

1 Current progress and future opportunities 2 in applications of **bioinformatics** for 3 biodefense and pathogen detection

4
5 Report from the Winter Mid-Atlantic Microbiome Meet-up, College Park, MD January 10th,
6 2018

7

8 Jacquelyn S. Meisel¹, Daniel J. Nasko¹, Brian Brubach¹, Victoria Cepeda Espinoza¹, Jessica
9 Chopyk², Héctor Corrada-Bravo¹, Marcus Fedarko¹, Jay Ghurye¹, Kiran Javkar¹, Nathan D.
10 Olson^{1,3}, Nidhi Shah¹, Sarah M. Allard², Adam L. Bazinet⁴, Nicholas H. Bergman⁴, Alexis
11 Brown⁵, J Gregory Caporaso⁶, Sean Conlan⁷, Jocelyne DiRuggiero⁸, Samuel P. Forry³, Nur A.
12 Hasan^{1,9}, Jason Kralj³, Paul M. Luethy¹⁰, Donald K. Milton¹¹, Brian D. Ondov^{1,7}, Sarah
13 Preheim¹², Shashikala Ratnayake⁴, Stephanie M. Rogers¹³, M.J. Rosovitz⁴, Eric G. Sakowski¹²,
14 Nils Oliver Schliebs¹⁴, Daniel D. Sommer⁴, Krista L. Ternus¹⁵, Gherman Uritskiy⁸, Sean X.
15 Zhang¹⁶, Mihai Pop¹, Todd J. Treangen^{1*∇}

16

17 ¹ Center for Bioinformatics and Computational Biology, University of Maryland College Park,
18 College Park, MD, USA

19 ² School of Public Health, University of Maryland College Park, College Park, MD, USA

20 ³ Material Measurement Laboratory, National Institute of Standards and Technology,
21 Gaithersburg, MD, USA

22 ⁴ National Biodefense Analysis and Countermeasures Center, Frederick, MD, USA

23 ⁵ Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD, USA

24 ⁶ The Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, AZ, USA

25 ⁷ National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA

26 ⁸ Department of Biology, Johns Hopkins University, Baltimore, MD, USA

27 ⁹ CosmosID Inc., Rockville, MD, USA
28 ¹⁰ Department of Pathology, University of Maryland School of Medicine, Baltimore, MD, USA
29 ¹¹ Maryland Institute for Applied Environmental Health, School of Public Health, University of
30 Maryland College Park, College Park MD, USA
31 ¹² Environmental Health and Engineering, Johns Hopkins University, Baltimore, MD, USA
32 ¹³ B.Next, In-Q-Tel, Inc., Arlington, VA, USA
33 ¹⁴ Department of Computer Science, University of Tübingen, Tübingen, Germany
34 ¹⁵ Signature Science, LLC, Arlington, VA, USA
35 ¹⁶ Division of Medical Microbiology, Department of Pathology, School of Medicine, Johns
36 Hopkins University, Baltimore, MD, USA
37 [▽] Current address: Department of Computer Science, Rice University, Houston, TX, USA

38
39
40
41
42
43
44
45
46
47
48
49
50
51

* Address correspondence to:

Todd Treangen
Department of Computer Science – MS-132
Rice University
P.O. Box 1892
Houston, TX 77005-1892
713-348-4724
treangen@rice.edu

52 **Abstract**

53

54 The Mid-Atlantic Microbiome Meet-up (M³) organization brings together academic, government,
55 and industry groups to share ideas and develop best practices for microbiome research. In
56 January of 2018, M³ held its fourth meeting, which focused on recent advances in biodefense,
57 **specifically those relating to infectious disease**, and the use of metagenomic methods for
58 pathogen detection. Presentations highlighted the utility of next-generation sequencing
59 technologies for identifying and tracking microbial community members across space and time.
60 However, they also stressed the current limitations of genomic approaches for biodefense,
61 including insufficient sensitivity to detect low abundance pathogens and the inability to quantify
62 viable organisms. Participants discussed ways in which the community can improve software
63 usability and shared new computational tools for metagenomic processing, assembly,
64 annotation, and visualization. Looking to the future, they identified the need for better
65 bioinformatics toolkits for longitudinal analyses, improved sample processing approaches for
66 characterizing viruses and fungi, and more consistent maintenance of database resources.
67 Finally, they addressed the necessity of improving data standards to incentivize data sharing.
68 Here, we summarize presentations and discussions from the meeting, identifying areas where
69 microbiome analyses have improved our ability to detect and manage biological threats and
70 infectious disease, as well as gaps of knowledge in the field that require future funding and
71 focus.

