

1 **Supplementary Materials**

2 This manuscript contains the following supplemental materials:

3 Supplement A – Detailed UltraSEQ Services (this document)

4 Supplement B – Other UltraSEQ Services (this document)

5 Supplement C – Sample-report_user guide (separate Excel Document)

6 Supplement D – Supplemental_File_Scores (separate Excel Document)

7 Supplement E – UltraSEQ Rules Engine Logic (this document)

8 Supplement F – Supplemental Results (this document)

9 Supplement A – Detailed UltraSEQ Services

10 **Preprocessing Service.** For datasets derived from sequencers, UltraSEQ's preprocessing
11 routine includes steps to trim low quality sequence regions, remove adapter sequences,
12 (optionally) merge paired end reads, and (optionally) remove host sequences to ensure optimal
13 reads remain for analysis. For Illumina and IonTorrent datasets, Trimmomatic v0.39 (Bolger et
14 al., 2014) was used for quality trimming/adaptor removal (settings include
15 ILLUMINACLIP:NexteraPE-PE.fa:2:30:10:2:keepBothReads, LEADING:3 TRAILING:3
16 SLIDINGWINDOW:5:20 MINLEN:50) and Fastp v0.23.0 (Chen et al., 2018) was used for paired
17 end read merging and deduplication (settings include -m (R1/R2 merging mode) when
18 applicable, -dedup, -Q -A (disable quality and adapter trimming). *[Note: subsequent to this
19 publication, UltraSEQ preprocessing routine was updated to use Fastp for quality
20 trimming/adaptor removal and merging in one step with the following setting: -m (R1/R2 merging
21 mode), -cut_front -cut_front_window_size 1 -cut_front_mean_quality 3 (mimics Trimmomatic
22 LEADING:3), -cut_tail -cut_tail_window_size 1 -cut_tail_mean_quality 3 (mimics Trimmomatic
23 TRAILING: 3), -cut_right -cut_right_window_size 5 -cut_right_mean_quality 20 (mimics
24 Trimmomatic SLIDINGWINDOW:4:20), -l 50. A second step was used to deduplicate the
25 dataset if necessary due to interference of Fastp's -cut_right setting on -dedup. We found that
26 this pipeline provides 1.5x speed increase, ~1.07x more usable reads, and automatic adaptor
27 removal (data not shown)].* Following trimming and adapter removal, Bowtie2 v2.3.5.1
28 (Langmead & Salzberg, 2012) was used with default settings to remove host reads that
29 produced an alignment to the human genome build GRCh38. For Nanopore datasets, Porechop
30 v0.2.4 (Wick et al., 2017) was used with default settings to remove adapters and MiniMap2
31 v2.24-r1122 (Li, 2018) was used with default settings to remove any reads that produced an
32 alignment to human genome build GRCh38. FastQC v0.11.9 (Andrews, 2010) and MultiQC
33 v1.10.1 (Ewels et al., 2016) were used to evaluate pre-processed and post-processed standard
34 data quality metrics and ensure preprocessing routines were effective.

35 To ensure the most informative reads are passed to the next UltraSEQ service, an additional
36 de-duplication step is performed by removing any duplicates that have an exact match for the
37 first 50 bases. Such duplication is known to occur during the library preparation step (Head et
38 al., 2014). Further, to reduce cloud compute costs, enhance run-times, and provide better
39 comparisons across datasets, subsampling was optionally performed prior to the alignment
40 service described below. Subsampling was performed by calculating the average number of bps
41 per read in a sample, then randomly sampling to the number of reads required to reach
42 10,100,000 total bps (note: the subsampling was performed by bps instead of number of reads
43 since some datasets, e.g., nanopore, have much longer reads). Subsampling was performed on
44 all datasets (if needed) with the exception of a second set of runs that was performed for
45 Yang et al. (Yang et al., 2019) datasets for antibiotic resistance genotyping. For these runs, a
46 separate UltraSEQ run was used in which the full sample was run without de-hosting to
47 enhance signals for antibiotic resistance genes (with the exception of the following samples that
48 were subsampled to 30,000 reads: Case 22, 10, 9, and 4; these were subsampled to 30,000
49 reads to reduce computational runtime).

50 **Aligner Service.** To avoid sequences with lengths longer than LAMBDA2's maximum query
51 length, sequences longer than 5,000 bps were chunked in pieces of maximum size 5,000.

52 LAMBDA2 enables alignment against both protein and nucleotide databases, but the results for
53 this study leveraged only protein databases, including the Uniref100 protein database (built April
54 2021) and Battelle’s Sequence of Concern protein database, which contains ~8,000 sequences
55 of concern (including virulence factors, toxins, bioregulators, pathways of concern, etc.), ~500
56 signatures of genetic engineering, and ~3,500 biological agents, including ~2,800 pathogens
57 (~2,600 human pathogens) as detailed here (Gemler et al., 2022). For selected runs, we also
58 leveraged our curated nucleotide database of human pathogens and select agents, although
59 these data are not shown here, as no improvement to results was noticed. For this study, the
60 following aligner settings were used: e-value = 1e-4, maximum number matches = 10, aligner's
61 seed-delta-increases-length flag = ON.

62 **Query Mapper Service:** This service maps regions within query sequences to identify high
63 quality alignment regions as well as chimeric reads / out-of-context DNA sequences. This
64 service processes the raw alignment results from the aligner service and identifies top alignment
65 results by first finding the top percent identity and then subsetting the raw alignment results to
66 alignments whose percent identity is within a tolerance, by default 1%, of the top percent
67 identity. The top alignment results are subsequently processed for positional information from
68 each database used, including protein and nucleotide databases. For each query position,
69 n_{counts} is defined as the total number of query alignment starts and query alignment stops
70 corresponding to that position. After n_{counts} has been populated at every position, a normalized
71 vector of counts (N_{counts}) is compiled according to equation 1:

72 Equation 1

73
$$N_{counts} = \frac{n_{counts} - \text{Min}(n_{counts})}{\text{Max}(n_{counts}) - \text{Min}(n_{counts})}$$

74 Following these calculations, a K-means clustering is performed for the N_{counts} values, and the
75 top cluster is used to define the region bounds. Specifically, kmeans++ (Arthur & Vassilvitskii,
76 2007) is used to set the initial centroids. For this application, the “furthest point” algorithm
77 sequentially selects initial centroids furthest from the ones in the previous iteration. Lloyd’s
78 Algorithm (Lloyd, 1982) is then used for clustering given those initial centroids. Finally, the
79 Elbow Method (Ng, 2012) is used to determine the best number of groups k. The top cluster is
80 defined as the cluster with the largest centroid. Further, as implemented within UltraSEQ, the
81 algorithm checks if the top two clusters’ centroids are within a particular tolerance (10% in the
82 case of the query mapper); if they are, the penultimate cluster is absorbed into the top cluster
83 (otherwise, the top cluster remains unaltered).

