

# Comparative genomics reveals plasmid-encoded genes associated with *Klebsiella* infection in colonized patients

**Jay Vornhagen**

University of Michigan Medical School

**Emily Roberts**

University of Michigan Medical School

**Lavinia Unverdorben**

University of Michigan Medical School

**Sophia Mason**

University of Michigan Medical School

**Alieysa Patel**

University of Michigan Medical School

**Ryan Crawford**

University of Michigan Medical School

**Caitlyn Holmes**

University of Michigan Medical School

**Yuang Sun**

University of Michigan Medical School

**Alexandra Teodorescu**

University of Michigan Medical School

**Evan Snitkin**

University of Michigan Medical School

**Lili Zhao**

University of Michigan

**Patricia Simner**

Johns Hopkins University School of Medicine

**Pranita Tamma**

Johns Hopkins University School of Medicine

**Krsihna Rao**

University of Michigan

**Keith Kaye**

Robert Wood Johnson Medical School

**Michael Bachman** (✉ [mikebach@med.umich.edu](mailto:mikebach@med.umich.edu))

**Article**

**Keywords:**

**Posted Date:** April 4th, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-1254643/v1>

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

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Nature Communications on August 1st, 2022. See the published version at <https://doi.org/10.1038/s41467-022-31990-1>.

1 **Title**

2 Comparative genomics reveals plasmid-encoded genes associated with *Klebsiella* infection in  
3 colonized patients  
4

5 **Authors**

6 Jay Vornhagen<sup>1,2</sup>, Emily K. Roberts<sup>3</sup>, Lavinia Unverdorben<sup>2</sup>, Sophia Mason<sup>1</sup>, Alieysa Patel<sup>1</sup>, Ryan  
7 Crawford<sup>4</sup>, Caitlyn L. Holmes<sup>1,2</sup>, Yuang Sun<sup>1</sup>, Alexandra Teodorescu<sup>1</sup>, Evan S. Snitkin<sup>2</sup>, Lili Zhao<sup>3</sup>,  
8 Patricia J. Simner<sup>5</sup>, Pranita D. Tamma<sup>6</sup>, Krishna Rao<sup>7</sup>, Keith S. Kaye<sup>7</sup>, Michael A. Bachman<sup>1,2\*</sup>  
9

10 **Affiliations**

11 <sup>1</sup>Department of Pathology, Michigan Medicine, University of Michigan, Ann Arbor, United States of  
12 America

13 <sup>2</sup>Department of Microbiology & Immunology, Michigan Medicine, University of Michigan, Ann Arbor,  
14 United States of America

15 <sup>3</sup>Department of Biostatistics, School of Public Health, University of Michigan, Ann Arbor, United States  
16 of America

17 <sup>4</sup>Department of Computational Medicine and Bioinformatics, Michigan Medicine, University of Michigan,  
18 Ann Arbor, United States of America

19 <sup>5</sup>Division of Medical Microbiology, Department of Pathology, Johns Hopkins University School of  
20 Medicine, Baltimore, United States of America

21 <sup>6</sup>Department of Pediatrics, Johns Hopkins University School of Medicine, Baltimore, United States of  
22 America

23 <sup>7</sup>Department of Internal Medicine/Infectious Diseases Division, Michigan Medicine, University of  
24 Michigan, Ann Arbor, United States of America

25 \*Corresponding author. Email: mikebach@med.umich.edu

26 **Abstract**

27 *Klebsiella* frequently colonizes the gut and colonization is associated with subsequent infection. To  
28 identify genes associated with progression from colonization to infection, we undertook a case-control  
29 comparative genomics study. Concordant cases (N = 85), where colonizing and invasive *Klebsiella*  
30 isolates were identical strain types, were matched to asymptotically colonizing controls (N = 160).  
31 Thirty-seven genes were significantly associated with infection, 27 of which remained significant  
32 following adjustment for patient variables and bacterial phylogeny. Infection-associated genes were not  
33 previously characterized virulence factors, but instead a diverse group of stress resistance, regulatory  
34 and antibiotic resistance genes, despite careful adjustment for antibiotic exposure. Many infection-  
35 associated genes were plasmid borne, and for some, the relationship with infection was mediated by  
36 gut dominance. Five infection-associated genes were validated in a geographically independent cohort  
37 of colonized patients. This study identifies novel genes that are reproducibly associated with  
38 progression to infection in patients colonized by diverse *Klebsiella*.

39

## 40 Introduction

41 The *Klebsiella pneumoniae* species complex (hereby referred to as *Klebsiella*) is comprised of several  
42 species capable of causing severe infections, including bacteremia, pneumonia, and urinary tract  
43 infection (UTI). The member species are *K. pneumoniae*, *K. variicola*, *K. quasipneumoniae*, *K.*  
44 *quasivariicola*, and *K. africana*, which are genetically distinct but often clinically indistinguishable  
45 (reviewed in<sup>1</sup>). *Klebsiella* infections are a serious public health concern because they are a leading  
46 cause of healthcare-associated infections<sup>2</sup>, can harbor multiple antimicrobial resistance (AMR)  
47 determinants, and can be hypervirulent. Antimicrobial resistant *Klebsiella* complicate disease  
48 treatment, leading to high mortality and healthcare costs whereas hypervirulent *Klebsiella* cause unique  
49 clinical manifestations such as pyogenic liver abscess and meningitis<sup>1,3,4</sup>. *Klebsiella* colonization is the  
50 primary risk factor for disease<sup>5,6</sup>, where the gut acts as a reservoir for disease-causing *Klebsiella*  
51 strains<sup>5,7</sup>. Patient variables, including comorbidities and baseline laboratory values, partially explain  
52 which patients will progress to infection<sup>8</sup>. Furthermore, high *Klebsiella* gut relative abundance is  
53 associated with an increased risk of infection in colonized patients<sup>9</sup>. For hypervirulent *Klebsiella*, the  
54 genes associated with their unique presentations have been well-defined by molecular epidemiology  
55 and experimental studies (reviewed in<sup>10</sup>); however, little is known about the *Klebsiella* genes contained  
56 by non-hypervirulent strains, including multi-drug resistant (MDR) strains, that increase the risk of  
57 disease in colonized patients more broadly. Given that most infections are caused by non-hypervirulent  
58 strains<sup>11</sup>, more data concerning the genetic determinants of infection are necessary for predicting,  
59 diagnosing, and experimentally evaluating these strains.

60  
61 In general, the ability of bacteria to progress from colonization to infection is dependent on the immune  
62 status of the host, environmental and iatrogenic exposures such as antimicrobials, and the virulence  
63 potential of the strain. Known *Klebsiella* virulence determinants are found in both the core genome,  
64 which are ~1,700 genes present in >95% of all sequenced strains and the accessory genome, which  
65 are genes that vary among sequenced strains<sup>12</sup>. The overall *Klebsiella* accessory genome is currently

66 estimated to include >100,000 unique genes<sup>1</sup>. With an average of >5,000 genes per genome,  
67 thousands of genes vary between strains. The size of the *Klebsiella* accessory genome provides a  
68 unique opportunity to perform comparative genomics studies to understand pathogenesis. Comparative  
69 genomics approaches are used to track hospital outbreaks<sup>13</sup>, identify in-host adaptation<sup>14</sup>, and infer  
70 conserved virulence determinants or pathways. The latter approach has been used to interrogate the  
71 pathogenicity of numerous bacterial pathogens, such as *Pseudomonas aeruginosa*<sup>15</sup>, *Streptococcus*  
72 *pneumoniae*<sup>16</sup>, *Escherichia coli*<sup>17</sup>, and invasive non-typhoidal *Salmonella enterica* spp.<sup>18</sup>. Previously,  
73 we used a comparative genomics approach to identify *Klebsiella* factors associated with infection in a  
74 small, hospital-wide cohort of *Klebsiella*-colonized patients<sup>19</sup>. Two infection-associated factors  
75 identified in this study, the psicose utilization locus and *ter* operon, were experimentally validated and  
76 characterized<sup>19,20</sup>, confirming the value of this approach for understanding *Klebsiella* pathogenesis.

77  
78 The strong association between gut colonization and *Klebsiella* disease indicates a potential point of  
79 intervention to prevent infection in the most at-risk patients, such as those in intensive care and  
80 hematology/oncology units. If patients with risk factors for infection are colonized by high-risk *Klebsiella*  
81 isolates, early detection and intervention may prevent disease. The development of these approaches  
82 is likely to be critical as therapeutic options become limited due to rising antimicrobial resistance. We  
83 previously performed a large cohort study of *Klebsiella*-colonized patients in the intensive care and  
84 hematology/oncology units and identified patient variables associated with subsequent infection<sup>8</sup>.  
85 Additionally, we identified an association between gut dominance by *Klebsiella* and infection, even after  
86 adjustment for the patient variables we identified<sup>9</sup>. Here, we aimed to rigorously identify *Klebsiella*  
87 genes that are associated with disease in colonized patients through whole-genome sequencing (WGS)  
88 in a nested case-control study from a cohort of over 1,900 colonized patients. This approach differed  
89 from other comparative genomics studies in that it focused on concordant infections (e.g., same  
90 colonizing and infecting isolate) after infection, directly compared cases of infection to well-matched  
91 controls from the same patient population and incorporated careful adjustment with patient variables

92 associated with infection. Through this approach, we identified 27 genes associated with infection in  
93 colonized patients. We identified which genes may act through gut dominance, and determined that  
94 many genes were found on large, conjugative plasmids. Finally, we demonstrate that five of these  
95 genes are associated with infection in a geographically independent cohort of intensive care patients  
96 colonized with *Klebsiella*.