72

73 **Keywords**

74

75 Microbiome; Metagenomics; Bioinformatics; Biodefense; Bio-threats; Pathogen Detection;
76 Longitudinal Analysis

77 **Abbreviations**

78

79 CBCB – Center for Bioinformatics and Computational Biology

80 CONSERVE – Center of Excellence at the Nexus of Sustainable Water Reuse, Food, and

81 Health

82 CRISPR – clustered regularly interspaced short palindromic repeats

83 CPU – central processing unit

84 FPGA – field programmable gate array

85 GPU – graphics processing unit

86 IQT – In-Q-Tel, Inc.

87 JHU – Johns Hopkins University

88 M³ – Mid-Atlantic Microbiome Meet-up

89 NAU – Northern Arizona University

90 NBACC – National Biodefense Analysis and Countermeasures Center

91 NGS – next-generation sequencing

92 NHGRI – National Human Genome Research Institute

93 NIH – National Institutes of Health

94 NIST– National Institute of Standards and Technology

95 RAM – random access memory

96 SPH – School of Public Health

97 UMD – University of Maryland

98

99

100 **Introduction**

101

102 Strong public health and biodefense research are essential for the prevention, detection, and
103 management of biological threats and infectious disease. Over the last century, the focus of
104 biodefense research has shifted in response to modern advances in biotechnology. Specifically,
105 a biological revolution is underway, generating promising new gene editing and synthetic
106 biology technologies that may transform modern medicine, but also present a threat to public
107 health if misappropriated [1]. As biotechnology becomes increasingly globalized, it is important
108 that we establish new strategies and tools for infectious disease detection and surveillance that
109 will help us protect against bioterrorism and manage disease outbreaks.

110 Rapid advances in next generation sequencing (NGS) technologies have helped advance
111 biodefense research by enabling the development of new methods for identifying and
112 characterizing pathogens. Amplification and sequencing of the 16S rRNA gene allows for high-
113 throughput detection of prokaryotic communities, while shotgun metagenomic sequencing
114 approaches capture the composition and functional potential of multi-domain populations.
115 Metagenomic analyses used for pathogen detection and identification are often time sensitive.
116 The results help inform high-stakes decision-making, such as choosing an appropriate medical
117 treatment, deciding if a food product should be recalled due to contamination, or determining if
118 an area should be shut down due to a suspected act of bioterrorism. In addition, geospatial and
119 temporal metagenomic analyses are essential for tracking the dynamic responses of microbial
120 populations to changes in environmental or human health. However, improvements in precision,
121 sensitivity, speed, cost, and accuracy of NGS and downstream analyses are necessary for
122 effective utilization in biodefense research [2-6].

123

124 On January 10, 2018, the Mid-Atlantic Microbiome Meet-up (M³) organization held a conference
125 aimed at understanding how the biodefense and pathogen detection fields are transformed by
126 new biological and computational technologies. **While biodefense was broadly discussed, the**
127 **participants focused primarily on emerging infectious disease applications.** The meeting took
128 place in the STAMP Student Union on the University of Maryland campus in College Park. The
129 M³ consortium brings together microbiome researchers from different sectors to discuss
130 challenges, develop standards and best practices, and help connect data generators with data
131 analysts [7]. The M³ community is constantly growing and, as of this publication, has 140
132 members from over 25 different institutions. The conference was attended by 67 participants
133 from academia, government and industry (**Fig 1**), with expertise in areas such as biodefense,
134 computer science, genomics, microbiology, and public health. There were two talks given by
135 invited speakers, 15 oral presentations selected from submitted abstracts, and several posters
136 displayed at the meeting (**Supplementary Table 1**) [8]. Additionally, there were three interactive
137 break-out sessions to address challenges of the field and encourage networking
138 (**Supplementary Table 2**). **The event was sponsored in part by CosmosID, Inc. but they did**
139 **participate in the organization of the event, nor in the selection of speakers and topics being**
140 **discussed.**

141 The tone for the meeting was set by the keynote address presented by Dr. Tara O'Toole,
142 Executive Vice President of the non-profit strategic investor In-Q-Tel, Inc.. Pointing to the
143 problems in detection, containment, and treatment during the recent H1N9 pandemic and Ebola
144 epidemic, Dr. O'Toole shared that current progress in the field is disappointing because
145 biodefense is not a priority for any single government agency, funding support is irregular, and
146 epidemics are becoming more common. Increasing international competition for biotechnology
147 advancements and leadership make it even more important to stimulate progress.

