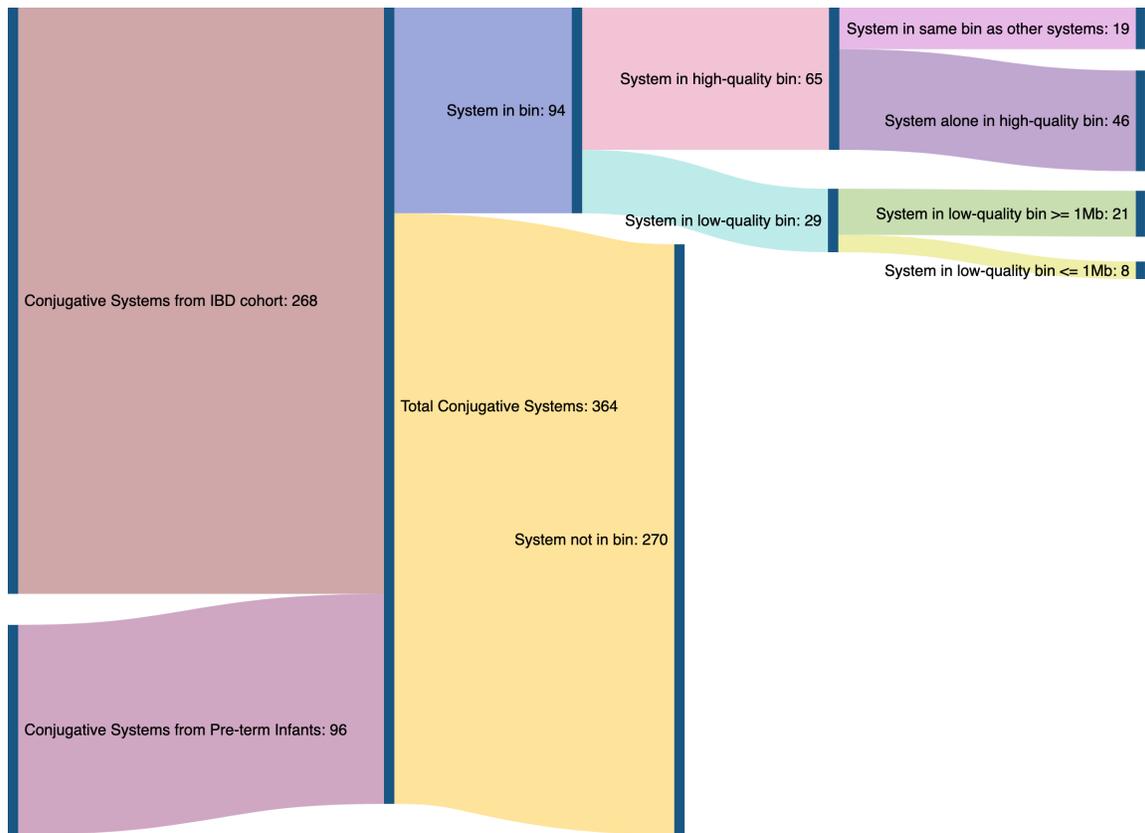
**Figure 1:** Example schematic of the gene organization of a bacterial conjugation system on the *Agrobacterium tumefaciens* pTi plasmid.