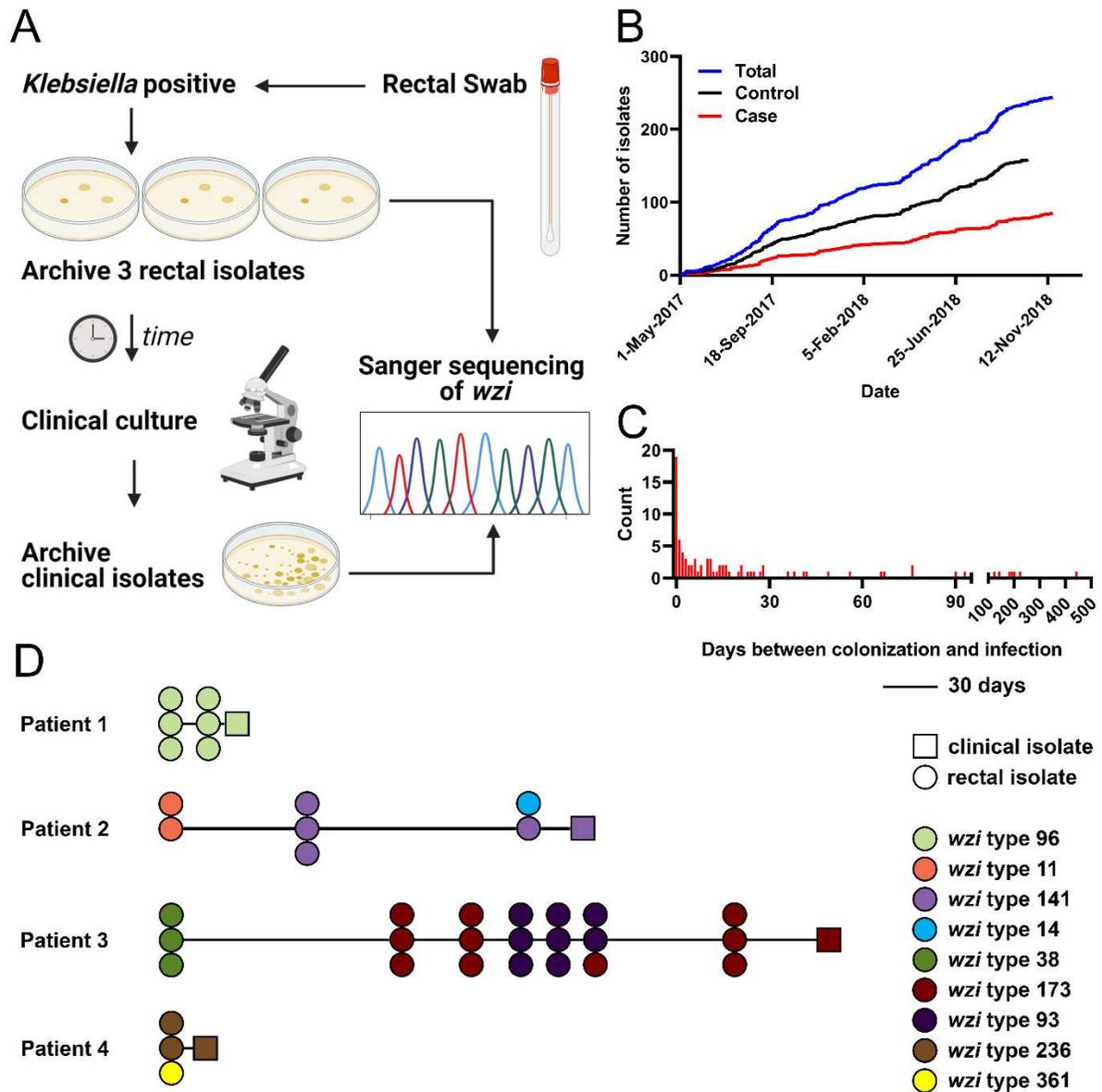
97

## 98 **Results**

### 99 **Selection of *Klebsiella* for WGS and description of the case-control study**

100 A nested case-control study was performed using a cohort of 1,978 *Klebsiella* colonized, intensive care,  
101 and hematology/oncology patients at a single academic medical center in Michigan, USA<sup>8</sup>. Over the  
102 course of this study patient rectal swabs were screened for the presence of *Klebsiella*. If a patient was  
103 colonized with *Klebsiella*, up to three rectal isolates were banked, and the patient was enrolled in the  
104 study. Subsequent rectal swab isolates and clinical culture *Klebsiella* isolates (blood, respiratory,  
105 urinary, and other) were also banked (Fig. 1A). Cases of infection were considered for inclusion based  
106 on chart review using strict clinical definitions<sup>8</sup>. To facilitate a rigorous comparative genomics approach,  
107 Sanger sequencing of the *wzi* gene<sup>21</sup> was used to screen for concordance between the clinical culture  
108 isolate and a preceding colonizing rectal isolate from the same patient (Fig. 1A). Cases and controls  
109 accumulated at a similar rate throughout the enrollment period (Fig. 1B). The median time between  
110 detection of *Klebsiella* colonization and subsequent infection was 10 days (range = 0-443 days, Fig.  
111 1C). In some instances, the progression from colonization to infection was simple, wherein a single *wzi*  
112 type was detected in all colonizing *Klebsiella* isolates and the subsequent infecting isolate (Fig. 1D,  
113 “Patient 1”). In other instances, the infection trajectory was more complicated, wherein some patients  
114 exhibited single or multiple strain replacements in the gut prior to infection (Fig. 1D, “Patient 2” and  
115 “Patient 3”), or mixed colonization where only one colonizing isolate caused the subsequent infection  
116 (Fig. 1D, “Patient 2” and “Patient 4”). Collectively, these data indicate concordant infections are  
117 detectable in colonized patients in the context of highly dynamic *Klebsiella* populations.

118



**Fig. 1. Description of *Klebsiella* isolate collection.**

(A) Schematic representation of *Klebsiella* isolate collection and *wzi*-based concordance testing for WGS selection. (B) Rectal isolate selected for WGS collection dates stratified by case status. (C) Histogram of the number of days between identification of a rectal *Klebsiella* isolate concordant with the clinical isolate from the same patient. (D) Trace diagrams of *wzi* typing in *Klebsiella* colonizing and infecting isolates in infected patients.

In total, 85 cases of infection with an identical *wzi* type to a colonizing rectal isolate (cases) were selected for WGS, including 39 bloodstream infections, 23 urinary tract infections, 19 pneumonias, and

129 4 other sites of infection (including two cases of cholecystitis, one case of empyema, and once case of  
130 ecthyma). Based on antibiotic susceptibility data from medical records, 44 of 85 (51.8%) of the clinical  
131 isolates from cases were pan-sensitive apart from intrinsic ampicillin resistance. The most common  
132 AMR phenotypes were resistance to first- and second-generation cephalosporins, the combinatorial  $\beta$ -  
133 lactam/ $\beta$ -lactamase-inhibitors ampicillin-sulbactam and piperacillin-tazobactam, and trimethoprim-  
134 sulfamethoxazole (Supplementary Fig. S1). Cases were matched to two colonized controls (who  
135 remained asymptomatic) based on sex, age, and date of rectal swab collection. In addition, controls  
136 had the same type of clinical culture collected as their corresponding case, but with no cultivatable  
137 bacteria. Two strata included only a single control because of the inability to find an appropriate match,  
138 and an additional 8 controls were removed from the study following WGS based on species  
139 classification outside of *Klebsiella*. In the final case-control design, 85 cases, including the clinical and  
140 rectal isolates, and 160 control rectal isolates were evaluated (Table 1).

141  
142 **Table 1. Patient demographics and select baseline characteristics**

Patient variable		Case (N = 85)	Control (N = 160)	Odds ratio (OR)	95% CI	P value
Sex*	Female	40 (47.1%)	76 (47.5%)			
	Male	45 (52.9%)	81 (50.6%)			
	Missing	0 (0.0%)	3 (1.88%)			
Median age (range)		60 (83-25)	59 (85-27)	1.02	0.963 to 1.08	0.527
Infection site	Blood	39 (45.9%)				
	Respiratory	19 (22.4%)				
	Urine	23 (27.1%)				
	Other	4 (4.71%)				
Mean weighted Elixhauser score (±SD)		22.2±11.6	19.7±11.9	1.02	0.994 to 1.04	0.127
Depression	Yes	30 (35.3%)	44 (27.5%)	1.31	0.761 to 2.26	0.327
	No	55 (64.7%)	113 (70.6%)			
	Missing	0 (0.0%)	3 (1.88%)			
Prior diuretics	Yes	30 (35.3%)	38 (23.8%)	1.69	0.905 to 3.05	0.101
	No	55 (64.7%)	119 (74.4%)			
	Missing	0 (0.0%)	3 (1.88%)			
Prior vitamin D	Yes	18 (21.2%)	19 (11.8%)	1.66	0.972 to 4.44	0.0591
	No	67 (78.8%)	138 (86.3%)			
	Missing	0 (0.0%)	3 (1.88%)			
Prior pressors/inotrope	Yes	20 (23.5%)	14 (8.75%)	3.06	1.40 to 6.68	<b>0.00502</b>
	No	65 (76.6%)	143 (89.4%)			
	Missing	0 (0.0%)	3 (1.88%)			
High-risk antibiotics	Yes	30 (35.3%)	32 (20.0%)	2.21	1.18 to 4.14	<b>0.013</b>
	No	55 (64.7%)	125 (78.1%)			
	Missing	0 (0.0%)	3 (1.88%)			
Albumin	< 2.5 g/dL	35 (41.2%)	34 (21.3%)	2.86	1.49 to 5.48	<b>0.00157</b>
	≥ 2.5 g/dL	47 (55.3%)	116 (72.5%)			
	Missing	3 (3.53%)	10 (6.25%)			

\*OR for sex was incalculable due to use of sex as a matching variable.

### Description of sequenced *Klebsiella* isolates

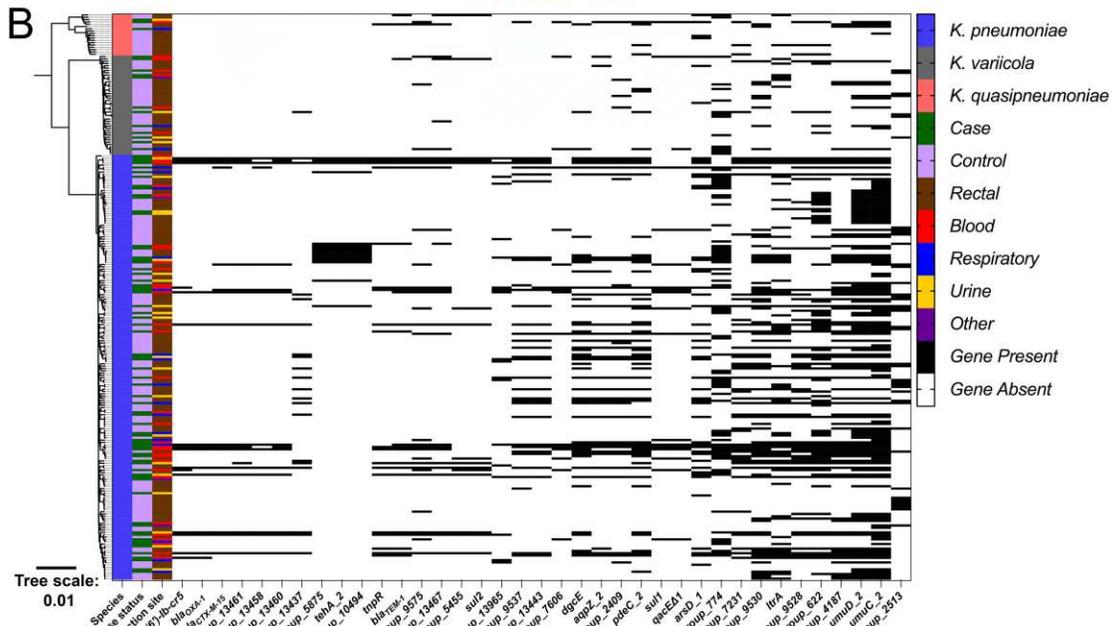
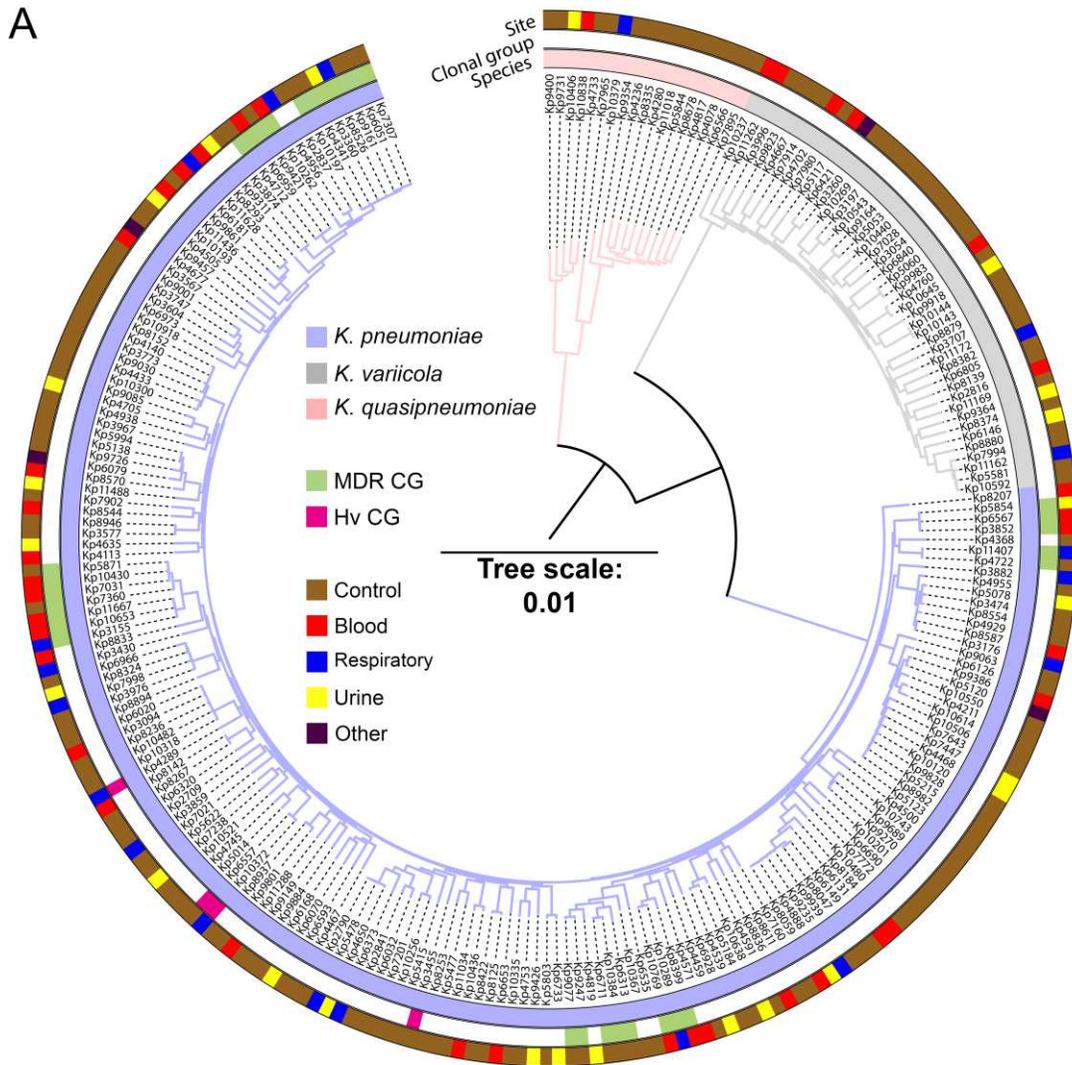
Following the determination of isolates for WGS, genomic DNA was extracted from all *Klebsiella* isolates and subjected to Illumina sequencing (Supplementary Table 1). Assembled genomes were generally of high quality, wherein the median number of contigs was 110.5 (range = 43-876 contigs) and the median N50 was 459,591 bp (range = 147,972-3,285,432 bp). All genomes were of expected length, wherein the median genome size was 5.50 Kb (range = 5.06-6.26 Kb). Sequencing quality was similar between rectal case and control isolates, which were used for most subsequent analyses (Supplementary Fig. S2). As expected, species assignment based on WGS included the three most

154 common members of the *Klebsiella pneumoniae* species complex (*K. pneumoniae*, *K. variicola*, *K.*  
155 *quasipneumoniae*). *K. pneumoniae*, *K. variicola*, and *K. quasipneumoniae* were evenly distributed  
156 amongst the four sites of infection and control rectal swabs (Supplementary Fig. S3A). Case and control  
157 rectal isolates within a stratum may be the same or different species, as reflected in either high (100-  
158 80%) or low (60-30%) overall nucleotide identity (Supplementary Fig. S3B). For the rectal and clinical  
159 isolate of a given case of infection, matching *wzi* types were highly predictive of overall nucleotide  
160 identity (median % nucleotide identity = 99.65%, range = 95.04-99.88%, Supplementary Fig. S3C). The  
161 time between detection of *Klebsiella* colonization and subsequent infection was negatively correlated  
162 with sequence similarity; however, the effect was minimal ( $r^2 = 0.06$ , Supplementary Fig. S3D).  
163 Interestingly, instances of lower nucleotide identity similarity (< 98%) between the rectal and clinical  
164 isolates from concordant cases of infection appeared to be due to large genomic losses or additions  
165 (Supplementary Fig. S3E), which may be driven by plasmids or other mobile genetic elements (MGEs).  
166 Indeed, these case clinical and rectal isolate pairs displayed differences in the number of annotations  
167 that corresponded to differences in genome size (Supplementary Fig. S3F).