84 To illustrate these calculations, consider the following example: one query sequence of length
85 150 bps aligns to 3 different subject sequences, with the following start and stop query positions
86 (and percent identities): Accession A, start position 1, end position 100 (percent identity = 100);
87 Accession B, start position 1, end position 50 (percent identity = 95); Accession C, start position
88 25, end position 100 (percent identity = 100). In this case, $n_{counts} = 195, 100, 95,$ and 200 and
89 $N_{counts} = 0.952, 0.0476, 0,$ $1.000,$ for positions 1, 25, 50, and 100, respectively. In this case, the
90 top K-means cluster includes N_{counts} values 0.95 and 1.000 associated with query positions 1
91 and 100, respectively. Query positions 1 and 100 then define the region bounds. The region

92 bounds are subsequently applied to the query sequence, and any overhangs of sufficient length
93 (default: 6 bps) can optionally be classified as their own region (overhangs that are less than the
94 sufficient length threshold are ignored). The default setting for this study was to generate
95 overhangs when possible. The query mapper also defines the region's type: if the region has
96 one or more alignments derived from a protein database, it is defined as a "translated" region; if
97 it only has alignments from a nucleotide database, it is defined as an "untranslated" region; if no
98 alignments are identified, it is defined as a "novel" region. Further, for translated regions, the
99 reading frame(s) is documented based on the alignment. UltraSEQ provides the option to re-
100 align novel regions for greater depth of analysis, but this option was not used in this study. In
101 the example presented here, two regions would be identified: one from query positions 1 to 100,
102 and the second from query positions 101 to 150.

103 **Context Services and Subservices.** These services generate contextual information and
104 passes information to downstream services. The **Metadata Service** maps metadata to
105 alignment results. For UniRef100 alignments, these metadata include Gene Ontology terms,
106 UniProt identifiers, UniRef100 identifiers (which are linked to proteins involved in genetic
107 engineering, housed within Battelle's SoC database), taxonomy identifiers (also linked to
108 Battelle's SoC database for agent metadata), and other. For SoC alignments, these metadata
109 further include tags such as coarse functionality (adherence, antibiotic resistance, etc.),
110 pathways, SoC groups, etc. as defined in (Gemler et al., 2022). For nucleotide alignments,
111 current metadata includes taxonomic identifiers. Other context services available for use but not
112 used in this study are described in the Supplement B - Other UltraSEQ Services.

113 **Rules Engine Service.** This service combines all of the above context and prediction services
114 for regions, sequences, and samples using user defined logic rules for rapid sequence triage.
115 UltraSEQ currently has 4 default rules engines to identify biothreats, controlled sequences for
116 DNA synthesis vendors, indicators of genetic engineering, and Metagenomics Diagnostics. The
117 first three are not described here as they are specific to various use cases. The fourth is
118 described in the main methods of the manuscript.

119 **Metagenomics Service.** This service provides sample level taxonomic composition based on
120 the regions identified from reads processed in the query mapper service in 3 steps: 1) filtering
121 out low quality reads, 2) scoring the remaining reads based on the information content of the
122 reads, and 3) predicting the taxonomic composition based on the scores. In the first step, the
123 default alignment quality filters used in this study include minimum alignment length of 48 base
124 pairs (16 amino acids set based on aligner seed length), 99% percent identity and 100% region
125 coverage for nucleotide alignments (note: no nucleotide databases were used in this study),
126 95% percent identity and 90% region coverage for protein alignments.

127 The metagenomics service works by estimating the information content of a read. That is, reads
128 that are unique to a protein from a specific organism contain the highest amount of information,
129 whereas reads that are found in proteins from across the tree of life contain less information.
130 The information content of a read is derived from the read's alignment data, in which the value
131 of its information content is inversely proportional to the product of the number of unique
132 accessions and taxonomies associated with high-quality alignments of a region – i.e., a region
133 that contains a single accession and taxonomy call is more useful than a region that contains

134 many accession and taxonomy calls. We note that the default protein reference database used
135 in this study, the UniRef100, clusters proteins with 100% similarity to each other into a single
136 reference accession. This clustering feature is important for the metagenomics service's
137 efficacy, since it prevents reference database duplication from incorrectly lowering the perceived
138 information content of a region (e.g., duplicates of the same protein are represented by a single
139 UniRef100 cluster, which would appear as a single subject accession in this study).

140 Sequence region-level taxonomy predictions are associated with confidence scores that are
141 based on alignment quality. For each unique taxonomy identified, the maximum confidence
142 score from alignments that are associated with it are assigned. Specifically, based on the results
143 of the query mapper service, all region alignments are compiled in a table ("query sequence
144 information table"), and scoring is initially performed on a per region, per agent (organism), per
145 accession basis. More specifically, each region (r), agent (a), and accession (acc) combination
146 is assigned the following score, $S_{a,r,acc}$:

$$147 \quad S_{a,r,acc} = \frac{Aqual_{acc,r}}{Na_r \times Nacc_r}$$

148 Where $Aqual_{acc,r}$ is the alignment quality (percent coverage x percent identity) in the region,
149 Na_r is the number of unique agents associated with the subject accessions in the region (score
150 is inversely proportional to the region's uniqueness) and $Nacc_r$ is the number of accessions
151 from the subject database that are associated with the region (score is inversely proportional to
152 the region's sequence complexity – higher complexity implies more specificity to a specific
153 protein). Subsequently, an agent region score, $S_{a,r}$, is calculated to be the score associated with
154 the highest scoring accession (or accessions in the case of a tie) for the given agent, region
155 combination:

$$156 \quad S_{a,r} = \max_{acc} S_{a,r,acc}$$

157 For each unique taxonomy across the sample, the agent scores $S_{a,r}$ are summed across all
158 regions for which each taxonomy is associated, and the sample-level or agent score (S_a) is
159 calculated:

$$160 \quad S_a = \sum_r S_{a,r}$$

161 At this point, the agent score (S_a) are rank ordered, and starting with highest sample-level
162 scoring taxonomy, all sequences associated with the highest scoring taxonomy are identified
163 and all other taxonomies associated with those sequences are removed from the query
164 sequence information table (as defined above). This process is iteratively repeated until all
165 taxonomies have been processed. The result is a pruned list of agent scores (S_a) and their
166 associated TaxIDs. From this list, a K-means cluster of the agent scores is performed by
167 domain (Bacteria, Archaea, Eukaryotes, and Viruses) using the same method as described
168 above for the Query Mapper Service, and the taxonomies associated with the top cluster in
169 each domain are set to be the final sample composition. As with the query mapper, the
170 algorithm checks if the top two clusters' centroids are within a particular tolerance, referred to as

171 the metagenomic clustering threshold (MCT) in the main body of the manuscript; if they are, the
172 penultimate cluster is absorbed into the top cluster (otherwise, the top cluster remains
173 unaltered). In this case, a 50% MCT was used for all UltraSEQ runs except during testing
174 phases as described in the Results Section. The final confidence associated with each agent,
175 C_a , is defined as the average alignment quality for all sequences used in the final (pruned) query
176 sequence information table for that agent. Note, due to the high abundance of phages, all
177 TaxIDs associated with phages and other similar non-human viruses were masked from these
178 calculations. This masking was accomplished by creating a removal list of all NCBI viral TaxID
179 associated with the following hosts: fungi, bacteria, algae, archaea, diatom, and protozoa.