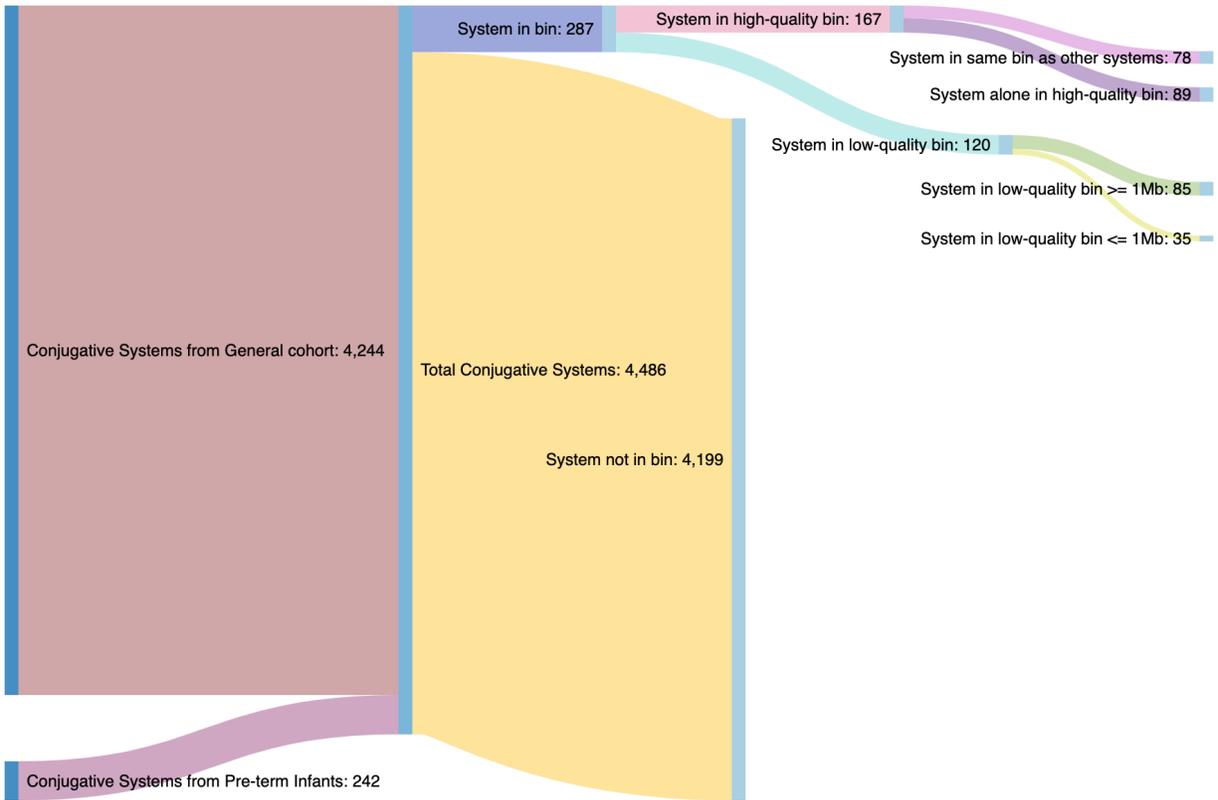
**Figure 2:** Overview of methods employed in this study. In the left panel is the workflow used to identify conjugative systems from previously assembled human gut bacterial genomes. The right panel outlines the workflow for the assembly of select North American samples and the use of pHMMs to identify the conjugative systems.



**Figure 3:** Venn diagram illustrating the overlap of the two methods of identifying type IV conjugative systems from the 101 assembled metagenomic samples.



**Figure 4:** Sankey diagram representing the flow of 364 contigs containing conjugative systems identified using pHMMs into bins generated by MetaBAT2 from assembled data.



**Figure 5:** Sankey diagram representing the flow of 4486 contigs containing conjugative systems identified using predicted protein alignments to the UniRef90 database into bins generated by MetaBAT2 from assembled data.

206 **Supplemental Table 1:** SRA accession numbers and cohorts for samples used in *de novo* assembly work-  
 207 flow.

## 208 **Declarations**

### 209 **Ethics approval and consent to participate**

210 Not applicable.

### 211 **Consent for publication**

212 Not applicable.

### 213 **Availability of data and materials**

214 All code needed to reproduce the results are available on Github, [https://github.com/bjoris33/](https://github.com/bjoris33/humanGutConj_Microbiome)  
 215 humanGutConj\_Microbiome

## 216 **Competing interests**

217 The authors declare that they have no competing interests.

## 218 **Funding**

219 Supported by CIHR Project Grant (PJT-159708) to D.R.E. and G.B.G. T.A.H. was supported by an NSERC  
220 PGS-D scholarship. B.R.J was supported by the Schulich School of Medicine Dean’s Research Scholarship.