168  
169 To confirm concordance between case clinical and rectal isolates, *wzi* strain typing was performed from  
170 the assembled genomes. All *wzi* types matched results from Sanger sequencing, and construction of  
171 a *wzi* approximately-maximum-likelihood tree revealed a total of 147 unique *wzi* types identified  
172 (Supplementary Fig. S4). Then, multi-locus sequence typing was performed across all isolates. One-  
173 hundred thirty-one unique sequence types (STs) were identified, and all case rectal and case clinical  
174 STs were identical. Rarefaction of STs indicated that cases and controls were equally rich  
175 (Supplementary Fig. S5). ST253 was the most abundant ST but still low frequency (10/245 rectal  
176 isolates, 4.08%), and 44 rectal isolates were non-typable (44/245, 17.96%). Overall, hypervirulent  
177 clonal groups (CGs)<sup>22</sup> were rare, comprising 2.35% of cases and 1.35% of controls. Only two isolates  
178 of hypervirulent ST23 were identified, and only one caused an infection. MDR CGs<sup>22</sup>, frequently  
179 containing extended-spectrum  $\beta$ -lactamases, carbapenemases, or both, were more common,

180 representing 18.8% of cases and 9.38% of controls. These included 1 ST11, 4 ST14, 3 ST15, 1 ST16,  
181 3 ST17, 5 ST37, 2 ST147, 1 ST152, 1 ST258, and 3 ST307. 85.7% (210/245) of the isolates in our  
182 collection were from non-hypervirulent, non-MDR CGs. Combined, MDR CGs were significantly  
183 associated with case status (odds ratio = 2.24, 95% C.I. 1.03 to 4.93, Fisher  $P = 0.043$ ). Construction  
184 of a core genome phylogeny demonstrated the presence of three *Klebsiella* species within our cohort  
185 and multiple lineages across colonizing and invasive isolates of each species (Supplementary Table 1;  
186 Fig. 2A). Together, these data indicate that sequenced *Klebsiella* isolates were highly diverse, without  
187 a dominant hypervirulent or MDR CG within the population.

188



189

190

191

**Fig. 2. *Klebsiella* rectal isolates are highly diverse and infection-associated genes span that diversity.**

192 (A) Core genome phylogeny tree of sequenced *Klebsiella* rectal isolates. The inner ring indicates the  
193 species, the middle ring indicates a clonal group of concern (hypervirulent [Hv] CG: clonal groups 23,  
194 25, 380; MDR CG: 14/15, 17/20, 37, 147, 258, 307), and the outer ring indicates site of infection for  
195 concordant clinical isolate. (B) Heatmap of infection-associated gene presence/absence organized  
196 vertically by core genome phylogeny and horizontally by the strength of association (Odds ratio)  
197 between gene presence and infection following adjustment for patient variables.

198  
199 To enable the identification of infection-associated genes in rectal samples, we assembled a  
200 pangenome (Supplementary Table 2). The rectal isolate pangenome consisted of 26,089 total  
201 annotations. The core genome (genes present in  $\geq 95\%$  genomes) was comprised of 3,921 genes, and  
202 the accessory genome (genes present in  $< 95\%$  genomes) consisted of 22,168 genes (Supplementary  
203 Fig. S6A). To assess the similarity of the accessory genome between *Klebsiella* isolates in our  
204 collection, pairwise Jaccard distances were calculated after excluding the core genome and  
205 summarized as a mean for each isolate (244 pairwise comparisons per isolate). Interestingly, the mean  
206 pairwise Jaccard distance for all rectal isolates was  $0.73 \pm 0.04$ , indicating that the accessory genome  
207 of each isolate was highly dissimilar from other isolates. No differences were observed between mean  
208 pairwise Jaccard distances for cases and controls (Supplementary Fig. S6B). Overall, the cohort of  
209 *Klebsiella* isolates exhibited high pan-genome diversity.

### 211 **Genes associated with *Klebsiella* infection**

212 To identify genes in the accessory genome associated with infection we used Scoary, collapsing linked  
213 genes that displayed identical presence/absence patterns across all sequenced isolates (e.g.,  
214 organized in operons, linked to MGEs, etc.) into units and comparing unit frequencies between cases  
215 and controls (Supplementary Table 3)<sup>23</sup>. Thirty-seven genes across 34 units were significantly  
216 associated with infection (Q value  $\leq 0.2$ , 2,166 tests). Two units contained genes with identical  
217 presence/absence patterns (*group\_5875-tehA\_2-group\_10494* and *group\_5455-sul2*); the remaining

218 units contained a single gene. Overall, 36 genes were associated with cases, one was associated with  
219 controls (Table 2).

220  
221 To control for patient variables in our comparative genomics analysis, we leveraged a previously  
222 constructed multivariate model for infection risk in the larger cohort from which this nested case-control  
223 study was derived. This model was comprised of mean weighted Elixhauser comorbidity score,  
224 depression, prior diuretic use, vitamin D use, pressors/inotrope use, high-risk antibiotic use, and low  
225 albumin. Prior medication use was defined as greater than 48 hours but less than 90 days prior to the  
226 colonizing swab and albumin was measured within 48 hours of the colonizing swab<sup>8</sup>. High-risk  
227 antibiotics, based on their impact on indigenous gut microbiota, were defined as  $\beta$ -lactam/ $\beta$ -lactamase  
228 inhibitor combinations, carbapenems, third- and fourth-generation cephalosporins, fluoroquinolones,  
229 clindamycin, and oral vancomycin<sup>24</sup>. Consistent with the explanatory value of these variables in the  
230 overall cohort, the variables used for adjustment trended towards or were significantly associated with  
231 infection in the nested case-control study (Table 1). We then used inverse probability of treatment  
232 weighting to adjust for these variables in measuring the association of each gene with the outcome of  
233 infection.

234  
235 Following adjustment for clinical variables, 28 of 37 genes remained significantly associated with  
236 infection ( $P$  value  $\leq 0.01$ ). These 28 infection-associated genes fell into four categories. First, many  
237 infection-associated genes encode AMR proteins even after correction for exposure to high-risk  
238 antibiotics. These included *aac(6')-Ib-cr5*, encoding an aminoglycoside acetyltransferase that also has  
239 activity against fluoroquinolones<sup>25</sup>, and *bla<sub>OXA-1</sub>*, *bla<sub>CTX-M-15</sub>*, *bla<sub>TEM-1</sub>* encoding broad-spectrum  $\beta$ -  
240 lactamases. Second, several infection-associated genes have predicted protein functions that may alter  
241 the physiology of the bacteria directly. These include the predicted aquaporin gene *aqpZ\_2*, tellurite-  
242 resistance gene *tehA\_2*, arsenate resistance gene *arsD\_1*, and diguanylate cyclase gene *dgcE*. Third,  
243 certain MGEs and plasmid maintenance genes were associated with infection (e.g., IS1380 family

244 transposase ISEcp1, Transposon Tn3 resolvase, IS6 family transposase IS6100, *umuDC\_2*). The  
245 remainder of infection-associated genes encode hypothetical proteins. Interestingly, certain genes were  
246 highly associated with one another, suggesting the presence of linked genes that may be acquired  
247 simultaneously through horizontal gene transfer (HGT) events (Supplementary Fig. S7).

248  
249 None of the *Klebsiella* genes associated with infection encoded known virulence factors. To directly  
250 assess the association of known virulence genes with infection, we used the Kleborate genotyping tool,  
251 which is an *in silico* platform for rapid genotyping of *Klebsiella* genomes (Supplementary Table 4)<sup>22</sup>.  
252 Consistent with the paucity of hypervirulent STs, most well-characterized *Klebsiella* virulence factors  
253 were rare in our dataset, including salmochelin (*iro*), aerobactin (*iuc*), and colibactin (*clb*) synthesis loci,  
254 and hypermucoviscosity regulators RmpA/A2 (Table 3). The yersiniabactin (*ybt*) synthesis locus was  
255 more frequent overall, and in cases versus controls, but was not associated with infection (odds ratio =  
256 1.7,  $P = 0.144$ ). Correspondingly, the virulence score, which is a composite metric of the virulence loci  
257 listed above, was not significantly different between cases and controls (Case virulence score =  $0.294$   
258  $\pm 0.784$ , Control virulence score =  $0.163 \pm 0.513$ ,  $P = 0.115$ , Student's t-test). Even when assessed  
259 directly and after controlling for clinical variables, canonical virulence factors are not associated with  
260 infection in this group of colonized, intensive care patients.