148 Dr. O'Toole outlined several keys to innovation and policy, which were echoed by the
149 presentations and discussions throughout the remainder of the meeting, including (1) the
150 willingness to think anew, (2) development of new tools and instruments, (3) implementation of
151 a technology-focused biodefense strategy, (4) delivery of near real-time situational awareness
152 for existing epidemics by leveraging modern data analytics and networked communications, and
153 (5) establishment of rich human networks and cross-sector partnerships between government
154 agencies, the private sector, and academia.

155 **Key Conclusions**

156

157 We start by highlighting the key conclusions and recommendations identified by the participants
158 in the meeting:

- 159 1. Sequencing-based assays frequently face challenges related to limits of detection and
160 technical biases, and culturing or other enrichment strategies remain necessary in many
161 applications. The accurate quantification of viable organisms or metabolic activity within
162 complex metagenomic samples remains an open challenge that is unlikely to be solved
163 through sequencing alone.
- 164 2. Current sample processing approaches tend to exclude viral and fungal/eukaryotic
165 components of microbial communities. In the case of viruses, this problem is
166 compounded by poor taxonomies and database resources.
- 167 3. Analytical approaches, community standards, and software for temporal data analysis
168 have lagged behind rapidly increased generation of such data.
- 169 4. Robust bioinformatics tools are critical for future progress. These tools must be
170 developed to better match the needs of end-users and must be subject to critical
171 validation.

172 5. Data standards are essential for ensuring the quality and usefulness of shared datasets,
173 but overly onerous reporting requirements discourage sharing. In cases where privacy is
174 a concern, we must also develop solutions that allow for secure storage and processing
175 of sensitive data.

176 These key recommendations are summarized in **Table 1** and more extensively discussed
177 below.

178

179 ***1. Sequencing-based assays frequently lack sensitivity***

180

181 While the biodefense community has benefited from high-throughput sequencing strategies,
182 these methods are not always as sensitive as required. In some cases, culturing is still the most
183 reliable method for detecting pathogens because standard sequencing pipelines are not always
184 available and achieving required sequencing depths may be cost-prohibitive. Dr. Sarah Allard
185 (UMD SPH) shared her work from CONSERVE (a Center of Excellence at the Nexus of
186 Sustainable Water Reuse, Food, and Health), whose mission is to enable the safe use of non-
187 traditional irrigation water sources on food crops [9]. Dr. Allard used both culture-based and
188 sequence-based methods to detect foodborne pathogens in water samples. She concluded that
189 culture-based techniques are currently the most sensitive pathogen detection strategies and
190 that sequencing analysis sensitivity and stringency vary strongly by method.

191

192 From a public health perspective, quantification of viable organisms contributing to disease is
193 essential, but cannot be achieved with metagenomic analysis alone. Culturing and other
194 approaches are important for gaining insight into the metabolic activity of the microbes in a
195 community [10]. Additionally, researchers must often make a trade-off between the sensitivity of
196 their detection methods and the computational costs of analyzing increasingly deep sequencing

197 datasets. Even partial culturing of select organisms or samples can help shift this trade-off. As
198 commented during a breakout session, “you can’t always sequence your way out of it.”

199

200 **2. *Few studies look beyond bacterial pathogens***

201

202 Shotgun metagenomics and a decrease in the cost of DNA sequencing have enabled
203 researchers to analyze the genetic potential of microorganisms directly from an environmental
204 sample. However, the majority of microbiome and metagenome studies focus only on the
205 prokaryotic component of the community, while few have explored the roles of fungi or viruses
206 in these microbial communities. This is due, in large part, to limitations in resources, laboratory
207 procedures and, in the case of viruses, the lack of a universally distributed marker gene.