180 **Reporting Services.** UltraSEQ provides several reports as well as described below. The text
181 below describes the details of these reports as of the writing of this manuscript, although we
182 anticipate additions/modifications as appropriate. Details for the *Top Alignment Report*,
183 *Taxonomy Report*, and *Default Report* are provided in the main section of the manuscript.
184 Additional details for the sample report are provided here.

185 **Sample Report.** As described in the main methods section of the manuscript, the ‘main report’
186 tab provides a list of all organisms identified from the above Metagenomics Service, the results
187 associated with the identified organisms, and the metadata associated with the organism from
188 Battelle’s SoC database. These results and metadata are used in a logical diagnostic rules
189 engine described in the ‘trigger-summary’ tabs. Statistics for the UltraSEQ run are provided in
190 the ‘sample-statistics’ tab. For each organism identified, the ‘VF’ tab provides a list of SoCs
191 identified, including virulence factors and antibiotic resistance genes from the Comprehension
192 Antibiotic Resistance Database (CARD) (Alcock et al., 2020). Specifically, if an alignment to one
193 of the proteins in the SoC database is contained within the Top Alignment Report and its TaxID
194 matches to the organism or one of its children, it is populated in the organism-specific ‘VF’ tab.
195 For antibiotic resistance profiles, only proteins in the CARD’s protein homology model are
196 currently used. These protein sequences are currently populated in Battelle’s SoC database, but
197 some metadata associated with these sequences and the drugs they confer resistance to
198 defined by CARD (e.g., in the aro.tsv and ro.tsv files provided by CARD downloads
199 <https://card.mcmaster.ca/download>). The ‘ABR’ tab pulls results from the ‘VF’ tab to provide
200 antibiotic resistance information. Information in this tab is organized by drug class and antibiotic
201 for easy interpretation. Other antibiotic resistance models (e.g., protein variant model, rRNA
202 gene variant model, and protein knockout model, etc.) are not currently used in UltraSEQ. Thus,
203 antibiotic resistance profiles are currently based solely on presence of genes that confer
204 antibiotic resistance (i.e., profiles are not based on point mutants that may help confer
205 resistance). In addition to the organism-specific antibiotic resistance profile, an organism-
206 specific agnostic profile is provided in the ‘CARD SoCs Report’ tab to further aid in antibiotic
207 resistance genotyping (in cases where antibiotic genes may map to the incorrect or many
208 different taxonomies).

209

210 **References for Supplemental Section A**

- 211 1. Alcock, B. P., Raphenya, A. R., Lau, T. T. Y., Tsang, K. K., Bouchard, M., Edalatmand, A., Huynh,
212 W., Nguyen, A. V., Cheng, A. A., Liu, S., Min, S. Y., Miroshnichenko, A., Tran, H. K., Werfalli, R. E.,
213 Nasir, J. A., Oloni, M., Speicher, D. J., Florescu, A., Singh, B., . . . McArthur, A. G. (2020). CARD
214 2020: antibiotic resistance surveillance with the comprehensive antibiotic resistance database.
215 *Nucleic Acids Res*, 48(D1), D517-D525. <https://doi.org/10.1093/nar/gkz935>
216 2. Andrews, S. (2010). *FastQC: A Quality Control Tool for High Throughput Sequence Data [Online]*.
217 Available online at: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
218 3. Arthur, D., & Vassilvitskii, S. (2007). K-means++: the advantages of careful seeding. SODA'07
219 Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms.
220 CityPhiladelphia, StatePA: placecountry. In: region SIAM Press.
221 4. Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina
222 sequence data. *Bioinformatics*, 30(15), 2114-2120.
223 <https://doi.org/10.1093/bioinformatics/btu170>
224 5. Chen, S., Zhou, Y., Chen, Y., & Gu, J. (2018). fastp: an ultra-fast all-in-one FASTQ preprocessor.
225 *Bioinformatics*, 34(17), i884-i890. <https://doi.org/10.1093/bioinformatics/bty560>
226 6. Ewels, P., Magnusson, M., Lundin, S., & Kaller, M. (2016). MultiQC: summarize analysis results
227 for multiple tools and samples in a single report. *Bioinformatics*, 32(19), 3047-3048.
228 <https://doi.org/10.1093/bioinformatics/btw354>
229 7. Gemler, B. T., C., M., A., H. C., D., H., Z., S., J., H. L., O., T., & C., B. (2022). Function-based
230 Classification of Hazardous Biological Sequences: Demonstration of a New Paradigm for
231 Biohazard Assessments (Submitted for Publication).
232 8. Head, S. R., Komori, H. K., LaMere, S. A., Whisenant, T., Van Nieuwerburgh, F., Salomon, D. R., &
233 Ordoukhanian, P. (2014). Library construction for next-generation sequencing: overviews and
234 challenges. *Biotechniques*, 56(2), 61-64, 66, 68, passim. <https://doi.org/10.2144/000114133>
235 9. Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat Methods*,
236 9(4), 357-359. <https://doi.org/10.1038/nmeth.1923>
237 10. Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18),
238 3094-3100. <https://doi.org/10.1093/bioinformatics/bty191>
239 11. Lloyd, S. (1982). Least squares quantization in PCM. *IEEE transactions on information theory*,
240 28(2), 129-137.
241 12. Ng, A. (2012). Clustering with the k-means algorithm. *Machine Learning*, 1-2.
242 13. Wick, R. R., Judd, L. M., Gorrie, C. L., & Holt, K. E. (2017). Completing bacterial genome
243 assemblies with multiplex MinION sequencing. *Microb Genom*, 3(10), e000132.
244 <https://doi.org/10.1099/mgen.0.000132>
245 14. Yang, L., Haidar, G., Zia, H., Nettles, R., Qin, S., Wang, X., Shah, F., Rapport, S. F., Charalampous,
246 T., Methe, B., Fitch, A., Morris, A., McVerry, B. J., O'Grady, J., & Kitsios, G. D. (2019).
247 Metagenomic identification of severe pneumonia pathogens in mechanically-ventilated
248 patients: a feasibility and clinical validity study. *Respir Res*, 20(1), 265.
249 <https://doi.org/10.1186/s12931-019-1218-4>

250

251

252

253 **Supplement B: Other UltraSEQ Services**

254 **Other Context Services and Subservices.** Other services include a **Genetic Engineering**
255 **(GE) Service** and a **Classifier Service**. The **GE Service enables** prediction of GE indicators,
256 including microservices for detection of GE proteins, GE signatures, codon optimization, and
257 codon re-coding. The Classifier Service includes artificial intelligence (AI) models to make
258 alignment-free predictions on amino acid sequences; the output is the probability that the input
259 is associated with a subset of threat metadata categories (the coarse functional categories)
260 described in Gemler et al. (Gemler et al., 2022). Information from the context services are
261 passed to the Prediction Services and Flagging System (Rules Engine) as described below.