**Table 2. Genes associated with patient case status.**

Gene	Annotation	Case frequency	Control frequency	Unadjusted			Adjusted for clinical variables		Mediation			
				OR	P value	Q value	OR	P value	Direct $\beta$	Direct P value	Indirect $\beta$	Indirect P value
<i>aac(6)-Ib-cr5</i>	fluoroquinolone-acetylating aminoglycoside 6'-N-acetyltransferase AAC(6)-Ib-cr5	0.176	0.006	30.68	<b>0.0011</b>	<b>0.00131</b>	52.41	<b>0.0003</b>	3.0601	<b>0.0045</b>	0.086	<b>0.01</b>
<i>bla<sub>OXA-1</sub></i>	oxacillin-hydrolyzing class D $\beta$ -lactamase OXA-1	0.153	0.006	25.54	<b>0.002</b>	<b>0.0053</b>	49.29	<b>0.0003</b>	2.9698	<b>0.0064</b>	0.077	0.03
<i>bla<sub>CTX-M-15</sub></i>	class A extended-spectrum $\beta$ -lactamase CTX-M-15	0.165	0.013	13.94	<b>0.0007</b>	<b>0.0053</b>	19.17	<b>0.0004</b>	2.3386	<b>0.004</b>	0.047	0.1567
<i>group_13461</i>	IS1380 family transposase ISEcp1	0.165	0.013	13.94	<b>0.0007</b>	<b>0.0053</b>	18.69	<b>0.0003</b>	2.482	<b>0.0022</b>	0.019	0.5833
<i>group_13458</i>	hypothetical protein	0.141	0.013	11.46	<b>0.0018</b>	<b>0.01959</b>	17.57	<b>0.0006</b>	2.0832	0.0112	0.034	0.3483
<i>group_13460</i>	hypothetical protein	0.153	0.013	12.69	<b>0.0011</b>	<b>0.01065</b>	16.73	<b>0.0005</b>	2.3118	<b>0.0047</b>	0.024	0.5217
<i>group_13437</i>	hypothetical protein	0.176	0.031	5.97	<b>0.001</b>	<b>0.05505</b>	7.28	<b>0.0014</b>	1.6137	<b>0.0054</b>	0.005	0.875
<i>group_5875-tehA_2-group_10494*</i>	HTH-type transcriptional regulatory protein GabR - tellurite resistance protein TehA - hypothetical protein	0.129	0.025	5.65	<b>0.004</b>	<b>0.1908</b>	7.09	<b>0.0029</b>	1.8251	<b>0.0053</b>	0.007	0.8383
<i>tnpR</i>	transposon Tn3 resolvase	0.176	0.031	5.97	<b>0.001</b>	<b>0.05505</b>	6.24	<b>0.0017</b>	1.6379	<b>0.0054</b>	0.023	0.45
<i>bla<sub>TEM-1</sub></i>	class A broad-spectrum $\beta$ -lactamase TEM-1	0.212	0.038	6.28	<b>0.0002</b>	<b>0.01236</b>	5.46	<b>0.0009</b>	1.6623	<b>0.0021</b>	0.035	0.1783
<i>group_9575</i>	tyrosine recombinase XerC	0.259	0.056	5.39	<b>0.0001</b>	<b>0.0053</b>	4.72	<b>0.0019</b>	1.3993	<b>0.0044</b>	0.048	0.05
<i>group_13467</i>	IS6 family transposase IS6100	0.224	0.044	5.75	<b>0.0002</b>	<b>0.01093</b>	4.58	<b>0.0026</b>	1.5338	<b>0.003</b>	0.023	0.4017
<i>group_5455-sul2*</i>	IS91 family transposase - sulfonamide-resistant dihydropteroate synthase Sul2	0.129	0.025	5.07	<b>0.0077</b>	<b>0.1908</b>	4.52	0.0258	2.0214	0.0053	-0.033	0.3617
<i>group_13965</i>	hypothetical protein	0.153	0.038	4.11	<b>0.0066</b>	<b>0.16685</b>	3.59	0.0264	1.0196	0.0806	0.074	0.01
<i>group_9537</i>	hypothetical protein	0.235	0.081	3.71	<b>0.0009</b>	<b>0.13482</b>	3.3	<b>0.0035</b>	1.0141	0.0216	0.049	0.025
<i>group_13443</i>	hypothetical protein	0.224	0.075	3.81	<b>0.001</b>	<b>0.16685</b>	3.23	<b>0.005</b>	1.0762	0.0178	0.06	<b>0.0083</b>
<i>group_7606</i>	hypothetical protein	0.141	0.025	6.26	<b>0.0021</b>	<b>0.10991</b>	3.07	0.0916	1.106	0.0975	0.097	0.0001
<i>dgcE</i>	hypothetical protein	0.341	0.138	3.37	<b>0.0003</b>	<b>0.07663</b>	3.06	<b>0.0016</b>	0.9197	0.0133	0.048	<b>0.0083</b>
<i>aqpZ_2</i>	aquaporin Z	0.247	0.094	3.96	<b>0.0005</b>	<b>0.16685</b>	2.97	<b>0.0085</b>	0.923	0.0385	0.074	<b>0.0001</b>
<i>group_2409</i>	ISL3 family transposase ISEc38	0.259	0.106	3.07	<b>0.002</b>	<b>0.1908</b>	2.67	<b>0.0083</b>	0.803	0.0487	0.05	<b>0.01</b>
<i>pdeC_2</i>	putative cyclic di-GMP phosphodiesterase PdeC	0.365	0.169	2.88	<b>0.0008</b>	<b>0.10991</b>	2.66	<b>0.0026</b>	0.7709	0.0296	0.041	0.015
<i>sul1</i>	sulfonamide-resistant dihydropteroate synthase Sul1	0.141	0.031	4.97	<b>0.0037</b>	<b>0.17478</b>	2.64	0.0939	1.102	0.0994	0.12	0.0001
<i>qacE<math>\Delta</math>1</i>	quaternary ammonium compound efflux SMR transporter QacE delta 1	0.153	0.031	5.46	<b>0.0019</b>	<b>0.12853</b>	2.63	0.1011	1.1732	0.0745	0.107	0.0001
<i>arsD_1</i>	arsenical resistance operon trans-acting repressor ArsD	0.318	0.144	3.19	<b>0.0006</b>	<b>0.17478</b>	2.62	<b>0.0043</b>	0.7263	0.0563	0.068	<b>0.0001</b>
<i>group_774</i>	hypothetical protein	0.459	0.25	2.55	<b>0.0013</b>	<b>0.1391</b>	2.55	<b>0.0016</b>	0.8601	<b>0.0088</b>	0.042	<b>0.0033</b>
<i>group_7231</i>	hypothetical protein	0.306	0.138	2.69	<b>0.0032</b>	<b>0.16685</b>	2.54	<b>0.0071</b>	0.7807	0.0373	0.026	0.1733
<i>group_9530</i>	hypothetical protein	0.459	0.238	2.64	<b>0.0009</b>	<b>0.08496</b>	2.38	<b>0.0037</b>	0.8449	<b>0.0094</b>	0.015	0.3533
<i>ltrA</i>	hypothetical protein	0.424	0.225	2.57	<b>0.0013</b>	<b>0.16683</b>	2.33	<b>0.0051</b>	0.7495	0.0229	0.037	0.015
<i>group_9528</i>	hypothetical protein	0.353	0.163	2.87	<b>0.001</b>	<b>0.12853</b>	2.31	0.0134	0.7418	0.0384	0.032	0.0917
<i>group_622</i>	hypothetical protein	0.482	0.269	2.71	<b>0.0006</b>	<b>0.12566</b>	2.24	<b>0.0063</b>	0.9086	<b>0.0049</b>	0.015	0.3683
<i>group_4187</i>	Antirestriction protein KlcA	0.424	0.219	2.54	<b>0.0017</b>	<b>0.12566</b>	2.24	<b>0.0092</b>	0.7747	0.0191	0.021	0.1983
<i>umuD_2</i>	Protein UmuD	0.6	0.369	2.42	<b>0.0016</b>	<b>0.10991</b>	2.1	0.0114	0.8272	0.0098	0.006	0.7
<i>umuC_2</i>	Protein UmuC	0.682	0.456	2.52	<b>0.0012</b>	<b>0.10991</b>	2.09	0.0141	0.8239	0.0107	0.006	0.71
<i>group_2513</i>	hypothetical protein	0.071	0.225	0.26	<b>0.0038</b>	<b>0.16685</b>	0.25	<b>0.0031</b>	-1.2577	0.0119	-0.02	0.3617

\*Loci have identical presence/absence patterns across all sequenced isolates and are clustered for analysis.

267 Next, we aimed to adjust for the effects of phylogeny on the association between each group of genes  
268 and infection. We used treeWAS to test for the confounding potential of clonality and recombination in  
269 the population structure by simulating a null genetic data set across the phylogenetic tree<sup>26</sup>. The output  
270 of treeWAS is three complimentary scores: The Terminal score counts the four possible combinations  
271 of genotype and phenotype without regard to phylogeny, whereas the Simultaneous and Subsequent  
272 scores account for population structure<sup>26</sup> in measuring the association between genotype and  
273 phenotype. Of the 28 infection-associated genes after adjustment (Table 2), 23 genes were also  
274 significant ( $P \leq 0.01$ ) by the Terminal score; *group\_774*, *group\_9530*, *ltrA*, *group\_622*, *group\_4187*,  
275 and *group\_2513* were not (Supplementary Fig. S8, Supplementary Table 5). Other than *group\_2513*,  
276 27 genes were significantly associated with infection by either the Simultaneous or Subsequent score,  
277 or both (Supplementary Fig. S8). For subsequent analyses, we focused on these 27 genes associated  
278 with infection both after adjustment for clinical variables and after at least one test of independence  
279 from bacterial phylogeny.

### 281 **AMR genes are associated with infection independent of prior exposure to antibiotics**

282 The identification of multiple AMR genes associated with infection after adjustment for exposure to  
283 high-risk antibiotics was unexpected. To investigate the association between AMR genes and infection  
284 further, we used two independent approaches. First, we identified classes of AMR genes as defined by  
285 Kleborate in our strain set. Overall, the prevalence of many of these AMR determinants was low (<25%).  
286 Consistent with results from Scoary, AMR determinants for aminoglycosides, fluoroquinolones,  
287 phenicols, sulfonamides, and trimethoprim were significantly associated with infection after adjustment.  
288 In addition, carbapenamase genes were also significantly associated with infection, and extended-  
289 spectrum  $\beta$ -lactamases (ESBL) approached significance ( $P = 0.011$ ; Table 3). Notably, one case rectal  
290 isolate (Kp10372, Supplementary Table 4) displayed a hypervirulent (ST23, *ybt*<sup>+</sup>, *col*<sup>+</sup>, *iuc*<sup>+</sup>, *rmpA/A2*<sup>+</sup>),  
291 MDR (KPC-3) genotype. The Kleborate Resistance Score was not significantly different between cases  
292 and controls (Case resistance score =  $0.2 \pm 0.483$ , Control resistance score =  $0.163 \pm 0.418$ ,  $P = 0.528$ ,

Student's t-test); however, this score is heavily weighted by the presence of carbapenemases in the context of colistin resistance, the latter of which was not detected in our dataset<sup>22</sup>. Next, we assessed phenotypic resistance across four antibiotic classes, using an agar dilution method with antibiotic concentrations at the Clinical & Laboratory Standards Institute breakpoint for resistance<sup>27</sup>. Screen results were confirmed for 12 randomly selected *Klebsiella* isolates using the broth microdilution method. Overall, phenotypic AMR rates were low (< 25%). Phenotypic resistance to chloramphenicol and gentamicin did not correlate well to genotypic resistance and was not associated with infection (Supplementary Table 6). Conversely, phenotypic resistance to ciprofloxacin, ceftazidime, cefuroxime, and cefepime was associated with infection following adjustment for clinical variables (Supplementary Table 6).

**Table 3. *Klebsiella* genotype association with infection.**

Genotype	Case frequency	Control frequency	Unadjusted		Adjusted for clinical variables	
			OR	P value	OR	P value
<b>Virulence determinants</b>						
<i>ybt</i>	0.207	0.133	1.7	0.144	1.78	0.155
<i>col</i>	0.024	0.007	3.72	0.286	0.629	0.709
<i>iuc</i>	0.024	0.007	3.72	0.286	0.436	0.507
<i>iro</i>	0.012	0.007	1.84	0.668	0.476	0.602
<i>rmpA</i>	0.024	0.007	3.72	0.286	0.436	0.507
<i>rmpA2</i>	0.012	0.007	1.84	0.668	0.937	0.964
<b>AMR determinants*</b>						
AGly acquired (aminoglycosides)	0.268	0.067	5.133	<b>7.01E-05</b>	3.77	<b>0.00427</b>
Flq (fluoroquinolones)	0.220	0.040	6.75	<b>1.14E-04</b>	6.69	<b>0.000507</b>
MLS (macrolides)	0.098	0.013	8.0	<b>0.010</b>	6.41	0.0649
Phe (phenicols)	0.171	0.027	7.51	<b>0.001</b>	6.35	<b>0.00332</b>
Rif (rifampin)	0.012	0.000	N/A	N/A	N/A	N/A
Sul (sulfonamides)	0.268	0.047	7.49	<b>1.22E-05</b>	5.19	<b>0.00226</b>
Tet (tetracyclines)	0.183	0.100	2.01	0.076	1.88	0.128
Tmt (trimethoprim)	0.244	0.053	5.73	<b>8.86E-05</b>	4.28	<b>0.00285</b>
Omp mutation	0.037	0.020	1.86	0.453	1.55	0.676
Bla (β-lactamase that is "other")	0.573	0.520	1.24	0.438	1.27	0.379
Bla Carb (carbapenemase)	0.061	0.007	9.68	0.040	36.2	<b>0.00151</b>
Bla ESBL (extended-spectrum β-lactamases)	0.207	0.087	2.76	0.011	2.94	0.0116
Bla broad (broad-spectrum β-lactamases)	0.354	0.347	1.03	0.915	0.902	0.734
Bla broad inhR (broad-spectrum β-lactamases with resistance to β-lactamase)	0.110	0.107	1.03	0.942	1.05	0.922

306 \*Kleborate predicted AMR genotypes absent in our dataset are not shown.

307  
308 Given the consistent association of AMR genes and phenotypes despite adjustment for exposure to  
309 high-risk antibiotics as a group, we evaluated whether their association with infection is due to exposure  
310 to specific cognate antibiotics. For example, *aac(6')-Ib-cr5* may be associated with exposure to  
311 aminoglycosides, a class that was not included in adjustment. Interestingly, prior antibiotic exposure by  
312 specific classes was not associated with the presence of most infection-associated genes in unadjusted  
313 analysis (Supplementary Fig. S7). Next, we explicitly adjusted for each class of antibiotic, in addition to  
314 clinical variables by IPTW, and assessed the association of each gene with infection. Even after  
315 individual adjustment for 18 classes of antibiotics (no patients had a prior exposure to daptomycin or  
316 polymyxin), all AMR genes remained significantly associated with infection (Supplementary Fig. S9).  
317 These data suggest an alternative explanation exists for the association between infection and AMR  
318 genes.

### 319 320 **Gut relative abundance mediates association between certain genes and *Klebsiella* infection**

321 We previously identified that high *Klebsiella* relative abundance in the gut is associated with infection<sup>9</sup>,  
322 and that genes associated with infection can enhance fitness in the gut<sup>20</sup>. Therefore, we hypothesized  
323 that the *Klebsiella* genes associated with infection may act by enabling gut dominance that in turn  
324 increases the risk of infection. *Klebsiella* relative abundance measurements from the rectal swabs most  
325 proximal to infection for cases and matched controls were available for 233 of 245 sequenced rectal  
326 isolates<sup>9</sup>. Reanalysis of these data was consistent with that of the original study, wherein cases had  
327 significantly higher *Klebsiella* relative abundance in the gut compared to controls (Fig. 3A). Mean  
328 relative abundance was not significantly different between species (Supplementary Fig. S10A).  
329 Abundance was higher in patients exposed to high-risk antibiotics, and when stratified by antibiotics  
330 exposure the abundance in cases was significantly higher than controls (Supplementary Fig. S10B-C).  
331 Analysis of *Klebsiella* gut relative abundance based on gene presence/absence (rather than case



341 gene presence/absence (mean displayed, \* $P$  value  $\leq 0.05$ , \*\* $P$  value  $\leq 0.005$ , \*\*\* $P$  value  $\leq 0.0005$ ,  
342 \*\*\*\* $P$  value  $\leq 0.00005$ , Student's t-test). Each data point represents a single patient rectal swab.