## 221 **Authors’ Contributions**

222 BRJ designed the experiments, analyzed and interpreted the data, and wrote the manuscript. TSB analyzed  
223 the data. TAH interpreted the data. DRE designed the experiments, interpreted the data, and provided  
224 funding. GBG designed the experiments, interpreted the data, edited the manuscript, and provided funding.

## 225 **Acknowledgements**

226 We thank Daniel Giguere for his input on the analyses and figures.

## 227 **References**

- 228 1. Smillie C, Garcillán-Barcia MP, Francia MV, Rocha EPC, Cruz F de la. Mobility of plasmids. *Microbiol*  
229 *Mol Biol Rev.* 2010;74:434–52.
- 230 2. Daccord A, Ceccarelli D, Burrus V. Integrating conjugative elements of the sxt/r391 family trigger the  
231 excision and drive the mobilization of a new class of vibrio genomic islands. *Mol Microbiol.* 2010;78:576–88.
- 232 3. Jiang X, Hall AB, Xavier RJ, Alm EJ. Comprehensive analysis of chromosomal mobile genetic elements  
233 in the gut microbiome reveals phylum-level niche-adaptive gene pools. *PLoS One.* 2019;14:e0223680.
- 234 4. Neil K, Allard N, Grenier F, Burrus V, Rodrigue S. Highly efficient gene transfer in the mouse gut  
235 microbiota is enabled by the incl2 conjugative plasmid tp114. *Commun Biol.* 2020;3:523.
- 236 5. Hamilton TA, Pellegrino GM, Therrien JA, Ham DT, Bartlett PC, Karas BJ, et al. Efficient inter-species  
237 conjugative transfer of a crispr nuclease for targeted bacterial killing. *Nat Commun.* 2019;10:4544.
- 238 6. Peters JM, Koo B-M, Patino R, Heussler GE, Hearne CC, Qu J, et al. Enabling genetic analysis of diverse  
239 bacteria with mobile-crispri. *Nat Microbiol.* 2019;4:244–50.
- 240 7. Citorik RJ, Mimee M, Lu TK. Sequence-specific antimicrobials using efficiently delivered rna-guided  
241 nucleases. *Nat Biotechnol.* 2014;32:1141–5.
- 242 8. Bikard D, Euler CW, Jiang W, Nussenzweig PM, Goldberg GW, Duportet X, et al. Exploiting crispr-cas  
243 nucleases to produce sequence-specific antimicrobials. *Nat Biotechnol.* 2014;32:1146–50.