262 **Region-based Taxonomy Prediction Subservice.** For other applications (such as forensic
263 applications), sequence and region-level taxonomic information can be useful. For this
264 prediction, a conservative and information-based approach is used that takes into consideration
265 strength of alignments, the number of times a TaxID appears across alignments, and the
266 taxonomic depth (species, genus, etc.). For this prediction, the TaxID frequency that each taxID
267 appears across the top alignments is calculated, and TaxID Depth is assigned as follows: 100
268 for species and below, 75 for genus, 50 for family, 30 for order, 20 for class, 15 for phylum, and
269 10 for domain and above. A normalized TaxID Depth is then calculated in the same manner as
270 the normalization defined in Equation 1 (**Supplemental Material A**). The TaxID score for each
271 TaxID is then calculated according to Eq 2.

272 Equation 2

$$273 \text{TaxID Score} = \sqrt{w_d (\text{Normalized TaxID Depth})^2 + w_f (\text{TaxID Frequency})^2}$$

274 Where default weight values: $w_d = 2.0$ and $w_f = 1.0$ are used (optimized weights based on
275 test/validation datasets, not shown). The final taxonomy predictions are then based on a 2-D K-
276 means clustering for the alignment confidence and TaxID score data. The TaxIDs in the top
277 cluster are considered the final predictions using the same K-means clustering methods as
278 described above (the rest of the TaxID predictions are discarded. Further, the confidence
279 associated with each TaxID prediction is reported as follows: taxonomy evidence is gathered for
280 the region from the alignments and the alignment scores (percent identity x percent region
281 coverage) are normalized by the max heuristic value (100 in the case of 100% identity over
282 100% of the region). For each TaxID identified, this normalized score is considered the final
283 confidence score associated with each TaxID.

284 **Region-based Function Prediction Subservice.** For each region, the function (gene ontology
285 terms) is calculated in a similar manner. Specifically, 1) function evidence is gathered for each
286 Region and alignment scores are normalized, 2) a 1-D K-means clustering is used for the
287 alignment confidence; the GO Term sets in the top cluster are selected, and 3) the final function
288 prediction confidence is calculated by averaging the alignment confidence values across the
289 members in the top cluster.

290 **Region-based Threat Prediction Subservice.** For each region, the threat metadata
291 associated with that region (damage, antibiotic resistance, adherence, etc. as defined in

292 (Gemler et al., 2022)) is tabulating the SoC alignment scores (percent identity x region
293 coverage) associated with threat metadata category, clustering the scores using K-means,
294 down selecting to only the top cluster, adding the alignment scores to 1/100th of the AI model,
295 then dividing by the maximum possible score.

296

297 **References for Supplemental Section B**

- 298 1. Gemler, B. T., C., M., A., H. C., D., H., Z., S., J., H. L., O., T., & C., B. (2022). Function-based
299 Classification of Hazardous Biological Sequences: Demonstration of a New Paradigm for
300 Biohazard Assessments (Submitted for Publication).

301

302

Bin	Disease Diagnosis Confidence	Qualifications
Bin 1A	Highest Confidence Agent	<ol style="list-style-type: none"> 1. SoC Filter* is True AND 2. SoC agent is a pathogen that infects human host AND 3. Relative abundance filters (based on predicted reads): 1%-5% for bacteria**; > 1%-2%** for fungi; >1% for other (protozoa, etc.) AND 4. SoC agent is contained in Battelle's human respiratory pathogen list (for respiratory datasets) or encephalitis/ meningitis pathogen list (for encephalitis datasets) AND*** 5. At least one SoC with the Active, Damage, Apoptosis, Inhibits, or Transmission threat category from the agent was used by UltraSEQ's metagenomics module for agent prediction
Bin 1B	High Confidence Agent	<ol style="list-style-type: none"> 1. SoC Filter* is False AND Conditions 2,3,4, and 5 above met
Bin 2	Medium Confidence Respiratory Agent	<ol style="list-style-type: none"> 1. Condition 1A and Condition1B = FALSE AND Conditions 2 and 3 above met AND 2. No SoCs identified from the above categories for that agent
Bin 3	Lowest Confidence Respiratory Agent	<ol style="list-style-type: none"> 1. Condition 1A and Condition1B = FALSE AND Conditions 2 and 3 above met AND 2. There are SoCs from the above categories for that agent, however none are found [Note: other SoCs may be identified such as antibiotic resistant SoCs and adherence SoCs]

304 * SoC Filter is a condition that is true when the UltraSEQ metagenomics service uses a
 305 UniRef100 cluster containing a SoC to trigger the taxonomy prediction.

306 ** For bacteria, a 5% threshold was used for the de Vries et al., PRJNA516289, Hasan et al.,
 307 PRJEB7888, PRJEB13360; a 1% bacteria filter was used for all other datasets; for fungi, a 2%
 308 filter was used PRJNA516582; a 1% fungi filter was used for all other datasets

309 *** At the time of this manuscript, Battelle's SoC database contained ~2,200 human pathogen
 310 species, ~150 of which are curated as potential contaminants (either from reagents used during
 311 sequencing and/or due to the biological sample such as normal skin flora; all of these
 312 annotations are provided in the sample report. Of the human pathogens, ~250 are contained
 313 within the encephalitis/ meningitis list and ~250 are contained within the respiratory list.

314

315

316 **Supplement F - Supplemental Results**
 317 **Encephalitis / meningitis**
 318 **PRJNA516289 (Miller et al. (Miller et al., 2019)).**

319 **Table F1. UltraSEQ Results for Miller Dataset**

Result	Parasites		Fungi		Bacteria		DNA Viruses		RNA Viruses	
	UltraSE Q	Mille r	UltraSE Q	Mille r	UltraSE Q	Mill er	UltraSE Q	Mille r	UltraSE Q	Mille r
TP	1	1	10	9	6	5	23	25	10	11
FP	0	0	0	0	0	1	0	0	0	0
FN	0	0	0	1	1	2	5	3	3	2
TN	4	4	38	38	51	50	16	16	5	5
PPA	100%	100 %	100%	90%	86%	71 %	82%	89%	77%	85%
NPA	100%	100 %	100%	100 %	100%	98 %	100%	100 %	100%	100 %
Accura cy	100%	100 %	100%	98%	98%	95 %	87%	93%	83%	89%

320 * As noted by Miller et al., the “truth” was considered the initial clinical result unless a confirmatory test was run (i.e., if
 321 a confirmatory test was run by Miller et al., the truth was considered to be the confirmatory test). For both SURPI and
 322 UltraSEQ, RNA viruses were reported using sequences derived from the RNA libraries, whereas all other organisms
 323 results were based on sequences from the DNA libraries.