208 Additional barriers to mycobiome and virome studies, include the ability to obtain sufficient
209 material from low biomass environments, high levels of host contamination, incomplete
210 databases, and a lack of available wet lab protocols and computational analysis pipelines. At
211 the meeting, it was noted that central repositories for shared protocols do exist (e.g. protocols.io
212 [11]) and a concerted effort in viral protocol sharing has been made by the Gordon and Betty
213 Moore Foundation, which funds VERVE Net [12]. Proposed goals to address other barriers
214 included providing financial and/or publication incentives for database curation and maintenance
215 and focusing work on gene function identification. Since the NCBI SRA already contains many
216 metagenomic sequencing datasets, it may be worthwhile to identify novel fungal and viral
217 genomes from existing datasets to optimize data usage, as this approach has been employed in
218 previous studies of environmental viruses [13].

219

220 Despite the aforementioned barriers to fungal and viral metagenomics, additional research in
221 this area can significantly contribute to biodefense. One such important topic is the spread of

222 viral pathogens. Invited seminar speaker Dr. Don Milton (UMD SPH) presented his work on the
223 transmission of the influenza virus in college dormitories [14]. The Centers for Disease Control
224 and Prevention (CDC) suggests that human influenza transmission mainly occurs by droplets
225 made when people with flu cough, sneeze, or talk. However, Dr. Milton explained that dueling
226 reviews have disputed the importance of airborne transmission [15-20]. He presented NGS data
227 showing that exhaled breath of symptomatic influenza cases contains infectious virus in fine
228 particles, suggesting that aerosol exposures are likely an important mode of transmission.

229

230 ***3. Tracking microbial communities across time and topography***

231

232 Temporal and biogeographic sequencing studies provide increased resolution of microbial
233 community shifts. In the context of biodefense, this is important for detecting and containing
234 outbreaks. Additionally, these studies provide insight into environmental changes, which may
235 contribute to epidemics by causing shifts in disease vectors and/or spurring human migration to
236 new regions or densely-populated urban areas. Several presentations at the meeting shared
237 spatiotemporal microbiome analyses of different environments. Dr. Sean Conlan (NIH, NHGRI)
238 presented his work using metagenomics to study outbreaks of nosocomial infections and
239 identified the transfer of plasmids from patients to the hospital environment [21,22]. Gherman
240 Uritskiy (JHU) and Dr. Sarah Preheim (JHU) used a combination of marker gene and
241 metagenomics approaches to characterize changes in environmental microbiomes in response
242 to perturbations. Uritskiy studied halite endoliths from the Atacama Desert in Chile over several
243 years and showed how they were significantly impacted by rainstorms. Dr. Preheim compared a
244 biogeochemical model to microbial communities' changes in a lake over the spring and summer
245 to reveal the influence of energy availability on microbial population dynamics.

246

247 While time-series datasets provide valuable information, they are much more difficult to analyze
248 with current statistical methods and models than cross-sectional sampling strategies [23,24].
249 Among other reasons, this is because it is difficult to identify the optimal sampling frequency, the
250 compositional nature of microbiome data frequently violates assumptions of statistical methods,
251 and commonly available software tools are often insufficient for required complex comparisons.
252 Addressing this, Dr. J Gregory Caporaso (NAU) presented QIIME 2 (<https://qiime2.org>) and
253 shared his team's QIIME 2 plugin, q2-longitudinal, which incorporates multiple methods for
254 characterizing longitudinal and paired-sample marker gene datasets [25].

255

256 **4. Development and application of metagenomic analysis tools is critical for progress**

257

258 Computational methods required for metagenomic analyses include taxonomic abundance
259 profiling, taxonomic sequence classification and annotation, functional characterization, and
260 metagenomic assembly. Many of the presentations at the meeting shared new and/or improved
261 tools for different aspects of microbiome studies. Victoria Cepeda (UMD) described how her
262 tool, MetaCompass, uses reference genomes to guide metagenome assembly [26] and
263 Gherman Uritskiy (JHU) presented his pipeline, metaWRAP, for the pre-processing and binning
264 of metagenomes [27]. Furthermore, Brian Ondov (UMD, NIH, NHGRI) shared his
265 implementation of the MinHash containment estimation algorithm to screen metagenomes for
266 the presence of genomes and plasmids [28]. Data visualization is important for accurately
267 interpreting microbiome data analyses and Dr. Héctor Corrada-Bravo (UMD) demonstrated how
268 to use his lab's tool, Metaviz [29], for interactive statistical analysis of metagenomes.