- 244 9. Gomaa AA, Klumpe HE, Luo ML, Selle K, Barrangou R, Beisel CL. Programmable removal of bacterial  
245 strains by use of genome-targeting crispr-cas systems. *mBio*. 2014;5:e00928–13.
- 246 10. Francia MV, Varsaki A, Garcillán-Barcia MP, Latorre A, Drainas C, Cruz F de la. A classification  
247 scheme for mobilization regions of bacterial plasmids. *FEMS Microbiol Rev*. 2004;28:79–100.
- 248 11. Nash RP, Habibi S, Cheng Y, Lujan SA, Redinbo MR. The mechanism and control of DNA transfer by  
249 the conjugative relaxase of resistance plasmid pCU1. *Nucleic Acids Res*. 2010;38:5929–43.
- 250 12. Becker EC, Meyer RJ. Recognition of oriT for DNA processing at termination of a round of conjugal  
251 transfer. *J Mol Biol*. 2000;300:1067–77.
- 252 13. Fronzes R, Christie PJ, Waksman G. The structural biology of type IV secretion systems. *Nat Rev*  
253 *Microbiol*. 2009;7:703–14.
- 254 14. Cabezón E, Ripoll-Rozada J, Peña A, Cruz F de la, Arechaga I. Towards an integrated model of bacterial  
255 conjugation. *FEMS Microbiol Rev*. 2015;39:81–95.
- 256 15. Bhatti M, Laverde Gomez JA, Christie PJ. The expanding bacterial type IV secretion lexicon. *Res*  
257 *Microbiol*. 164:620–39.
- 258 16. Guglielmini J, Cruz F de la, Rocha EPC. Evolution of conjugation and type IV secretion systems. *Mol*  
259 *Biol Evol*. 2013;30:315–31.
- 260 17. Guglielmini J, Quintais L, Garcillán-Barcia MP, Cruz F de la, Rocha EPC. The repertoire of ICE in  
261 prokaryotes underscores the unity, diversity, and ubiquity of conjugation. *PLoS Genet*. 2011;7:e1002222.
- 262 18. Shterzer N, Mizrahi I. The animal gut as a melting pot for horizontal gene transfer. *Can J Microbiol*.  
263 2015;61:603–5.
- 264 19. Kaufman JH, Terrizzano I, Nayar G, Seabolt E, Agarwal A, Slizovskiy IB, et al. Integrative and con-  
265 jugative elements (ice) and associated cargo genes within and across hundreds of bacterial genera. *bioRxiv*.  
266 2020. doi:10.1101/2020.04.07.030320.
- 267 20. Arredondo-Alonso S, Willems RJ, Schaik W van, Schürch AC. On the (im)possibility of reconstructing  
268 plasmids from whole-genome short-read sequencing data. *Microb Genom*. 2017;3:e000128.
- 269 21. Almeida A, Mitchell AL, Boland M, Forster SC, Gloor GB, Tarkowska A, et al. A new genomic blueprint  
270 of the human gut microbiota. *Nature*. 2019;568:499–504.
- 271 22. Krawczyk PS, Lipinski L, Dziembowski A. PlasFlow: Predicting plasmid sequences in metagenomic data  
272 using genome signatures. *Nucleic Acids Res*. 2018;46:e35.
- 273 23. Maguire F, Jia B, Gray KL, Lau WYV, Beiko RG, Brinkman FSL. Metagenome-assembled genome  
274 binning methods with short reads disproportionately fail for plasmids and genomic islands. *Microb Genom*.

275 2020;6.

276 24. Suzek BE, Wang Y, Huang H, McGarvey PB, Wu CH, UniProt Consortium. UniRef clusters: A compre-  
277 hensive and scalable alternative for improving sequence similarity searches. *Bioinformatics*. 2015;31:926–32.

278 25. Hall AB, Yassour M, Sauk J, Garner A, Jiang X, Arthur T, et al. A novel ruminococcus gnavus clade  
279 enriched in inflammatory bowel disease patients. *Genome Med*. 2017;9:103.

280 26. Gibson MK, Wang B, Ahmadi S, Burnham C-AD, Tarr PI, Warner BB, et al. Developmental dynamics  
281 of the preterm infant gut microbiota and antibiotic resistome. *Nat Microbiol*. 2016;1:16024.

282 27. Bushnell B, Rood J, Singer E. BBMerge - accurate paired shotgun read merging via overlap. *PLoS One*.  
283 2017;12:e0185056.

284 28. Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for illumina sequence data. *Bioinform-*  
285 *atics*. 2014;30:2114–20.

286 29. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. MetaSPAdes: A new versatile metagenomic assem-  
287 bler. *Genome Res*. 2017;27:824–34.

288 30. Eren AM, Esen ÖC, Quince C, Vineis JH, Morrison HG, Sogin ML, et al. Anvi'o: An advanced analysis  
289 and visualization platform for 'omics data. *PeerJ*. 2015;3:e1319.

290 31. Eddy SR. Accelerated profile HMM searches. *PLoS Comput Biol*. 2011;7:e1002195.

291 32. Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: Prokaryotic gene  
292 recognition and translation initiation site identification. *BMC Bioinformatics*. 2010;11:119.

293 33. Menzel P, Ng KL, Krogh A. Fast and sensitive taxonomic classification for metagenomics with Kaiju.  
294 *Nat Commun*. 2016;7:11257.

295 34. Robertson J, Bessonov K, Schonfeld J, Nash JHE. Universal whole-sequence-based plasmid typing and  
296 its utility to prediction of host range and epidemiological surveillance. *Microb Genom*. 2020;6.

297 35. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat Methods*.  
298 2015;12:59–60.

299 36. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*. 2012;9:357–9.

300 37. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence alignment/map format  
301 and SAMtools. *Bioinformatics*. 2009;25:2078–9.