324
 325 **PRJNA516582 (Saha et al. (Saha et al., 2019)).**

326 **Table F2: Summary of Results for Saha and UltraSEQ**

Result	Saha Results				UltraSEQ Results			
	All samples	Culture Only	All confirmed Cases	CHIKV cases	All samples	Culture Only	All confirmed Cases	CHIKV cases
TP	52	7	24	17	53	7	25	17
FP	NR*	NR*	NR*	NR*	0	0	0	0
FN	12	1	12	0	11	1	11	0
TN	29	N/A	N/A	N/A	29	N/A	N/A	N/A
PPA	81%	88%	67%	100%	83%	88%	72%	100%
NPA	N/A	N/A	N/A	N/A	100%	N/A	N/A	N/A
ACC	87%	88%	67%	100%	88%	88%	72%	100%

327 * NR= not reported; N/A=not applicable (i.e., could not be calculated)

328 ** As detailed in the methods, UltraSEQ identified *E. coli* in nearly every sample despite the fact
 329 that *E. coli* was only identified by clinical tests in 2 samples. By using the UltraSEQ logic as
 330 defined in the Methods section without any background sample subtraction (as required by
 331 Saha), UltraSEQ was able to remove all *E. coli* false positives.

332

333 'CSF_metagenomics' from idseq.net (Hasan et al. (Hasan et al., 2020)).

334 Table F3: Table of Species Identified by UltraSEQ for Sample CW322

Taxonomy Name	TaxID	NCBI TaxID Rank	Type	Confidence	Relative Abundance TaxID + Children (vs PREDICTED ONLY reads)
<i>Neisseria meningitidis</i>	487	species	Bacteria	99.4	99.96
Human alphaherpesvirus 2	10310	species	Virus	99.0	0.0025

335

336

337

CW322 Share Download

Sample Details

Metagenomic

Taxon name Name Type: Scientific Background: CSF_Metagenomics_BG Categories: 3 Threshold filters Read Specificity: All

Bacteria Viruses Phage

178 rows passing the above filters, out of 208 total rows. Clear All Filters

Taxon	Score	Z Score	rPM	r	contig	contig r	%id	L	E value	NT NR
> Neisseria (22 bacterial species: ● 2)	158,368,898	99.0 99.0	8,437.2 8,276.6	128,225 125,784	1,227 1,226	119,585 119,545	99.1 98.9	2,851.3 375.8	10 ⁻²⁹² 10 ⁻²⁵⁹	
> Cutibacterium (2 bacterial species)	-0	-0.2 -0.2	4.1 2.7	62 41	0 0	0 0	99.5 99.9	155.3 54.7	10 ⁻⁹⁵ 10 ⁻²⁹⁷	
> Morococcus (1 bacterial species)	26,320	0.0 100.0	0.0 2.6	0 40	0 1	0 40	0.0 97.5	0.0 157.0	0 10 ⁻¹⁰⁵	
> Escherichia (1 bacterial species: ● 1)	-0	-0.2 -0.2	2.0 1.9	30 29	0 0	0 0	99.9 98.9	306.3 91.2	10 ⁻²⁰⁸ 10 ⁻²¹¹	
> Pseudomonas (14 bacterial species)	3	-0.2 -0.2	1.1 1.1	16 16	0 0	0 0	99.4 100.0	186.7 63.8	10 ⁻¹²¹ 10 ⁻²⁷⁹	
> Staphylococcus (4 bacterial species: ● 1)	0	-0.2 -0.2	0.8 1.0	12 15	0 1	0 4	99.9 99.8	148.2 76.7	10 ⁻⁸⁹ 10 ⁻²⁵²	
> all taxa with neither family nor genus classification (7 species)	6	-0.2 -0.2	0.7 0.9	10 14	2 0	8 0	100.0 88.8	423.7 87.9	10 ⁻²⁷⁷ 10 ⁻²⁰⁵	
> Clostridioides (1 bacterial species: ● 1)	2,133	-100.0 52.0	0.0 0.8	0 12	0 0	0 0	0.0 96.2	0.0 36.8	0 10 ⁻³⁰⁸	
> Propionibacterium (6 bacterial species)	0	-0.2 -0.2	0.3 0.5	4 7	0 0	0 0	99.9 100.0	202.0 51.9	10 ⁻¹²⁶ 10 ⁻³⁰⁸	
▼ Simplexvirus (1 viral species: ● 1)	2,916	10.8 100.0	0.3 0.3	4 4	0 0	0 0	99.8 99.7	216.0 64.0	10 ⁻¹³⁵ 10 ⁻³⁰⁸	
Human alphaherpesvirus 2 Known Pathogen	2,916	100.0 100.0	0.3 0.3	4 4	0 0	0 0	99.8 99.7	216.0 64.0	10 ⁻¹³⁵ 10 ⁻³⁰⁸	
> Mycoplasma (2 bacterial species)	61	-0.2 15.2	0.1 0.3	2 4	0 1	0 4	95.8 98.5	24.0 66.0	10 ⁻³ 10 ⁻⁴⁰	
> Francisella (1 bacterial species)	-0	-100.0 -0.2	0.0 0.3	0 4	0 0	0 0	0.0 100.0	0.0 33.0	0 10 ⁻³⁰⁸	
> Zhizhongheella (1 bacterial species)	0	0.0 0.4	0.0 0.1	0 2	0 0	0 0	0.0 90.5	0.0 47.0	0 10 ⁻³⁰⁸	