269

270 Conventional metagenomic analyses often reflect the most abundant elements from a complex
271 sample and cannot detect rare elements with confidence. Dr. Nicholas Bergman (NBACC)

272 shared a more sensitive single cell metagenomics approach that allows for increased detection
273 of all elements of a community sample. Dr. Bergman’s talk also emphasized the necessity of
274 improving sensitivity, preventing contamination, eliminating biases, and increasing efficiency for
275 sequencing-based techniques.

276

277 *Bioinformatics tools should better match the needs of end-users*

278

279 Many discussions at the meetings focused on how the field can optimize tool utility. It was
280 agreed that scientists should always carefully evaluate the strengths and weakness of available
281 methods, either via existing “bake-off” studies or through the available documentation, to ensure
282 they are using the best tools to address their specific problem. Tool developers should disclose
283 the limits of their methods and advise on the types of data their software is best suited to
284 analyze. Developers should also work towards producing software that is easy to download and
285 install, providing comprehensive documentation for their tools, and ensuring open-access for the
286 academic community. As a community, we should encourage that publications list not only
287 cases and data types where methods perform best, but also where they under-perform or even
288 fail. Additional studies, like the Critical Assessment of Metagenome Interpretation (CAMI)
289 [30,31], Microbiome Quality Control project [32], or challenges run under the aegis of
290 PrecisionFDA [33] should be conducted to help characterize strengths and weaknesses of
291 different approaches and evaluate their impact on data analysis and interpretation.

292

293 Some meeting attendees are currently contributing to these goals. Dr. Nathan Olson (UMD,
294 NIST) presented his evaluation of different 16S rRNA marker gene survey bioinformatic
295 pipelines using mixture samples. Additionally, Dr. Daniel Nasko (UMD) characterized how
296 genomic database growth affects study findings, showing that different versions of the RefSeq

297 database strongly influenced species-level taxonomic classifications from metagenomic
298 samples [34]. Because the version of software and databases used can significantly affect
299 findings, this information should be reported more consistently in the literature. Furthermore, we
300 should consider strategies to preserve previous software and database versions to enable
301 future replication of analyses.

302

303 *Bioinformatics tools must better navigate the trade-off between speed and accuracy*

304

305 Metagenomic analysis methods vary in central processing unit (CPU) time, memory, and disk
306 resource usage and this is not always clearly reported in software publications. Additionally,
307 method scalability relative to size or type of input data also varies considerably. Optimizing
308 speed and accuracy is especially important for biodefense applications. For instance,
309 improvements in NGS analysis allowing for collection and analysis of samples in a clinically
310 relevant time frame can help effectively track hospital outbreaks and prevent the spread of
311 infection [35]. Furthermore, confidence in the accuracy of these analyses is required to execute
312 appropriate plans of action and prevent panic. Recently, findings of *Bacillus* strains on the
313 International Space Station that were genomically similar to pathogenic *B. anthracis* required
314 more detailed characterization to ensure that their presence was not a concern for the health of
315 the crew [36-38]. *B. anthracis* was also initially reported to be found in the NYC subway system,
316 along with *Yersinia pestis*, the pathogen responsible for the plague [39]. After public attention
317 prompted further analysis, the authors found no evidence that these organisms were present
318 and found no evidence of pathogenicity [40,41], again highlighting the importance of careful
319 evaluation and interpretation of results, especially those with severe public health
320 consequences.

321

322 Many different strategies for speeding up analyses were discussed at the meeting, including
323 hardware, software, and algorithm choice. Some hardware considerations for speed of analyses
324 include balancing CPUs with co-processors such as graphics processing units (GPUs) or field
325 programmable gate arrays (FPGAs), server configuration in terms of amount of random access
326 memory (RAM), or disk storage type and speed. Programs and algorithms vary in accuracy as
327 well as ease of parallelization. Often a slower yet parallelizable algorithm is preferred to one that
328 is not parallelizable. If a program supports parallelism, consideration should be given to the type
329 of hardware required. For example, some available options include large multicore servers for
330 multithreaded applications, cluster nodes for distribution of compute jobs, or cloud computing
331 solutions. Other strategies might involve analyzing only a subset of the data or using a smaller,
332 application-specific reference database.