343  
344 To test the hypothesis that infection-associated genes may act through gut dominance we undertook a  
345 mediation analysis, a statistical approach that estimates the strength of both the direct associations  
346 between each gene and the outcome of infection and indirect association through gut relative  
347 abundance. For seven gene units, the association with infection was mediated, at least in part, by gut  
348 abundance (Table 2; indirect  $P$  value, indicated in bold). The AMR gene *aac(6')-Ib-cr5* had both a strong  
349 direct association with infection and an indirect mediation by abundance ( $P$  value  $\leq 0.01$ ). The  
350 hypothetical protein *group\_774* also had both significant direct and indirect effects. The remaining five  
351 genes (hypothetical protein *group\_13443*, *dgcE*, *aqpZ\_2*, transposase *group\_2409*, and *arsD\_1*) had  
352 only significant indirect effects, indicating that their association with infection is mediated by gut relative  
353 abundance. Collectively, these findings begin to suggest biological roles for *Klebsiella* genes  
354 associated with infection, in which some genes may act through high gut relative abundance, but others  
355 likely have separate functions.

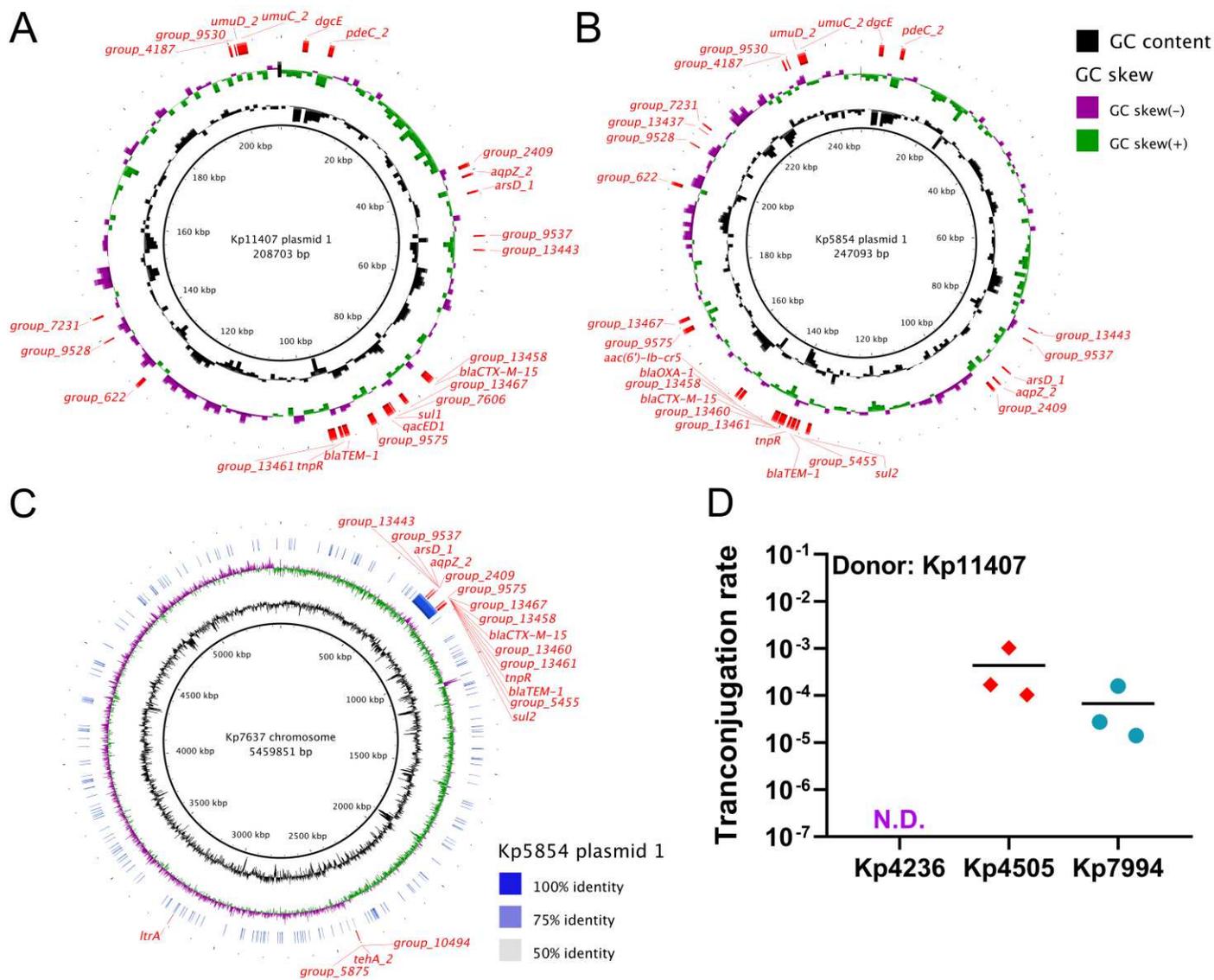
### 356 357 **Most infection-associated genes are carried on plasmids**

358 Given the diversity of functions and frequencies of the infection-associated genes, we sought to  
359 determine the genetic context of these genes as part of the accessory genome. Based on nucleotide  
360 homology, many of these genes have been previously detected on plasmids<sup>28-30</sup>. Consistent with  
361 plasmid carriage, certain infection-associated genes, including both AMR and non-AMR genes, are  
362 highly associated with one another (Supplementary Fig. S7) and appear to be linked (Fig. 2D).  
363 Furthermore, we detected the loss or gain of groups of infection-associated genes between some case  
364 clinical and rectal pairs, rather than individual genes (strata 211, 214, 701, 713, 716, Supplementary  
365 Fig. S2F). To determine the association between plasmids and infection-associated genes, we used  
366 the PlasmidFinder tool to detect replicons in the rectal isolates in our dataset<sup>31</sup>. We detected a replicon

367 in 163 of 245 isolates (66.5%, 72.94% of case rectal isolates, 63.13% of control isolates), with a median  
368 of 2 replicons per isolate (Supplementary Table 7, range = 0-10 replicons). There was no significant  
369 difference between the number of replicons per isolate between cases and controls (Supplementary  
370 Fig. S11A). The three most common replicons were Col(pHAD28), IncFIB(K), and IncFII(K), and these  
371 replicons were distributed across the core genome phylogeny (Supplementary Fig. S11B); however,  
372 no replicon was significantly associated with infection (Supplementary Table 8).

373  
374 To determine whether infection-associated genes are located on plasmids or the chromosome, we  
375 performed long-read sequencing and hybrid assembly with existing short read data on 10 isolates from  
376 cases of infection (Fig. 4A-C, S12). All of these isolates contained plasmids with infection-associated  
377 genes, and no two plasmids containing infection-associated genes were identical in terms of size,  
378 organization and gene content. At one extreme, the rectal case isolate Kp11407 contained a 208,703  
379 bp hybrid IncFIB(K)/IncFII(K) plasmid encoding all 18 infection-associated loci present in this isolate  
380 (Fig. 4A). This plasmid closely aligns to the conjugative plasmid pUUH239.2 (CP002474.1), which was  
381 linked to a nosocomial *K. pneumoniae* outbreak<sup>32</sup>. Notably, this plasmid also contains several heavy  
382 metal resistance operons, including *ars*, *pco*, and *si*<sup>β2</sup>, which may enhance resistance to innate immune  
383 responses<sup>33</sup>. A second case rectal isolate, Kp5854, encoded 23 of 27 infection-associated loci on a  
384 247,093 bp hybrid IncFIB(K)/IncFII(K) plasmid (Fig. 4B). This plasmid closely aligns to pF16KP0070-1  
385 (CP052586.1), which also encodes *ars*, *pco*, and *sil*. Interestingly, in the corresponding clinical isolate,  
386 this plasmid appears to have integrated into the chromosome with the loss of some infection-associated  
387 genes (Fig. 4C; stratum 713, Supplementary Fig. S3G). These isolates also contained four additional  
388 infection-associated genes that mapped to the chromosome. The number of plasmid-encoded,  
389 infection-associated genes varied widely between plasmids. In contrast to Kp11407 and Kp5854,  
390 Kp11162 contained one plasmid with two infection-associated genes and a second plasmid with four  
391 (Supplementary Fig. S12G-H). Across these isolates 23 of all 27 infection-associated genes were  
392 found on at least one plasmid (Supplementary Fig. S12). Only *group\_774* and the *tehA\_2* gene cluster

393 were not found on at least one plasmid. This suggests that these large plasmids facilitate acquisition of  
 394 infection-associated loci and are likely drivers of the association between many genes and infection.  
 395



396  
 397 **Fig. 4. Infection-associated genes are present on large, conjugative plasmids**  
 398 (A) Plasmid map of Kp11407 plasmid 1. (B) Plasmid map of Kp5854 plasmid 1. (C) Alignment of  
 399 Kp7637 chromosome and Kp5854 plasmid 1. (D) Transconjugation rates using Kp11407 as a donor  
 400 strain to Kp4236, Kp4505, and Kp7994, which are *K. quasipneumoniae*, *K. pneumoniae*, and *K.*  
 401 *variicola*, respectively, as recipient strains (mean displayed, N.D. = None detected). For D, each data  
 402 point represents a single biological replicate.  
 403

#### 404 **Infection-associated genes are transmissible by plasmids**

405 To determine if infection-associated genes could be horizontally transmitted, we performed  
406 transconjugation assays between rectal isolates. Kp11407 was selected as a donor strain, as it is a *K.*  
407 *pneumoniae* case rectal isolate that contains various infection-associated genes, including *bla*<sub>CTX-M-15</sub>  
408 and *dgcE*, as well as the Col(pHAD28), IncFIB(K), and IncFII(K) replicons. Three cefuroxime-sensitive  
409 control isolates lacking these genes, *K. pneumoniae* (Kp4505), *K. variicola* (Kp7994), and *K.*  
410 *quasipneumoniae* (Kp4236), were selected as recipient strains. Kp7994 and Kp4236 lack plasmid  
411 replicons, and Kp4505 contains an IncFIB(K) replicon. Conjugation was detected between Kp11407  
412 and Kp4505 and Kp7994 (Fig. 4D). All transconjugants were cefuroxime resistant, indicating the  
413 transfer of *bla*<sub>CTX-M-15</sub>, and were PCR-positive for *dgcE*. Therefore, we concluded that this large plasmid  
414 is conjugative resulting in the transfer of infection-associated genes.

#### 416 **Validation of infection-associated genes in a geographically distinct patient cohort**

417 Finally, we aimed to validate our findings across multiple geographical locations. To this end, we  
418 screened a cohort of *Klebsiella* colonized patients from a single academic medical center in Maryland,  
419 USA for six infection-associated genes: *aac(6')-Ib-cr5*, *bla*<sub>CTX-M-15</sub>, *group\_13461*, *group\_13467*, *dgcE*,  
420 and *aqpZ\_2*. These genes were selected because they were representative of the many facets of our  
421 findings, encapsulating known AMR genes (*aac(6')-Ib-cr5*, *bla*<sub>CTX-M-15</sub>), MGEs (*group\_13461*,  
422 *group\_13467*), and those where the relationship between gene presence and infection was (*aac(6')-Ib-*  
423 *cr5*, *dgcE*, *aqpZ\_2*) and was not (*bla*<sub>CTX-M-15</sub>, *group\_13461*, *group\_13467*) significantly mediated by  
424 *Klebsiella* gut relative abundance. Moreover, these genes are found on conjugative plasmids (Figure  
425 4). This second cohort consisted of 55 cases of *Klebsiella* infection and 432 asymptotically colonized  
426 control patients. Since *wzi* concordance data was not available, we screened both case isolates and a  
427 rectal isolate from each patient for each gene. Rectal isolates were not available for three infected  
428 patients and a clinical isolate was not available for one patient. In total, 538 *Klebsiella* isolates (54  
429 clinical isolates, 52 case rectal isolates, 432 control rectal isolates) were screened for the presence or

absence of target genes using a multiplex RT-PCR assay. As an internal control for *Klebsiella* detection, we included a previously validated assay for the *fiu* gene<sup>9</sup>. Based on either clinical or rectal isolates from cases, we observed a significant association with infection for all genes except *aqpZ\_2* (Table 4). Overall, 82.4% (42/51 of available pairs) of clinical-rectal pairs returned concordant RT-PCR results, which was highly consistent with our previous study using *wzi* sequencing to determine concordance<sup>8</sup>. As in the primary cohort (Supplementary Fig. S7), the presence of these genes in case rectal isolates were associated with one another (Supplementary Fig. S13), suggesting that these genes in the Maryland isolates may also be linked through conjugative plasmids. The strongest associations were between *dgcE* and *apqZ\_2* and between *aac(6')-Ib-cr5*, *bla<sub>CTX-M-15</sub>*, *group\_13461* and *group\_13467*. Overall, five of six genes were reproducibly associated with case status in colonized patients in two geographically distinct cohorts.

