1    **Average Nucleotide Identity**

2    The .fsa files from prokka, containing all contigs $\geq$ 500 bp in length from our *Klebsiella* cohort, were
3    added to a separate directory with the publicly available suspected *K. variicola* sequences and the
4    genomes for *K. pneumoniae* CAV1042, *K. pneumoniae* HS11286, and *K. quasipneumoniae* ATCC
5    7000603. We used the command line ANI tool pyANI (https://github.com/widdowquinn/pyani) with the
6    mummer method to compute the pairwise average nucleotide identity among the 207 genomes. The
7    resulting ANIm matrix was viewed as a clustermap using seaborn
8    (https://seaborn.pydata.org/index.html).

9    To further confirm the delineation of *K. variicola* into a major and minor clade, we repeated the ANI
10   analysis using the Jspecies webserver (http://jspecies.ribohost.com/jspeciesws/#analyse) BLAST method
11   in April 2018 between KvMX2 and Yh43 (1) (Table S3).

12   **Core-genome alignment and visualizations**

13   Analysis of K. variicola in relation to other species in the Klebsiella genus was performed by running
14   roary v 3.8.0 with 90% identity on .gff output of prokka from 4 *K. variicola* strains and a strain from *K.*
15   *pneumoniae*, *K. quasipnuemoniae*, *K. quasivariicola*, *K. grimontii*, *K. aerogenes*, *K. oxytoca*, and *K.*
16   *michiganensis*, with *Kluyvera georgiana* used as an outgroup.   The 1,262 genes were aligned in roary
17   with PRANK v1.0 and convereted into a newick file with FastTree v2.1.10 (Table S4).

18   Initial delineation of the population structure for *K. variicola* was performed by moving the .gff file
19   output from prokka into a separate directory and using roary v3.8.0 to identify the core-genomes (those
20   shared by 100% of all genomes) and PRANK v1.0 to align them (Table S5) (2, 3). The resulting alignment
21   file was constructed into a newick tree using FastTree v2.1.10 and viewed in ITOL (4, 5).  To quantify
22   recombination between the strains and identify lineages in a phylogenetic independent manner, we
23   performed FastGear on the roary alignment file of the 3,430 genes shared by the 145 *K. variicola*
24   genomes (Table S6).

25   To understand the population structure of the 143 genomes in the second lineage, we used parSNP
26   within the harvest suite to construct phylogenetic trees from the scaffolds.fasta file with and without
27   recombination (6) (Fig. S1) (Table S7). As an orthologous method, we created an approximate-
28   maximum-likelihood tree of 3,500 core-genes shared by the 143 *K. variicola* isolates in the second
29   lineage from roary and FastTree (Table S8) (Fig. S2). The resulting newick tree and alignment file were
30   used to identify clusters of isolates with ClusterPicker v1.3 (7).  Clusters were identified using an initial
31   and main support threshold of .9, a genetic distance threshold of 4.5, and large cluster threshold of 10.
32   The clusters identified were then visualized on the parSNP tree without recombination with 100%
33   concordance (26/26 clusters). To alternatively view the population structure of *K. variicola*, just the SNP
34   locations were identified by performing snp-sites on the roary alignment file of the 143 genomes in the
35   major clade (8).  This file was visualized as an unrooted equal angle Nearest Neighbor phylogenetic
36   network in SplitsTrees v4 (9, 10). To improve the resolution on the highly related isolates in cluster 21,
37   which contained WUSM_KV_10 and 6 isolates from an investigation of infectious agents in an ICU, we
38   used roary v.3.8.0 to identify the 4,867 core-genes for these 7 genomes at 95%
39   identity{26230489}(Table S9). Single Nucleotide Polymorphisms in these core-genes were identified
40   using SNP-sites{28348851}(Table S10).

41

42  A final phylogenetic tree was created by performing roary and PRANK on *Klebsiella aerogenes* KCTC
43  2190, *K. quasipneumoniae* ATCC 700603, and *K. pneumoniae* ATCC BAA-1705 at >90% identity. The
44  2,932 genes shared by all isolates were used to construct a newick file using FastTree (Table S11).