333
334 Finally, strategies discussed for speeding up time-critical analyses included employing a multi-
335 tiered approach (e.g., a quick first pass followed by more detailed analyses [42]), and
336 considering the suitability of various sequencing platforms for certain applications. Interventions
337 or optimizations were discussed with regard to their impact on analysis accuracy and
338 interpretation of results. Preferred solutions are the ones that provide both the desired speed
339 and accuracy, though more often than not there is a trade-off between the two. The optimal
340 balance also depends on the use case. Assessment and validation methods are required to
341 characterize a method's speed and accuracy. It will be up to the subject matter experts to
342 determine the desired accuracy level for each case and the extent to which they can sacrifice
343 accuracy for speed.

344

345 **5. Data needs to be moved out of private silos and into public repositories**

346

347 Data sharing is continually a challenge that gets raised within the biological community,
348 especially as DNA/RNA sequencing becomes more ubiquitous and tangible outside of core
349 facilities [43]. This challenge is prevalent across multiple scientific disciplines, and was recently
350 highlighted by the National Research Council as a priority for microbial forensics [44]. There are
351 numerous reasons data are not being shared, including the need to protect personally
352 identifiable information or intellectual property rights prior to publication, and the lack of
353 sufficient infrastructure or manpower to upload at scale. However, leveraging this diversity and
354 breadth of data will be important for an effective biodefense capacity, as well as other
355 bioscience applications like healthcare, pharmaceuticals, agriculture, and industry. In order to
356 incentivize data sharing, we need to evaluate and improve publicly available resources for
357 storing and processing data.

358

359 Inherent altruism or obligation to share data should be met with as little friction as possible, and
360 we need to incentivize openness. One incentive is academic credit through authorship on
361 publications, though this will require combined efforts of researchers, journal editors, and
362 funding agencies to better define what contributions constitute data authorship and what
363 responsibilities data authors have [45,46]. Another potential incentive is the availability of free
364 software for data analysis, and meeting participants debated the desirability and sustainability of
365 service-based options (e.g. MG-RAST [47]) compared to user-installable software options (e.g.
366 QIIME [48], mothur [49]). At the meeting, Dr. Nur A. Hasan (CosmosID, Inc.) highlighted the
367 cloud-based metagenome tools and databases his company has to offer. There are also strong
368 movements toward software sharing, such as the Astrophysics Source Code Library [50] and
369 the Materials Resource Registry at NIST [51].

370

371 It is expected that some quality standard is needed to maintain useable, open repositories.
372 Where that standard is set can affect how much data is shared. For example, a high bar may

373 ensure high quality sequences and comprehensive metadata but minimize sharing, while a
374 lower quality bar will more likely move data out of silos. The solution may be a combination of
375 repositories with varying standards, or a single repository which allows for varying degrees of
376 annotation completeness and allows the user to modify searches based on that feature. It is
377 important to note that a single repository may be difficult to reliably curate and manage at scale.
378 Another option is distributed but federated systems, like used by the US Virtual Astronomical
379 Observatory [52]. Groups like the Genomic Standards Consortium [53,54] are working towards
380 improving data quality by supporting projects such as Minimum Information about any
381 Sequence (MIxS) [55], which establishes standards for describing genomic data and provides
382 checklists to help with annotation. We need to build a community consensus on how much
383 metadata is required to make reporting less onerous for data providers but ensure data usability
384 by others in the field.

385
386 Incentivizing open data sharing should not be the only solution, as some sensitive data cannot
387 be openly shared due to privacy regulations (e.g. human genomes and Health Insurance
388 Portability and Accountability Act regulations). Other sectors, such as the financial industry,
389 have long been working on solutions to enable storage, transit, and operations of protected
390 data. These solutions include software-based approaches (e.g., homomorphic encryption, Yao's
391 protocol, secure fault-tolerant protocols, oblivious transfer) and hardware-based approaches
392 (e.g. AES full disk encryption for data storage, Intel® Software Guard Extension for secure
393 operations). Dr. Stephanie Rogers presented the GEMStone 2.0 project from B. Next, an IQT
394 Lab, called SIG-DB, which explores homomorphic encryption and Intel Software Guard
395 Extension (SGX) to securely search genomic databases [56]. Early results of applying these
396 solutions to biological data are promising and should be explored more fully.