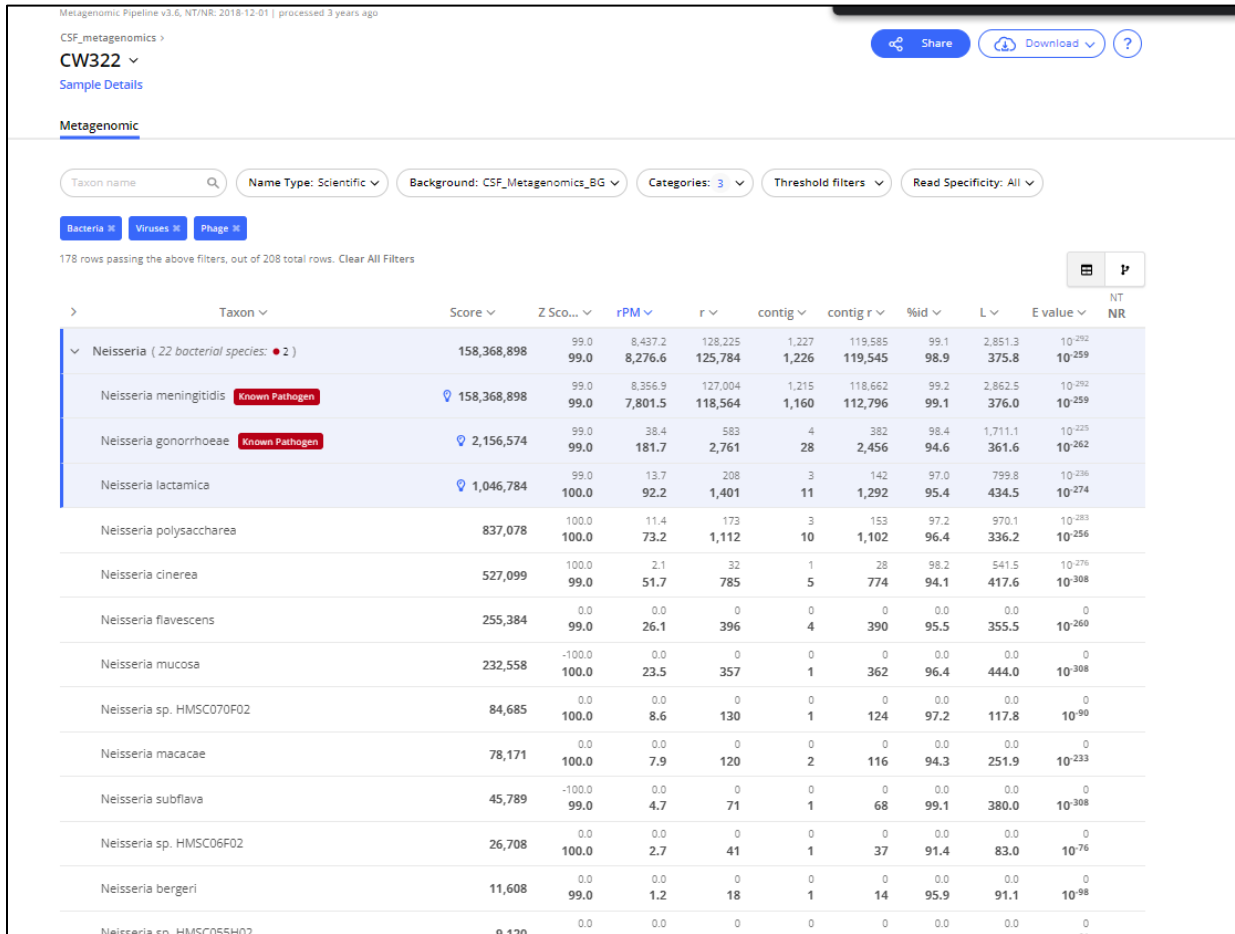
338

339 **Figure F1: CW322 Results Showing Bacteria and Viruses Identified by IdSeq (Note that several**
 340 **more rows of genera were identified and not shown here).**

341

342

343



344 **Figure F2: CW322 Results Showing all Neisseria Species Identified (Note that several more**
 345 **rows of species were identified and not shown here).**

347

348 **Respiratory disease: Influenza**

349 **PRJEB7888 (Fischer et al. (Fischer et al., 2015)).**

350 **Table F4: Summary of Results for Fischer and UltraSEQ**

Pipeline	TP	FN	TN	FP	PPA	NPA	Accuracy
Explify	16	3	5	0	84%	100%	83%
Fisher	15	4	5	0	79%	100%	88%

UltraSE Q	14	5	5	0	74%	100%	79%
--------------	----	---	---	---	-----	------	-----

351

352

353 **Respiratory disease: ventilator associated pneumonia (VAP)**

354 ***PRJNA554856 (Watts et al. (Watts et al., 2019)).***

355

Table F5: AbR Report for SRR9693434 (Patient 2, Day 1)

fluoroquinolone antibiotic (ARO:0000001)		['major facilitator superfamily (MFS) antibiotic efflux pump (ARO:0010002)', 'Staphylococcus aureus norA', ['ARO:3004667'], 31, 0.016, '0'] ['major facilitator superfamily (MFS) antibiotic efflux pump (ARO:0010002)', 'Staphylococcus aureus norA', ['ARO:3004667'], 31, 0.016, '0'] ['major facilitator superfamily (MFS) antibiotic efflux pump (ARO:0010002)', 'Staphylococcus aureus norA', ['ARO:3004667'], 31, 0.016, '0'] ['major facilitator superfamily (MFS) antibiotic efflux pump (ARO:0010002)', 'Staphylococcus aureus norA', ['ARO:3004667'], 31, 0.016, '0']	ciprofloxacin (ARO:0000036) enoxacin (ARO:0000023) ofloxacin (ARO:3000663) norfloxacin (ARO:3000662)	['Staphylococcus aureus norA', ['ARO:3004667'], 31, 0.016, '0']	['major facilitator superfamily (MFS) antibiotic efflux pump (ARO:0010002)', 'Staphylococcus aureus norA', ['ARO:3004667'], 31, 0.016, '0']
phosphonic acid antibiotic (ARO:0000025)		['fosfomycin thiol transferase (ARO:3000133)', 'FosD', ['ARO:3004674'], 1, 0.001, '0']			

glycylcycline (ARO:0000042)	['mepA', [ARO:3000026'], 44, 0.023, '0']	['multidrug and toxic compound extrusion (MATE) transporter (ARO:3000112)', 'mepA', [ARO:3000026'], 44, 0.023, '0']	tigecycline (ARO:0000030)	['mepA', [ARO:3000026'], 44, 0.023, '0']	['multidrug and toxic compound extrusion (MATE) transporter (ARO:3000112)', 'mepA', [ARO:3000026'], 44, 0.023, '0']
lincosamide antibiotic (ARO:0000017)		['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'poxTA', [ARO:3004470'], 6, 0.003, '0'] ['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'vgaA', [ARO:3002829'], 1, 0.001, '0'] ['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'vgaC', [ARO:3002831'], 1, 0.001, '0']			
macrolide antibiotic (ARO:0000000)		['macrolide phosphotransferase (MPH) (ARO:3000333)', 'mphC', [ARO:3000319'], 72, 0.037, '0'] ['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'poxTA', [ARO:3004470'], 6, 0.003, '0'] ['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'vgaA', [ARO:3002829'], 1, 0.001, '0'] ['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'vgaC', [ARO:3002831'], 1, 0.001, '0']	spiramycin (ARO:3000156) clarithromycin (ARO:0000065) roxithromycin (ARO:0000027) tylosin (ARO:3000145) oleandomycin (ARO:3000867) azithromycin (ARO:3000158) erythromycin (ARO:0000006) dirithromycin (ARO:3000176) telithromycin (ARO:0000057)	['mphC', [ARO:3000319'], 72, 0.037, '0']	['macrolide phosphotransferase (MPH) (ARO:3000333)', 'mphC', [ARO:3000319'], 72, 0.037, '0']
mupirocin (ARO:3000554)		['antibiotic-resistant isoleucyl-tRNA synthetase (ileS) (ARO:3000446)', 'mupA', [ARO:3000521'], 10, 0.005, '1'] ['antibiotic-resistant isoleucyl-tRNA synthetase (ileS) (ARO:3000446)', 'mupB',	mupirocin (ARO:3000554)	['mupA', [ARO:3000521'], 10, 0.005, '1'] ['mupB', [ARO:3000510'], 2, 0.001, '1']	['antibiotic-resistant isoleucyl-tRNA synthetase (ileS) (ARO:3000446)', 'mupA', [ARO:3000521'], 10, 0.005, '1'] ['antibiotic-resistant isoleucyl-tRNA synthetase (ileS) (ARO:3000446)',