**Table 4. Gene association with positive clinical cultures in a geographically independent cohort of *Klebsiella* colonized patients.**

Gene	Sample	Frequency	OR	95% CI	P value
<i>aac(6')-Ib-cr5</i>	Case clinical	0.167	4.88	1.96 to 11.1	<b>0.0009</b>
	Case rectal	0.115	3.18	1.22 to 8.41	<b>0.0276</b>
	Control rectal	0.039			
<i>dgcE</i>	Case clinical	0.481	2.28	1.29 to 3.99	<b>0.0074</b>
	Case rectal	0.5	2.46	1.37 to 4.41	<b>0.0039</b>
	Control rectal	0.289			
<i>group_13461</i>	Case clinical	0.204	4.55	2.06 to 9.6	<b>0.0005</b>
	Case rectal	0.192	4.23	1.93 to 9.24	<b>0.0042</b>
	Control rectal	0.053			
<i>group_13467</i>	Case clinical	0.296	3.81	1.94 to 7.29	<b>0.0002</b>
	Case rectal	0.25	3.02	1.48 to 6.1	<b>0.0043</b>
	Control rectal	0.1			
<i>bla<sub>CTX-M-15</sub></i>	Case clinical	0.185	4.94	2.19 to 11.4	<b>0.0005</b>
	Case rectal	0.173	4.55	1.83 to 10.9	<b>0.0013</b>
	Control rectal	0.044			
<i>aqpZ_2</i>	Case clinical	0.296	1.58	0.83 to 2.93	0.164
	Case rectal	0.327	1.82	0.98 to 3.38	0.076
	Control rectal	0.213			

## 446 Discussion

447 In this study, we have identified novel *Klebsiella* genes associated with infection in colonized patients.  
448 The findings from this diverse collection of *Klebsiella* isolates, carefully adjusted for patient risk factors  
449 and phylogeny, and validated in a second geographically distinct cohort, provides important insights  
450 into the broader dynamics of *Klebsiella* pathogenesis. Most of these genes were found across a diverse  
451 set of plasmids, suggesting that these genes are either direct mediators of fitness and virulence or  
452 indirect markers of large, conjugative plasmids that increase infection risk as a unit. These genes can  
453 be separated into two categories: those whose association with infection is mediated through increased  
454 abundance in the gut, and those that have a direct association with infection. Genes whose effect on  
455 infection risk is not mediated through intestinal dominance may act at a subsequent step of  
456 pathogenesis. These genes associated with infection are readily mobilized within and across members  
457 of the *Klebsiella pneumoniae* species complex, and changes in infection-associated gene content and  
458 genetic context were detectable within colonized patients who progressed to infection. Collectively, this  
459 study provides a novel set of genetic signatures of infection that have potential for predicting risk of  
460 infection in patients colonized by diverse *Klebsiella*.

461  
462 The majority of infection-associated loci in our dataset are present on plasmids, which can vary in size,  
463 gene content, copy number, and conjugative ability (Fig. 4A-C, Supplementary Fig. S12). This plasmid  
464 carriage creates genetic linkage between genes that are associated with infection, raising the possibility  
465 that some are direct effectors of fitness or virulence whereas others are robust markers for these  
466 effectors. In some instances, clear relationships exist between linked genes. For example, nucleotide  
467 cyclases are often present in tandem with cyclic nucleotide phosphodiesterases, thus the linkage  
468 of *dgcE* and *pdeC\_2* is expected. *dgcE* and *pdeC\_2* may regulate expression of type 1 and 3 pili  
469 through modulation of cyclic di-GMP levels. Expression of the type 1 pilus reduces virulence in the lung  
470 and spleen but enhances gut colonization in murine models<sup>34,35</sup>. Therefore, *dgcE* may enhance fitness  
471 in the gut while *pdeC\_2* may enhance fitness in extra-intestinal sites. Likewise, the linkage of *bla*<sub>CTX-M</sub>

472 15 and *group\_13461* is expected since ISEcp1 transposases are known to mobilize and promote the  
473 expression of CTX-M-type  $\beta$ -lactamases<sup>36,37</sup>. Conversely, relationships between other linked genes,  
474 such as *arsD\_1* and the various AMR genes, are more difficult to understand. It may be the case that  
475 these infection-associated genes act at different stages of infection or one or both are robust genetic  
476 markers for fitness genes elsewhere on the plasmid.

477  
478 The strong association between multiple AMR genes and infection following adjustment for high-risk  
479 antibiotic exposure indicates that the relationship between antibiotic resistance and infection is complex  
480 in *Klebsiella*. Antibiotic resistance genes could be associated with infection by enabling intestinal  
481 outgrowth and dominance after exposure to a cognate antibiotic. However, the ESBL determinants  
482 *bla*<sub>CTX-M-15</sub>, *bla*<sub>OXA-1</sub>, and *bla*<sub>TEM-1</sub> were associated with infection independent of antibiotic exposure and  
483 their association was not significantly mediated by gut relative abundance. Potential fitness benefits of  
484 AMR genes in the absence of antibiotics are intriguing, yet highly AMR *Klebsiella* STs are generally  
485 less virulent in experimental models than their hypervirulent counterparts, and specific mutations that  
486 confer AMR can have deleterious effects on *Klebsiella* fitness<sup>38-41</sup>. On the other hand, plasmid  
487 acquisition of AMR is less costly than chromosomal mutations<sup>42</sup>. Therefore, AMR genes may be linked  
488 to plasmid-encoded genes that enhance infection risk, and then complicate treatment of infections after  
489 they are initiated.

490  
491 The identification of non-AMR genes associated with infection as mediated by gut dominance adds to  
492 our understanding of bacterial genotypes and the risk for infection in *Klebsiella* colonized patients. We  
493 previously demonstrated that the *ter* operon, annotated as a tellurite resistance operon, is a  
494 microbiome-dependent gut fitness factor<sup>20</sup>. Similarly, the *arsD\_1* gene identified in this study is part of  
495 an arsenic resistance (*ars*) operon. Tellurite and arsenic are both toxic metalloids, implying that  
496 the *ars* operon may be playing an analogous role to *ter*. Consistent with this, gut relative abundance  
497 mediates the association between *arsD\_1* and infection (Table 2). Conversely, we did not find evidence

498 that gut relative abundance mediates the association between the tellurite resistance gene *tehA\_2* gene  
499 and infection, although it is predicted to encode a different molecular function than *ter*. These data  
500 highlight the complex interplay between bacterial genotype and gut fitness, and the need to understand  
501 the role of colonization fitness to both assess infection risk and prophylactic decolonization  
502 therapies. The risk of *Klebsiella* disease in colonized patients may not warrant prophylactic antibiotic  
503 administration, especially in the context of rising rates of AMR; however, non-antibiotic options are  
504 currently being explored (reviewed in<sup>43</sup>). A better understanding of how bacterial pathogens, such as  
505 *Klebsiella*, successfully colonize and dominate the gut is needed to ensure the efficacy of any  
506 decolonization approach.

507  
508 Although this study advances our understanding of infection in patients colonized by diverse *Klebsiella*,  
509 it is not without its limitations. First, this study is retrospective and relies on the availability of clinical  
510 data in the electronic medical record. Prospective studies with longer timeframes aimed at evaluating  
511 the association between genes identified in this study and infection would be useful for validating our  
512 findings. Second, many infection-associated genes are collinear and plasmid-borne. This finding is  
513 important; however, it complicates the evaluation of the causal relationship between individual genes  
514 and infection. Experimental studies are necessary to determine the causal relationship between each  
515 gene and disease. Third, the *Klebsiella* genes associated with infection likely differ by patient  
516 population. We found different genes associated with infection (*ter*, psicose utilization locus) in a  
517 hospital-wide comparative genomics study and experimentally validated them as encoding fitness  
518 factors<sup>19,20</sup>. That patient population may have different risk factors and stronger barriers to infection.  
519 Among intensive care and hematology/oncology patients, the genes associated with infection were  
520 highly reproducible in a geographically independent cohort. Finally, this study focused on gene  
521 presence in the accessory genome and does not consider the role of mutations within the core or  
522 accessory genomes.

524 This study provides insights into infection risk from diverse lineages of the *K. pneumoniae* species  
525 complex. The well-characterized virulence genes of hypervirulent clonal groups, many of which are  
526 plasmid-borne<sup>22</sup>, were rare and not informative in assessing infection risk in this population. MDR  
527 isolates were more frequent representing six clonal groups, but overall, 85.7% (210/245) of the isolates  
528 in our collection were from non-hypervirulent, non-MDR CGs. Given the predominance of non-  
529 hypervirulent, or “classical” strains worldwide<sup>22</sup>, efforts to understand the pathogenesis and enact  
530 preventative therapies is likely to have a substantial payoff. The varied combinations of infection-  
531 associated genes found across case isolates suggests that a diverse accessory genome acquired  
532 through HGT drives infection risk in *Klebsiella*. Certain combinations have become fixed; some  
533 hypervirulent *Klebsiella* clones have carried their plasmids for decades. However, new combinations  
534 can arise within a colonized patient (Fig. 4; Supplementary Fig. S12). Even in seemingly fixed plasmid-  
535 strain combinations, exchange has occurred leading to multiple instances of the convergence of  
536 virulence and MDR determinants<sup>22</sup>. In total, this study adds to the growing understanding of the  
537 *Klebsiella* accessory genome as a critical determinant of infection risk.

538  
539 The insights from this study can be used to predict infection risk in colonized, hospitalized patients and  
540 trigger infection prevention interventions. We have previously identified certain patient risk factors  
541 associated with infection, and that gut dominance is associated with infection independently of these  
542 clinical factors<sup>8,9</sup>. Here we identified infection-associated genes independent of patient risk factors and  
543 that part of the association is mediated by gut dominance. We have also developed the tools to measure  
544 these risks, including a robust qPCR assay to measure intestinal dominance and the genotyping PCR  
545 assays implemented in this study. The fact that AMR genes and phenotypes are strongly associated  
546 with infection could also enable the use of standard culture approaches to detect these genes, and  
547 therefore, high-risk colonizing strains. It is possible that these approaches can be combined with  
548 queries of the electronic medical record for clinical variables to generate an integrated risk score for a  
549 colonized patient. In the near term, this score could be used to identify patients for trials of novel

550 decolonization therapies. Eventually, integrated risk models could identify at-risk patients who could  
551 most benefit from effective prevention interventions.

552

## 553 **Methods**

### 554 **Ethics statement**

555 Patient enrollment and sample collection at the University of Michigan was approved by and performed  
556 in accordance with the Institutional Review Boards (IRB) of the University of Michigan Medical School  
557 (Study number HUM00123033). Patient enrollment and sample collection at Johns Hopkins University  
558 was approved and performed in accordance with the IRB of the Johns Hopkins University (Study  
559 number IRB00129775).

### 561 **Sample collection and strain selection**

562 At Michigan Medicine, *Klebsiella* isolates were selected from a large cohort of *Klebsiella* colonized  
563 patients. Cohort identification, enrollment, clinical data extraction, chart review, and case definitions are  
564 described in detail elsewhere<sup>8</sup>. During the study period, patient rectal swabs were screened for the  
565 presence of *Klebsiella* by plating on MacConkey agar followed by taxonomic identification using MALDI  
566 TOF mass spectrometry. Up to three *Klebsiella* isolates per rectal swab, and any subsequent *Klebsiella*  
567 positive clinical cultures, were archived. Minimum inhibitory concentration data from case clinical  
568 isolates was extracted from the Michigan Medicine clinical microbiology lab. To identify concordant  
569 infections, *Klebsiella* rectal swab and case clinical isolates were subjected to *wzi* PCR and Sanger  
570 sequencing, and the *wzi* type was assigned by uploading the consensus *wzi* sequence to the BigsDB  
571 database (<https://bigsdb.pasteur.fr>). Each case was matched to two controls based on sex, age ( $\pm$  10  
572 years), and date of rectal swab collection ( $\pm$  90 days). Several cases could only be matched to a control  
573 swab within 120 days ( $n = 6$ ). Case matching was conducted in R, version 4.0.3.

574  
575 At the Johns Hopkins University, adult patients admitted to medical or surgical ICUs are routinely  
576 screened for VRE rectal colonization on admission and weekly thereafter until ICU discharge. Residual  
577 Amies broth from rectal ESwabs (COPAN Diagnostics, Inc., Murrieta, California) were inoculated on a

578 MacConkey agar plate using the WASP DT (COPAN) automated specimen processor, and isolates  
579 were archived as above.