45  **Antibiotic Resistance Gene, Virulence Gene identification**

46  Acquired ARGs were identified from the .ffn output of prokka using the command line version of
47  ResFinder with default parameters against all available database classes (11). Similarity, the plasmid
48  replicons were identified using the command line version of PlasmidFinder against the
49  Enterobacteriaceae database{24777092}. The number of isolates with > or <= the median number of
50  ARGs and plasmid replicons were tallied and used as input for a Chi-Square test calculator
51  (https://www.socscistatistics.com/tests/chisquare2/Default2.aspx) (Table S12). Virulence genes were
52  annotated by downloading the BIGSDB (http://bigsdb.pasteur.fr/klebsiella/klebsiella.html) list of
53  virulence genes in January 2018 and making them into a custom blast nucleotide database. BLASTN was
54  used to query the .ffn output of prokka against the virulence gene database, with hits requiring 95%
55  identity (Table S12).

56  Acquired ARGs in the WUSM *K. variicola* cohort were viewed as a network diagram in Cytoscape v 3.4.0
57  by constructing a text file where each source node is represented by a unique ARG and the target is the
58  isolate genome with an edge weight of 1 (12) (Table S12).

59

60  *fim* **operon visualization**

61  The complete *fim* operon sequence was obtained from the draft genome of TOP52, a model
62  uropathogenic *K. pneumoniae* strain (13, 14). BLASTN was used to extract the *fim* operon containing
63  contigs from the draft genomes of the strains used for mouse infections and *in vitro* experiments. The
64  contigs were reannotated using prokka, and the GenBank files were visualized in EasyFig without any
65  pairwise BLAST identity values (15). ORFs were colorized based on suspected function. To visualize any
66  SNPs between the different operons, the complete sequence containing the operon was aligned using
67  MUSCLE and visualized in JALView (16).

68  **Usher analysis**

69  Putative usher sequenced were obtained from the pan-genome of the WUSM *K. variicola* isolates by
70  searching the gene_presence_absence.csv output of roary for genes or annotations containing the
71  phrase "outer membrane usher". To ensure that all possible ORFs were identified, the
72  pan_genome_refence.fa containing representatives of every gene in the pan-genome, was compared
73  against the *fimD* usher nucleotide sequencing using the BLASTN webserver in April 2018. To determinate
74  if any *K. variicola* usher sequences were already described, we used protein BLAST to compare the
75  amino acid sequences against *fim, mrk, kpa, kpb, kpc, kpd, kpe, kpf, kpg,* and *kpj* (17, 18) The amino acid
76  sequence for every usher sequence was obtained and added into a multifasta containing representative
77  usher sequences from various Gram-negative phyla described by Nuccio and Baumler (19) (Table S13).
78  The multifasta was aligned using MUSCLE and then converted into a newick tree using FastTree (20, 21).

79    The resulting tree was viewed in ITOL had clades annotated from the Nuccio and Baumler scheme and
80    terminal branches from *K. variicola* labeled by the operon name.

81    The distribution of all usher sequences in the WUSM *K. variicola* pan-genome was surveyed by creating
82    a presence/absence matrix for all usher genes and then hierarchically clustering the matrix in Seaborn.
83    The resulting heatmap was annotated by name of operon, name of isolate, and conservation within *K.*
84    *variicola*. Suspected sequences with truncated ushers were inspected when prokka annotation
85    identified two adjacent usher ORFs and manually annotated on the heatmap.

86    The *K. variicola* specific usher sequences identified were submitted against the nonredundant protein
87    sequences database in April 2018 and had the top hit blast identity values recorded (Table S14). Given
88    that KvhC had the lowest amino acid percent identity of the newly characterized usher proteins, we
89    obtained all the amino acid sequences for blast hits greater than 49% identity and 99% the query length
90    to construct a phylogenetic tree. The amino acid sequences were aligned using MUSCLE and rooted to
91    the nearest usher sequence in our collection, KvaC. The alignment file was made into a newick tree using
92    FastTree and then viewed in ITOL with percent identity and query length values added to each node.

93    The contig containing the *kvh* operon was extracted from the genome of strain WUSM_KV_52 and
94    reannotated with prokka. The resulting GenBank file was viewed in easyfig to observe genes syntenic
95    with the operon. All ORFs were submitted to BLASTP in April 2018 against the nonredundant protein
96    database to identity putative functions. ORFs with suspected roles in transposase and prophage activity
97    were specifically marked.