302 38. Kang DD, Li F, Kirton E, Thomas A, Egan R, An H, et al. MetaBAT 2: An adaptive binning algorithm  
303 for robust and efficient genome reconstruction from metagenome assemblies. *PeerJ*. 2019;7:e7359.

304 39. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: Assessing the quality of  
305 microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res*. 2015;25:1043–55.

- 306 40. Moss EL, Maghini DG, Bhatt AS. Complete, closed bacterial genomes from microbiomes using nanopore  
307 sequencing. *Nature Biotechnology*. 2020. doi:10.1038/s41587-020-0422-6.
- 308 41. Giguere DJ, Bahcheli AT, Joris BR, Paulssen JM, Gieg LM, Flatley MW, et al. Complete and validated  
309 genomes from a metagenome. *bioRxiv*. 2020. doi:10.1101/2020.04.08.032540.
- 310 42. Beaulaurier J, Zhu S, Deikus G, Mogno I, Zhang X-S, Davis-Richardson A, et al. Metagenomic binning  
311 and association of plasmids with bacterial host genomes using dna methylation. *Nat Biotechnol*. 2018;36:61–  
312 9.

SRR3131692	pre-term infant
SRR3131694	pre-term infant
SRR3131696	pre-term infant
SRR3131698	pre-term infant
SRR3131700	pre-term infant
SRR3131702	pre-term infant
SRR3131704	pre-term infant
SRR3131706	pre-term infant
SRR3131708	pre-term infant
SRR3131710	pre-term infant
SRR3131712	pre-term infant
SRR3131714	pre-term infant
SRR3131716	pre-term infant
SRR3131718	pre-term infant
SRR3131720	pre-term infant
SRR3131722	pre-term infant
SRR3131724	pre-term infant
SRR3131726	pre-term infant
SRR3131728	pre-term infant
SRR3131730	pre-term infant
SRR3131732	pre-term infant
SRR3131734	pre-term infant
SRR3131736	pre-term infant
SRR3131738	pre-term infant
SRR3131740	pre-term infant
SRR3131742	pre-term infant
SRR3131744	pre-term infant
SRR3131746	pre-term infant
SRR3131748	pre-term infant
SRR3131750	pre-term infant
SRR3131752	pre-term infant
SRR3131754	pre-term infant
SRR3131756	pre-term infant
SRR3131758	pre-term infant
SRR3131760	pre-term infant
SRR3131762	pre-term infant
SRR3131764	pre-term infant
SRR3131766	pre-term infant
SRR3131768	pre-term infant
SRR3131770	pre-term infant
SRR3131772	pre-term infant
SRR3131774	pre-term infant
SRR3131776	pre-term infant
SRR3131778	pre-term infant
SRR3131780	pre-term infant
SRR3131782	pre-term infant
SRR3131784	pre-term infant
SRR3131786	pre-term infant
SRR3131788	pre-term infant
SRR3131790	pre-term infant
SRR3131812	pre-term infant
SRR5650110	general
SRR5650111	general
SRR5650112	general

SRR5650113	general
SRR5650114	general
SRR5650115	general
SRR5650116	general
SRR5650117	general
SRR5650118	general
SRR5650119	general
SRR5650120	general
SRR5650121	general
SRR5650122	general
SRR5650123	general
SRR5650124	general
SRR5650125	general
SRR5650126	general
SRR5650127	general
SRR5650128	general
SRR5650129	general
SRR5650130	general
SRR5650131	general
SRR5650132	general
SRR5650133	general
SRR5650134	general
SRR5650135	general
SRR5650136	general
SRR5650137	general
SRR5650138	general
SRR5650139	general
SRR5650140	general
SRR5650141	general
SRR5650142	general
SRR5650143	general
SRR5650144	general
SRR5650145	general
SRR5650146	general
SRR5650147	general
SRR5650148	general
SRR5650149	general
SRR5650150	general
SRR5650151	general
SRR5650152	general
SRR5650153	general
SRR5650154	general
SRR5650155	general
SRR5650156	general
SRR5650157	general
SRR5650158	general
SRR5650159	general