### 581 **Whole genome sequencing, assembly, and annotation**

582 Archived *Klebsiella* isolates were cultured in Luria-Bertani (LB, Becton, Dickinson and Company,  
583 Franklin Lakes, NJ) broth at 37° C with shaking for genomic DNA extraction. Genomic DNA was  
584 extracted using the DNeasy UltraClean Microbial Kit (Qiagen, Hilden, Germany) per manufacturer's  
585 instructions. Extracted genomic DNA was then checked for quality and purity and submitted to the  
586 University of Michigan Advanced Genomics Core, where it was sequenced on an Illumina NovaSeq  
587 using a 300 cycle S4 flow cell (Illumina, San Diego, USA) to generate approximately 6,000,000 150 bp  
588 paired-end reads. Resulting reads were trimmed using Trimmomatic v0.39<sup>44</sup> and assembled into  
589 scaffolds using SPAdes v3.13.0<sup>45</sup>. Genome assemblies within each stratum were compared using  
590 QUAST v5.0.2<sup>46</sup>, wherein each genome assembly was compared to the case rectal isolate genome  
591 assembly. Genome annotation was performed using Prokka v1.14.6<sup>47</sup>, and a pangenome was  
592 assembled using Roary v3.13.0<sup>48</sup>, with settings "-i 90" and "-s." Finally, assembly genotyping was  
593 performed using Kleborate v1.0.0<sup>22</sup> and PlasmidFinder v2.0.1<sup>31</sup>. For long-read sequencing, genomic  
594 DNA was checked for quality and purity and submitted to the University of Michigan Advanced  
595 Genomics Core, where it was sequenced on a GridION X5 platform using a MinION flow cell (Oxford  
596 Nanopore Technology, Oxford, UK). Low-quality read ends were trimmed with NanoFilt v2.5.0<sup>49</sup> and  
597 assembled using Unicycler v0.4.8 using the "bold" setting<sup>50</sup>. The core gene alignment was created using  
598 Illumina assemblies with cognac using the default parameters<sup>51</sup>. The *wzi* gene sequences from each  
599 Illumina assemblies were identified by BLAST v2.9.0<sup>52</sup>. The BLAST output was parsed to extract the  
600 aligned sequence from each query, and the gene sequences were aligned with MAFFT v7.310.  
601 Approximate maximum likelihood trees were generated from the *wzi* alignment and cognac alignment  
602 with fastTree v2.1.10<sup>353</sup>. Phylogenetic trees were visualized with the APE R package v5.3<sup>54</sup>.  
603 Rarefaction was performed using iNEXT<sup>55</sup>, with setting " $q = 0$ ." All analyses were performed in R v4.0.5.

604 Plasmid maps and corresponding alignments were constructed and visualized using BRIG<sup>56</sup>. All code  
605 developed for this manuscript is available at [https://github.com/rdcrawford/Bachman\\_CU8](https://github.com/rdcrawford/Bachman_CU8).

### 607 **Phenotypic antibiotic screen**

608 To screen for phenotypic antibiotic resistance, all *Klebsiella* rectal isolates were arrayed into 96-well  
609 plates and stored at -80° C in 20% LB-glycerol. We were unable to locate one isolate in our freezer  
610 following WGS, thus it was omitted from phenotypic analysis. Then *Klebsiella* rectal isolates were sub-  
611 cultured in LB broth at 37° C overnight and replica plated on antibiotic containing LB-agar. Antibiotic  
612 concentrations (36 mg/mL ampicillin, 32 mg/mL chloramphenicol, 16 mg/mL gentamicin, 1 mg/mL  
613 ciprofloxacin, 32 mg/mL cefazolin, 32 mg/mL cefuroxime, 16 mg/mL cefepime) were selected in  
614 accordance with CLSI criteria for resistance. All assays were performed in triplicate.

### 616 **Transconjugation assay**

617 Single colonies of Kp4236, Kp4505, Kp7994 (recipient strains; chloramphenicol resistant), and  
618 Kp11407 (donor strain; cefuroxime resistant) were sub-cultured overnight in 1 mL of LB broth at 37° C  
619 with shaking. The next day, overnight cultures were harvested by centrifugation and resuspended in 1  
620 mL sterile phosphate-buffered saline (PBS). 50 µL of Kp4236, Kp4505, or Kp7994 resuspension was  
621 mixed with 50 mL of Kp11407 resuspension and spotted on LB-agar. 50 mL of unmixed resuspension  
622 of each strain was spotted in parallel. Bacterial spots were grown overnight at 37° C, then harvested  
623 into 1 mL sterile PBS by scraping and serially plated on differential growth media. Transconjugation  
624 rates were calculated by quantifying the bacterial density of dual-resistant colonies divided by the total  
625 density of the recipient strain. Transconjugates were confirmed by qPCR targeting the *dcgE* gene using  
626 unmixed strains as a reference (assay details below). All assays were performed in triplicate, and 3  
627 transconjugate colonies and 3 unmixed recipient and donor colonies were selected per assay for PCR  
628 confirmation.

## RT-PCR design and assay

To screen for the presence of *aac(6')-Ib-cr5*, *bla<sub>CTX-M-15</sub>*, *group\_13461*, *group\_13467*, *dgcE*, and *aqpZ\_2*, multiplex real-time PCR assays were designed using PanelPlex (DNA Software, Ann Arbor, USA) and combined with a previously validated *Klebsiella*-specific *fiu* assay<sup>9</sup>. All primer and probe sequences can be found in Supplementary Table 9. Clinical and rectal isolates were arrayed and subcultured in 150  $\mu$ L LB overnight at 37° C. To extract genomic DNA, overnight cultures were harvested by centrifugation and resuspended in 100  $\mu$ L of PCR-grade water, boiled at 95° C for 5 minutes, and pelleted. Then, 5  $\mu$ L of the supernatant was used as a PCR template. Primer (Integrated DNA Technologies, Coralville, USA) and probe (Thermo Fisher Scientific, Waltham, USA) concentrations were as previously described<sup>9</sup>. Real-time PCR was performed using a QuantStudio 3 real-time thermocycler (Thermo Fisher Scientific, Waltham, USA) using the following protocol: 50°C for 2 minutes, 95°C for 15 minutes, then 40 cycles of 94°C for 1 minute and 60°C for 1 minute. All assays were performed in duplicate unless replicates conflicted, wherein a third replicate was performed, and all isolate arrays included 3 positive controls and 3 randomly dispersed negative controls. The target gene was assigned as present in an isolate if a reaction was 1) positive for *fiu*, 2) amplified within 2 Ct of the range of the positive control Cts, and 3) results were concordant between all replicates.

## Statistical analysis

For clinical variable comparisons, odds ratios and associated 95% confidence intervals were calculated using conditional logistic regression in R (version 3.6.3) and a *P* value of less than 0.05 was considered statically significant. For infection-associated gene discovery, a *P* value of less than 0.05 and Q value less than 0.2 was used as a cutoff for further exploration. For these analyses, Scoary v1.6.16<sup>23</sup> was used with the "--collapse" and "-n" settings with our custom core genome phylogeny, limited to genes that were 5-95% frequent in the dataset. For analyses considering population structure, treeWAS v1.1<sup>26</sup> was used with the phen.type = "discrete" setting and our custom core genome phylogeny. For gene-gene associations, odds ratios were calculated using Fishers exact test in R Studio v1.2.5001. To

656 account for the inflated type I error rate due to the large number of comparisons, a  $P$  value of less than  
657 0.01 was considered statistically significant for subsequent clinical variable adjustment, mediation,  
658 population structure, and genotyping analyses.

659  
660 Next, we considered adjustment of clinical variables. Inverse probability of treatment weighting method  
661 was used for this purpose<sup>57</sup>. The analysis is performed separately for each gene. In the propensity  
662 score model, a logistic regression was used with the clinical variables as predictors and binary gene  
663 presence/absence as the outcome. The predicted probabilities estimated from the regression model  
664 were used to compute weights and then included in the outcome model. The outcome model is a logistic  
665 regression with *Klebsiella* infection as the outcome and gene (presence/absence) as the predictor,  
666 while incorporating the matched nature of the data. Following IPTW analysis, we assessed if each gene  
667 was associated with infection after individually adjusting for exposure to each class of antibiotics. We  
668 repeated the described IPTW analysis using the same weights, controlling for each prior exposure as  
669 a covariate in the outcome model one at a time. The “survey” package in R was used to obtain the  
670 robust standard error.

671  
672 We then conducted a mediation analysis on each gene that was significantly associated with infection  
673 to explore if the relationship between gene presence/absence and infection is mediated through  
674 abundance. Specifically, we assessed 1) the direct effect defined as the effect of gene presence on  
675 infection, adjusting for abundance as a continuous variable and clinical factors, 2) the indirect effect of  
676 the gene on infection mediated through abundance after adjusting clinical factors. To account for the  
677 matched case-control design, a generalized linear mixed model was used with a random effect for the  
678 strata defining matched case-controls. A significant mediation (i.e., indirect) effect indicates evidence  
679 that the effect of the gene on infection is mediated through abundance, and direct effects were also  
680 reported. The  $P$  values for these effects were assessed to identify genes to investigate further. We  
681 conducted the analyses using the ‘mediation’ and ‘lme4’ packages<sup>58,59</sup> in R v3.6.3.

682

683 For *in vitro* experiments, one-sample *t*-test, unpaired Student's *t*-tests, and ANOVA followed by  
684 indicated post-hoc test was used to determine significant differences between groups following log  
685 transformation of experimental values. For *in vitro* experiments, a *P* value of less than 0.05 was  
686 considered statistically significant. Experimental replicates represent biological replicates. Statistical  
687 analysis of *in vitro* experiments was performed using Prism 8 (GraphPad Software, La Jolla, CA).

688

689 **Data availability**

690 All sequencing files are available from the Sequence Read Archive (SRA) database in BioProject  
691 PRJNA789565. All remaining experimental data are available in the main text or the supplementary  
692 materials. Deidentified data from human subjects may be made available upon request, pending  
693 approval from the University of Michigan Institutional Review Board.

694

## References

- 1 Wyres, K. L., Lam, M. M. C. & Holt, K. E. Population genomics of *Klebsiella pneumoniae*. *Nature reviews. Microbiology* **18**, 344-359, doi:10.1038/s41579-019-0315-1 (2020).
- 2 Magill, S. S. *et al.* Multistate Point-Prevalence Survey of Health Care–Associated Infections. *New England Journal of Medicine* **370**, 1198-1208, doi:10.1056/NEJMoa1306801 (2014).
- 3 CDC. The Direct Medical Costs of Healthcare-Associated Infections in U.S. Hospitals and the Benefits of Prevention. (2009).
- 4 O'Neill. *Tackling Drug-Resistant Infections Globally: final report and recommendations*, <<https://amr-review.org>> (2016).
- 5 Martin, R. M. *et al.* Molecular Epidemiology of Colonizing and Infecting Isolates of *Klebsiella pneumoniae*. *mSphere* **1** (2016).
- 6 Collingwood, A. *et al.* Epidemiological and Microbiome Associations Between *Klebsiella pneumoniae* and Vancomycin-Resistant *Enterococcus* Colonization in Intensive Care Unit Patients. *Open forum infectious diseases* **7**, ofaa012, doi:10.1093/ofid/ofaa012 (2020).
- 7 Gorrie, C. L. *et al.* Gastrointestinal Carriage Is a Major Reservoir of *Klebsiella pneumoniae* Infection in Intensive Care Patients. *Clinical infectious diseases : an official publication of the Infectious Diseases Society of America* **65**, 208-215, doi:10.1093/cid/cix270 (2017).
- 8 Rao, K. *et al.* Risk Factors for *Klebsiella Infections* among Hospitalized Patients with Preexisting Colonization. *mSphere*, e0013221, doi:10.1128/mSphere.00132-21 (2021).
- 9 Sun, Y. *et al.* Measurement of *Klebsiella* Intestinal Colonization Density To Assess Infection Risk. *mSphere*, e0050021, doi:10.1128/mSphere.00500-21 (2021).
- 10 Russo, T. A. & Marr, C. M. Hypervirulent *Klebsiella pneumoniae*. *Clinical microbiology reviews* **32**, doi:10.1128/CMR.00001-19 (2019).
- 11 Wyres, K. L. *et al.* Genomic surveillance for hypervirulence and multi-drug resistance in invasive *Klebsiella pneumoniae* from South and Southeast Asia. *Genome Med* **12**, 11, doi:10.1186/s13073-019-0706-y (2020).

- 721 12 Martin, R. M. & Bachman, M. A. Colonization, Infection, and the Accessory Genome of *Klebsiella*  
722 *pneumoniae*. *Frontiers in cellular and infection microbiology* **8**, 4, doi:10.3389/fcimb.2018.00004  
723 (2018).
- 724 13 Snitkin, E. S. *et al.* Tracking a hospital outbreak of carbapenem-resistant *Klebsiella pneumoniae*  
725 with whole-genome sequencing. *Sci Transl Med* **4**, 148ra116,  
726 doi:10.1126/scitranslmed.3004129 (2012).
- 727 14 Marvig, R. L., Sommer, L. M., Molin, S. & Johansen, H. K. Convergent evolution and adaptation  
728 of *Pseudomonas aeruginosa* within patients with cystic fibrosis. *Nat Genet* **47**, 57-64,  
729 doi:10.1038/ng.3148 (2015).
- 730 15 Allen, J. P. *et al.* A comparative genomics approach identifies contact-dependent growth  
731 inhibition as a virulence determinant. *Proceedings of the National Academy of Sciences of the*  
732 *United States of America* **117**, 6811-6821, doi:10.1073/pnas.1919198117 (2020).
- 733 16 Aggarwal, S. D. *et al.* Function of BriC peptide in the pneumococcal competence and virulence  
734 portfolio. *PLoS pathogens* **14**, e1007328, doi:10.1371/journal.ppat.1007328 (2018).
- 735 17 Mageiros, L. *et al.* Genome evolution and the emergence of pathogenicity in avian *Escherichia*  
736 *coli*. *Nature communications* **12**, 765, doi:10.1038/s41467-021-20988-w (2021).
- 737 18 Bornstein, K. *et al.* Genetic changes associated with the temporal shift in invasive non-typhoidal  
738 *Salmonella* serovars in Bamako Mali. *PLoS Negl Trop Dis* **13**, e0007297,  
739 doi:10.1371/journal.pntd.0007297 (2019).
- 740 19 Martin, R. M. *et al.* Identification of Pathogenicity-Associated Loci in *Klebsiella pneumoniae* from  
741 Hospitalized Patients. *mSystems* **3**, doi:10.1128/mSystems.00015-18 (2018).
- 742 20 Vornhagen, J. *et al.* A plasmid locus associated with *Klebsiella* clinical infections encodes a  
743 microbiome-dependent gut fitness factor. *PLoS pathogens* **17**, e1009537,  
744 doi:10.1371/journal.ppat.1009537 (2021).