		['ARO:3000510'], 2, 0.001, '1']			'mupB', ['ARO:3000510'], 2, 0.001, '1']
oxazolidinone antibiotic (ARO:3000079)		['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'poxtA', ['ARO:3004470'], 6, 0.003, '0'] ['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'vgaA', ['ARO:3002829'], 1, 0.001, '0'] ['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'vgaC', ['ARO:3002831'], 1, 0.001, '0']	linezolid (ARO:0000072)	['poxtA', ['ARO:3004470'], 6, 0.003, '0']	['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'poxtA', ['ARO:3004470'], 6, 0.003, '0']
penam (ARO:3000008)		['methicillin resistant PBP2 (ARO:3001208)', 'mecA', ['ARO:3000617'], 35, 0.018, '0']	methicillin (ARO:0000015)	['mecA', ['ARO:3000617'], 35, 0.018, '0']	['methicillin resistant PBP2 (ARO:3001208)', 'mecA', ['ARO:3000617'], 35, 0.018, '0']
phenicol antibiotic (ARO:3000387)		['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'poxtA', ['ARO:3004470'], 6, 0.003, '0'] ['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'vgaA', ['ARO:3002829'], 1, 0.001, '0'] ['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'vgaC', ['ARO:3002831'], 1, 0.001, '0']	chloramphenicol (ARO:3000385) florfenicol (ARO:3000461)	['poxtA', ['ARO:3004470'], 6, 0.003, '0']	['ABC-F ATP-binding cassette ribosomal protection protein (ARO:3004469)', 'poxtA', ['ARO:3004470'], 6, 0.003, '0']

356

357 **PRJNA554461 (Yang et al. (Yang et al., 2019)).**

358

359 **Table F6.** Comparison of UltraSEQ and WIMP Results for the Yang et al. (PRJNA554461) VAP
360 Dataset

Platform	TPs	FNs	TNs	FPs	PPA	NPA	Accuracy
----------	-----	-----	-----	-----	-----	-----	----------

361
362

					[TP / (TP+FN)]	[TN / (TN +FP)]	
WIMP (Author)	10	2	6	7	83%	46%	64%
UltraSEQ	9	3	5	3	75%	63%	70%

Table F7. UltraSEQ's Antibiotic Genotype Profiles Agree with Phenotypic Profiles.

Case	AbR profile by culture*	UltraSEQ AbR summary for Identified Pathogen w/ # reads** Drug class Antibiotic	Author results
1	R: ticarcillin/ clavulanic acid R: ceftazidime I: levofloxacin NT: Tetracycline	Cephalosporin (beta lactam): 2 reads Note: Others identified as well; fluoroquinolone identified in agent agnostic report	blaTEM-4,blaTEM-112, blaTEM-157, blaACT-5, oqxB, tetC
2	R: methicillin R: erythromycin, clindamycin	Penam (beta lactam): 2 reads methicillin: 10 reads Macrolide: Roxithromycin, oleandomycin, telithromycin, spiramycin, azithromycin, clarithromycin, erythromycin, tylosin, dirithromycin: 17 reads lincosamide: clindamycin, lincomycin: 17 reads Note: several others identified as well	mecA ermA, erm tet38,ant(4')-Ib, tetC, blaTEM-4
3	I: tetracycline	Tetracycline: Tetracycline: 3 reads Tigecycline: 3 reads No Methicillin resistance identified Note: several others identified as well but all 2 reads or less; included all strains of S. aureus identified	tetK, tet38, tetQ
4	R: tetracycline R: Trimethoprim-sulfamethoxazole R: ciprofloxacin, levofloxacin	Tetracycline: Tetracycline: 14 reads Minocycline, demeclocycline, oxytetracycline, chlortetracycline, doxycycline: 4 reads Diaminopyrimidine: Trimethoprim: 1 read	tetX sul1 dfrA acrF, pare, mdf mphA, aadA5, vgaC, blaACT-5, blaACT-14, mefA, mel

		Note: Aminoglycosides and macrolides identified as well (>10 reads each); fluoroquinolone identified in agent agnostic report	
5	S: all tested agents	None No Methicillin resistance identified	None
6	S: all tested agents	Note: 7 classes identified, all with 2 reads or less No Methicillin resistance identified	Tet38, blaTEM4
7	S: all tested agents	N/A	None
8	NT	Note: 56 reads to drug efflux protein conferring resistance to multiple drug classes	tetM, isaC, sul1, tetQ, mphA, aadA5

363 * R=resistant, I=intermediate, S=Susceptible, NT=Not tested; N/A = not applicable

364 ** Only appropriate true positives and true negatives are listed (full AbR phenotype is unknown).

365 Results in **green** font indicate that for the identified pathogen, UltraSEQ identified the same
366 antibiotic or class as the phenotype data; those in **blue** denote that UltraSEQ identified a closely
367 related class; those in **orange** indicate that that the antibiotic was only identified in the agent
368 agnostic report; those in *italics* were not phenotypically tested.