98

99

100   1.    Richter M, Rossello-Mora R, Oliver Glockner F, Peplies J. 2016. JSpeciesWS: a web server for
101         prokaryotic species circumscription based on pairwise genome comparison. Bioinformatics
102         32:929-31.
103   2.    Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, Fookes M, Falush D, Keane JA,
104         Parkhill J. 2015. Roary: rapid large-scale prokaryote pan genome analysis. Bioinformatics
105         31:3691-3.
106   3.    Loytynoja A. 2014. Phylogeny-aware alignment with PRANK. Methods Mol Biol 1079:155-70.
107   4.    Price MN, Dehal PS, Arkin AP. 2010. FastTree 2--approximately maximum-likelihood trees for
108         large alignments. PLoS One 5:e9490.
109   5.    Letunic I, Bork P. 2016. Interactive tree of life (iTOL) v3: an online tool for the display and
110         annotation of phylogenetic and other trees. Nucleic Acids Res 44:W242-5.
111   6.    Treangen TJ, Ondov BD, Koren S, Phillippy AM. 2014. The Harvest suite for rapid core-genome
112         alignment and visualization of thousands of intraspecific microbial genomes. Genome Biol
113         15:524.
114   7.    Rose R, Lamers SL, Dollar JJ, Grabowski MK, Hodcroft EB, Ragonnet-Cronin M, Wertheim JO,
115         Redd AD, German D, Laeyendecker O. 2017. Identifying Transmission Clusters with Cluster Picker
116         and HIV-TRACE. AIDS Res Hum Retroviruses 33:211-218.
117   8.    Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, Harris SR. 2016. SNP-sites: rapid
118         efficient extraction of SNPs from multi-FASTA alignments. Microb Genom 2:e000056.
119   9.    Huson DH, Bryant D. 2006. Application of phylogenetic networks in evolutionary studies. Mol
120         Biol Evol 23:254-67.

121    10.    Nascimento M, Sousa A, Ramirez M, Francisco AP, Carrico JA, Vaz C. 2017. PHYLOViZ 2.0:
122            providing scalable data integration and visualization for multiple phylogenetic inference
123            methods. Bioinformatics 33:128-129.
124    11.    Kleinheinz KA, Joensen KG, Larsen MV. 2014. Applying the ResFinder and VirulenceFinder web-
125            services for easy identification of acquired antibiotic resistance and E. coli virulence genes in
126            bacteriophage and prophage nucleotide sequences. Bacteriophage 4:e27943.
127    12.    Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T.
128            2003. Cytoscape: a software environment for integrated models of biomolecular interaction
129            networks. Genome Res 13:2498-504.
130    13.    Johnson JG, Spurbeck RR, Sandhu SK, Matson JS. 2014. Genome Sequence of Klebsiella
131            pneumoniae Urinary Tract Isolate Top52. Genome Announc 2.
132    14.    Rosen DA, Pinkner JS, Jones JM, Walker JN, Clegg S, Hultgren SJ. 2008. Utilization of an
133            intracellular bacterial community pathway in Klebsiella pneumoniae urinary tract infection and
134            the effects of FimK on type 1 pilus expression. Infect Immun 76:3337-45.
135    15.    Sullivan MJ, Petty NK, Beatson SA. 2011. Easyfig: a genome comparison visualizer. Bioinformatics
136            27:1009-10.
137    16.    Waterhouse AM, Procter JB, Martin DM, Clamp M, Barton GJ. 2009. Jalview Version 2--a
138            multiple sequence alignment editor and analysis workbench. Bioinformatics 25:1189-91.
139    17.    Khater F, Balestrino D, Charbonnel N, Dufayard JF, Brisse S, Forestier C. 2015. In silico analysis of
140            usher encoding genes in Klebsiella pneumoniae and characterization of their role in adhesion
141            and colonization. PLoS One 10:e0116215.
142    18.    Wu CC, Huang YJ, Fung CP, Peng HL. 2010. Regulation of the Klebsiella pneumoniae Kpc fimbriae
143            by the site-specific recombinase KpcI. Microbiology 156:1983-92.
144    19.    Nuccio SP, Baumler AJ. 2007. Evolution of the chaperone/usher assembly pathway: fimbrial
145            classification goes Greek. Microbiol Mol Biol Rev 71:551-75.
146    20.    Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
147            Nucleic Acids Res 32:1792-7.
148    21.    Price MN, Dehal PS, Arkin AP. 2009. FastTree: computing large minimum evolution trees with
149            profiles instead of a distance matrix. Mol Biol Evol 26:1641-50.

150