- 745 21 Brisse, S. *et al.* wzi Gene sequencing, a rapid method for determination of capsular type for  
746 *Klebsiella* strains. *Journal of clinical microbiology* **51**, 4073-4078, doi:10.1128/JCM.01924-13  
747 (2013).
- 748 22 Lam, M. M. C. *et al.* A genomic surveillance framework and genotyping tool for *Klebsiella*  
749 *pneumoniae* and its related species complex. *Nature communications* **12**, 4188,  
750 doi:10.1038/s41467-021-24448-3 (2021).
- 751 23 Brynildsrud, O., Bohlin, J., Scheffer, L. & Eldholm, V. Rapid scoring of genes in microbial pan-  
752 genome-wide association studies with Scoary. *Genome Biol* **17**, 238, doi:10.1186/s13059-016-  
753 1108-8 (2016).
- 754 24 Baggs, J. *et al.* Risk of Subsequent Sepsis Within 90 Days After a Hospital Stay by Type of  
755 Antibiotic Exposure. *Clinical infectious diseases : an official publication of the Infectious*  
756 *Diseases Society of America* **66**, 1004-1012, doi:10.1093/cid/cix947 (2018).
- 757 25 Robicsek, A. *et al.* Fluoroquinolone-modifying enzyme: a new adaptation of a common  
758 aminoglycoside acetyltransferase. *Nature medicine* **12**, 83-88, doi:10.1038/nm1347 (2006).
- 759 26 Collins, C. & Didelot, X. A phylogenetic method to perform genome-wide association studies in  
760 microbes that accounts for population structure and recombination. *PLoS Comput Biol* **14**,  
761 e1005958, doi:10.1371/journal.pcbi.1005958 (2018).
- 762 27 CLSI. in *CLSI Document M07-A10* (Clinical and Laboratory Standards Institute, Wayne, PA,  
763 2015).
- 764 28 Alcock, B. P. *et al.* CARD 2020: antibiotic resistome surveillance with the comprehensive  
765 antibiotic resistance database. *Nucleic acids research* **48**, D517-D525, doi:10.1093/nar/gkz935  
766 (2020).
- 767 29 Wu, J. & Rosen, B. P. The *arsD* gene encodes a second trans-acting regulatory protein of the  
768 plasmid-encoded arsenical resistance operon. *Molecular microbiology* **8**, 615-623,  
769 doi:10.1111/j.1365-2958.1993.tb01605.x (1993).

- 770 30 Serfiotis-Mitsa, D. *et al.* The structure of the KlcA and ArdB proteins reveals a novel fold and  
771 antirestriction activity against Type I DNA restriction systems *in vivo* but not *in vitro*. *Nucleic acids*  
772 *research* **38**, 1723-1737, doi:10.1093/nar/gkp1144 (2010).
- 773 31 Carattoli, A. & Hasman, H. PlasmidFinder and In Silico pMLST: Identification and Typing of  
774 Plasmid Replicons in Whole-Genome Sequencing (WGS). *Methods Mol Biol* **2075**, 285-294,  
775 doi:10.1007/978-1-4939-9877-7\_20 (2020).
- 776 32 Sandegren, L., Linkevicius, M., Lytsy, B., Melhus, A. & Andersson, D. I. Transfer of an  
777 *Escherichia coli* ST131 multiresistance cassette has created a *Klebsiella pneumoniae*-specific  
778 plasmid associated with a major nosocomial outbreak. *J Antimicrob Chemother* **67**, 74-83,  
779 doi:10.1093/jac/dkr405 (2012).
- 780 33 Djoko, K. Y., Ong, C. L., Walker, M. J. & McEwan, A. G. The Role of Copper and Zinc Toxicity  
781 in Innate Immune Defense against Bacterial Pathogens. *The Journal of biological chemistry* **290**,  
782 18954-18961, doi:10.1074/jbc.R115.647099 (2015).
- 783 34 Khater, F. *et al.* In silico analysis of usher encoding genes in *Klebsiella pneumoniae* and  
784 characterization of their role in adhesion and colonization. *PloS one* **10**, e0116215,  
785 doi:10.1371/journal.pone.0116215 (2015).
- 786 35 Rosen, D. A., Twentyman, J. & Hunstad, D. A. High Levels of Cyclic Di-GMP in *Klebsiella*  
787 *pneumoniae* Attenuate Virulence in the Lung. *Infection and immunity* **86**, doi:10.1128/IAI.00647-  
788 17 (2018).
- 789 36 Bonnet, R. Growing group of extended-spectrum beta-lactamases: the CTX-M enzymes.  
790 *Antimicrobial agents and chemotherapy* **48**, 1-14, doi:10.1128/aac.48.1.1-14.2004 (2004).
- 791 37 Poirel, L., Decousser, J. W. & Nordmann, P. Insertion sequence ISEcp1B is involved in  
792 expression and mobilization of a *bla(CTX-M)* beta-lactamase gene. *Antimicrobial agents and*  
793 *chemotherapy* **47**, 2938-2945, doi:10.1128/aac.47.9.2938-2945.2003 (2003).

- 794 38 Bruchmann, S., Feltwell, T., Parkhill, J. & Short, F. L. Identifying virulence determinants of  
795 multidrug-resistant *Klebsiella pneumoniae* in *Galleria mellonella*. *Pathog Dis* **79**,  
796 doi:10.1093/femspd/ftab009 (2021).
- 797 39 Xiong, H. *et al.* Distinct Contributions of Neutrophils and CCR2+ Monocytes to Pulmonary  
798 Clearance of Different *Klebsiella pneumoniae* Strains. *Infection and immunity* **83**, 3418-3427,  
799 doi:10.1128/IAI.00678-15 (2015).
- 800 40 Olonisakin, T. F. *et al.* CD36 Provides Host Protection Against *Klebsiella pneumoniae*  
801 Intrapulmonary Infection by Enhancing Lipopolysaccharide Responsiveness and Macrophage  
802 Phagocytosis. *The Journal of infectious diseases* **214**, 1865-1875, doi:10.1093/infdis/jiw451  
803 (2016).
- 804 41 Wong, J. L. C. *et al.* OmpK36-mediated Carbapenem resistance attenuates ST258 *Klebsiella*  
805 *pneumoniae* *in vivo*. *Nature communications* **10**, 3957, doi:10.1038/s41467-019-11756-y  
806 (2019).
- 807 42 Vogwill, T. & MacLean, R. C. The genetic basis of the fitness costs of antimicrobial resistance:  
808 a meta-analysis approach. *Evol Appl* **8**, 284-295, doi:10.1111/eva.12202 (2015).
- 809 43 Gargiullo, L., Del Chierico, F., D'Argenio, P. & Putignani, L. Gut Microbiota Modulation for  
810 Multidrug-Resistant Organism Decolonization: Present and Future Perspectives. *Front Microbiol*  
811 **10**, 1704, doi:10.3389/fmicb.2019.01704 (2019).
- 812 44 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence  
813 data. *Bioinformatics* **30**, 2114-2120, doi:10.1093/bioinformatics/btu170 (2014).
- 814 45 Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to single-  
815 cell sequencing. *J Comput Biol* **19**, 455-477, doi:10.1089/cmb.2012.0021 (2012).
- 816 46 Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUASt: quality assessment tool for genome  
817 assemblies. *Bioinformatics* **29**, 1072-1075, doi:10.1093/bioinformatics/btt086 (2013).
- 818 47 Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068-2069,  
819 doi:10.1093/bioinformatics/btu153 (2014).

820 48 Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**,  
821 3691-3693, doi:10.1093/bioinformatics/btv421 (2015).

822 49 De Coster, W., D'Hert, S., Schultz, D. T., Cruts, M. & Van Broeckhoven, C. NanoPack: visualizing  
823 and processing long-read sequencing data. *Bioinformatics* **34**, 2666-2669,  
824 doi:10.1093/bioinformatics/bty149 (2018).

825 50 Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial genome  
826 assemblies from short and long sequencing reads. *PLoS Comput Biol* **13**, e1005595,  
827 doi:10.1371/journal.pcbi.1005595 (2017).

828 51 Crawford, R. D. & Snitkin, E. S. cognac: rapid generation of concatenated gene alignments for  
829 phylogenetic inference from large, bacterial whole genome sequencing datasets. *BMC*  
830 *bioinformatics* **22**, 70, doi:10.1186/s12859-021-03981-4 (2021).

831 52 Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search  
832 tool. *Journal of molecular biology* **215**, 403-410, doi:10.1016/S0022-2836(05)80360-2 (1990).

833 53 Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree: computing large minimum evolution trees with  
834 profiles instead of a distance matrix. *Mol Biol Evol* **26**, 1641-1650, doi:10.1093/molbev/msp077  
835 (2009).

836 54 Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in R  
837 language. *Bioinformatics* **20**, 289-290, doi:10.1093/bioinformatics/btg412 (2004).

838 55 Hsieh, T. C., Ma, K. H. & Chao, A. iNEXT: an R package for rarefaction and extrapolation of  
839 species diversity (Hill numbers). *Methods in Ecology and Evolution* **7**, 1451-1456,  
840 doi:<https://doi.org/10.1111/2041-210X.12613> (2016).

841 56 Alikhan, N. F., Petty, N. K., Ben Zakour, N. L. & Beatson, S. A. BLAST Ring Image Generator  
842 (BRIG): simple prokaryote genome comparisons. *BMC Genomics* **12**, 402, doi:10.1186/1471-  
843 2164-12-402 (2011).

844 57 Haukoos, J. S. & Lewis, R. J. The Propensity Score. *Jama* **314**, 1637-1638,  
845 doi:10.1001/jama.2015.13480 (2015).

846 58 Tingley, D., Yamamoto, T., Hirose, K., Keele, L. & Imai, K. mediation: R Package for Causal  
847 Mediation Analysis. *2014* **59**, 38, doi:10.18637/jss.v059.i05 (2014).  
848 59 Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting Linear Mixed-Effects Models Using lme4.  
849 *2015* **67**, 48, doi:10.18637/jss.v067.i01 (2015).

850

851

## 852 **Acknowledgements**

853 The authors would like to thank the University of Michigan Advanced Genomics Core of the University  
854 of Michigan Medical School's Biomedical Research Core Facilities for their assistance with whole-  
855 genome sequencing. Additionally, we acknowledge support from the Bioinformatics Core of the  
856 University of Michigan Medical School's Biomedical Research Core Facilities.

857

## 858 **Funding**

859 This work was supported by funding from National Institution of Health (<https://www.nih.gov/>) grants  
860 R01AI125307 to MAB and 1R01AI148259-01 to ESS. JV was supported by the Postdoctoral  
861 Translational Scholar Program (NIH UL1TR002240). RDC was supported by the University of Michigan  
862 bioinformatics training grant (NIH T32GM070449). The funders had no role in study design, data  
863 collection and analysis, decision to publish, or preparation of the manuscript.

864

## 865 **Contributions**

866 Conceptualization: JV, LZ, PJS, KR, KSK, MAB

867 Methodology: ER, RC, ESS, LZ, KR, MAB

868 Investigation: JV, ER, LU, AP, RC, SM, CLH, YS, AC, AT

869 Visualization: JV, RC

870 Funding acquisition: MAB

871 Project administration: MAB

872 Supervision: MAB

873 Writing – original draft: JV, ER, LU, SM, MAB

874 Writing – review & editing: JV, ER, LU, AP, RC, SM, CLH, YS, AC, AT, ESS, LZ, PJS, KR, KSK, MAB

875

## 876 **Competing interests**

877 KR is supported in part from an investigator-initiated grant from Merck & Co, Inc. and has consulted for  
878 Bio-K+ International, Inc., Roche Molecular Systems, Inc., and Seres Therapeutics. All other authors  
879 declare that they have no competing interests.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [NatCommsSupplementaryInformation.docx](#)
- [SupplementaryTable1.xlsx](#)
- [SupplementaryTable2.xlsx](#)
- [SupplementaryTable3.xlsx](#)
- [SupplementaryTable4.xlsx](#)
- [SupplementaryTable5.xlsx](#)
- [SupplementaryTable6.xlsx](#)
- [SupplementaryTable7.xlsx](#)
- [SupplementaryTable8.xlsx](#)
- [SupplementaryTable9.xlsx](#)
- [00694rsflat.pdf](#)