369

370 **Illumina RNASeq Dataset: Respiratory viruses**

371 **PRJEB13360 (Graf, Flygare (Flygare et al., 2016; Graf et al., 2016)).** Detailed results are
372 provided in **Supplement D – Supplemental_File_Scores.xlsx.**

373 **Illumina RNASeq Dataset: nasopharyngeal swabs for SARS-CoV-2 diagnosis**

374 **PRJNA634356 (Babiker et al (Babiker et al., 2020)).**

375

376 **Table F8.** Comparison of UltraSEQ Results to Babiker et al. (KrakenUniq) and Explify

Platform	TPs	FNs	TNs	FPs	PPA [TP / (TP+FN)]	NPA [TN / (TN +FP)]	Accuracy
Author (KrakenUniq)	26	1	17	1	96%	94%	96%
Explify	20	3	16	0	86%	100%	93%
UltraSEQ	27	0	16	2	100%	89%	96%

377

378 **Mixed:**

379 **de Vries et al. Dataset.** (<https://veb.lumc.nl/CliniMG>)

380 **Table F9.** UltraSEQ Results for the de Vries Dataset Compared to Other Pipelines as Reported
 381 in (de Vries et al., 2021)

Pipeline	Positive predictive value (PPV) [%]*	PPA (Sensitivity) [%]*
UltraSEQ	100	92
Centrifuge	100	92
DAMIAN	100	77
DIAMOND	93	85
DNASTAR	71	100
FEVIR	88	100
Genome Detective	100	85
Jovian	100	77
MetaMIC	100	77
metaMix	100	100
One Codex	100	77
RIEMS	81	85
Taxonomer	100	85
VirMet	93	92

382 * PPV and PPA for UltraSEQ results were determined as described in the Methods Section.
 383 PPA and PPV for all other datasets determined as reported in Supplemental Table 2 and 4,
 384 respectively in (de Vries et al., 2021).

385

386 **Mixed Illumina RNASeq Dataset: In-house COVID-19 Saliva Study**

387 **Battelle (PRJNA856680).**

388 Detailed results provided in Supplement D – Supplemental_File_Scores.xlsx.

389

390 **References for Supplemental Section F**

- 391 1. Babiker, A., Bradley, H. L., Stittleburg, V. D., Ingersoll, J. M., Key, A., Kraft, C. S., Waggoner, J. J., &
392 Piantadosi, A. (2020). Metagenomic Sequencing To Detect Respiratory Viruses in Persons under
393 Investigation for COVID-19. *J Clin Microbiol*, 59(1). <https://doi.org/10.1128/JCM.02142-20>
- 394 2. de Vries, J. J. C., Brown, J. R., Fischer, N., Sidorov, I. A., Morfopoulou, S., Huang, J., Munnink, B.
395 B. O., Sayiner, A., Bulgurcu, A., Rodriguez, C., Gricourt, G., Keyaerts, E., Beller, L., Bachofen, C.,
396 Kubacki, J., Samuel, C., Florian, L., Dennis, S., Beer, M., . . . Claas, E. C. J. (2021). Benchmark of
397 thirteen bioinformatic pipelines for metagenomic virus diagnostics using datasets from clinical
398 samples. *J Clin Virol*, 141, 104908. <https://doi.org/10.1016/j.jcv.2021.104908>
- 399 3. Fischer, N., Indenbirken, D., Meyer, T., Lutgehetmann, M., Lellek, H., Spohn, M., Aepfelbacher,
400 M., Alawi, M., & Grundhoff, A. (2015). Evaluation of Unbiased Next-Generation Sequencing of
401 RNA (RNA-seq) as a Diagnostic Method in Influenza Virus-Positive Respiratory Samples. *J Clin*
402 *Microbiol*, 53(7), 2238-2250. <https://doi.org/10.1128/JCM.02495-14>
- 403 4. Flygare, S., Simmon, K., Miller, C., Qiao, Y., Kennedy, B., Di Sera, T., Graf, E. H., Tardif, K. D.,
404 Kapusta, A., Rynearson, S., Stockmann, C., Queen, K., Tong, S., Voelkerding, K. V., Blaschke, A.,
405 Byington, C. L., Jain, S., Pavia, A., Ampofo, K., . . . Schlager, R. (2016). Taxonomer: an interactive
406 metagenomics analysis portal for universal pathogen detection and host mRNA expression
407 profiling. *Genome Biol*, 17(1), 111. <https://doi.org/10.1186/s13059-016-0969-1>
- 408 5. Graf, E. H., Simmon, K. E., Tardif, K. D., Hymas, W., Flygare, S., Eilbeck, K., Yandell, M., &
409 Schlager, R. (2016). Unbiased Detection of Respiratory Viruses by Use of RNA Sequencing-
410 Based Metagenomics: a Systematic Comparison to a Commercial PCR Panel. *J Clin Microbiol*,
411 54(4), 1000-1007. <https://doi.org/10.1128/JCM.03060-15>
- 412 6. Hasan, M. R., Sundararaju, S., Tang, P., Tsui, K. M., Lopez, A. P., Janahi, M., Tan, R., & Tilley, P.
413 (2020). A metagenomics-based diagnostic approach for central nervous system infections in
414 hospital acute care setting. *Sci Rep*, 10(1), 11194. <https://doi.org/10.1038/s41598-020-68159-z>
- 415 7. Miller, S., Naccache, S. N., Samayoa, E., Messacar, K., Arevalo, S., Federman, S., Stryke, D., Pham,
416 E., Fung, B., Bolosky, W. J., Ingebrigtsen, D., Lorizio, W., Paff, S. M., Leake, J. A., Pesano, R.,
417 DeBiasi, R., Dominguez, S., & Chiu, C. Y. (2019). Laboratory validation of a clinical metagenomic
418 sequencing assay for pathogen detection in cerebrospinal fluid. *Genome Res*, 29(5), 831-842.
419 <https://doi.org/10.1101/gr.238170.118>
- 420 8. Saha, S., Ramesh, A., Kalantar, K., Malaker, R., Hasanuzzaman, M., Khan, L. M., Mayday, M. Y.,
421 Sajib, M. S. I., Li, L. M., Langelier, C., Rahman, H., Crawford, E. D., Tato, C. M., Islam, M., Juan, Y.
422 F., de Bourcy, C., Dimitrov, B., Wang, J., Tang, J., . . . DeRisi, J. L. (2019). Unbiased Metagenomic
423 Sequencing for Pediatric Meningitis in Bangladesh Reveals Neuroinvasive Chikungunya Virus
424 Outbreak and Other Unrealized Pathogens. *mBio*, 10(6). <https://doi.org/10.1128/mBio.02877-19>
- 425 9. Watts, G. S., Thornton, J. E., Jr., Youens-Clark, K., Ponsero, A. J., Slepian, M. J., Menashi, E., Hu,
426 C., Deng, W., Armstrong, D. G., Reed, S., Cranmer, L. D., & Hurwitz, B. L. (2019). Identification
427 and quantitation of clinically relevant microbes in patient samples: Comparison of three k-mer
428 based classifiers for speed, accuracy, and sensitivity. *PLoS Comput Biol*, 15(11), e1006863.
429 <https://doi.org/10.1371/journal.pcbi.1006863>
- 430 10. Yang, L., Haidar, G., Zia, H., Nettles, R., Qin, S., Wang, X., Shah, F., Rapport, S. F., Charalampous,
431 T., Methe, B., Fitch, A., Morris, A., McVerry, B. J., O'Grady, J., & Kitsios, G. D. (2019).
432 Metagenomic identification of severe pneumonia pathogens in mechanically-ventilated
433 patients: a feasibility and clinical validity study. *Respir Res*, 20(1), 265.
434 <https://doi.org/10.1186/s12931-019-1218-4>