397

398 **Conclusions**

399

400 Overall, this meeting successfully brought together scientists from academia, government and
401 industry to present their research and discuss how high-throughput genomics methods have
402 stimulated interest and progress in biodefense and pathogen detection. Notably, meeting
403 participants used NGS tools to identify the transfer of microbes from patients to their hospital
404 environments, track the transmission of influenza in a community living space, study
405 environmental shifts over time, and evaluate the safety of using non-traditional water sources on
406 food crops. These studies, and others, have been partly driven by cheaper, more reliable
407 sequencing technologies and improvements in computational analysis tools. Open-source
408 software for sequence processing and quality control, taxonomic annotation, metagenomic
409 assembly and binning, and data visualization have been essential for growth. Continued
410 development of these resources will result in significant scientific advances.

411

412 Despite this progress, there are several limitations to using NGS approaches for biodefense
413 problems. First and foremost, sequencing methods are unable to accurately quantify viable
414 organisms from metagenomic samples, which is essential for identifying potential threats to
415 public health. Beyond that, applications for which NGS approaches are well-suited still present
416 many challenges. Although sequencing costs are steadily declining, it remains expensive to
417 process, computationally analyze, and store the increasingly large datasets that are generated.
418 Confident detection of infectious, but potentially rare pathogens in a community often requires
419 very deep sequencing and scientists must make the appropriate speed, cost, and accuracy
420 trade-offs to best answer their research questions. In many cases, sequencing experiments may
421 need to be complemented with culturing, enrichment, or other targeted approaches. **Because of**
422 **these limitations, and others, researchers must be extremely careful when interpreting data to**
423 **identify biothreats; reporting false positives without critical validation can have significant fiscal**
424 **and public health consequences.** Developing the capacity to identify not only when a potential

425 pathogen is present, but also at what levels it is actively contributing to infectious disease will
426 greatly improve our response to biothreats. Another area that requires further investigation is
427 the detection of antimicrobial resistance. While only briefly highlighted in the meeting talks about
428 influenza and noscomial tracing, antimicrobial resistance poses a significant threat to public
429 health and biodefense. Current metagenomic sequencing methods allow us to identify
430 antimicrobial resistance genes from different environments, however these techniques cannot
431 determine whether these genes are actively being expressed and are currently not practical for
432 wide-spread adoption in clinical settings [57].

433

434 To date, few microbiome studies have focused on viral and fungal/eukaryotic organisms,
435 despite their potentially important community interactions and roles in pathogenesis. In order to
436 generate relevant virome and mycobiome datasets, we must improve sample processing
437 techniques and dedicate resources to effectively curate and maintain publicly available
438 databases. We also need to develop advanced statistical toolkits for analyzing longitudinal
439 studies. In general, tool developers should focus on creating user-friendly, adaptable resources,
440 with comprehensive documentation and clear descriptions of default settings and optional
441 parameters. These tools must be critically evaluated for their appropriate use-cases; however,
442 when looking for emerging threats, it will be necessary to develop validation approaches that do
443 not require the use of gold standards.

444

445 In order to encourage additional growth, the greater scientific community should invest in
446 expanding and enforcing clear standards for genomic datasets. If set appropriately, these
447 standards will help incentivize data sharing and improve the quality and usability of public
448 repositories. Additional focus should be on strengthening best practices and solutions for
449 handling sensitive datasets that are subject to privacy regulations. Moving forward, active

450 conversations between researchers and policy makers will be essential to expand and
451 implement these ideas in biodefense.

452

453 **Declarations**

454

455 **Ethics approval and consent to participate**

456

457 Not applicable

458

459 **Consent for publication**

460

461 Not applicable

462

463 **Availability of data and material**

464

465 Not applicable

466

467 **Competing interests**

468

469 The conference was partly supported through funding from CosmosID, Inc.. The sponsor did not
470 participate in the organization of the event or the selection of speakers and topics.

471

472 **Funding**

473

474 The meeting was supported in part by the Center for Health-related Informatics and Bioimaging,
475 a Center organized under the MPowering the State Partnership between the University of

476 Maryland Baltimore and College Park campuses. JSM, BB, VCE, MF, JG, KJ, NS, and MP were
477 supported in part by grants to MP, including grant R01-AI-100947 from the NIH and grant IIS-
478 1513615 from the NSF. DJN and TJT were supported in part by the FunGCAT program from the
479 Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects
480 Activity (IARPA), via the Army Research Office (ARO) under Federal Award No. W911NF-17-2-
481 0089. HC was supported by the NIH, R01 grant GM114267. The views and conclusions
482 contained herein are those of the authors and should not be interpreted as necessarily
483 representing the official policies or endorsements, either expressed or implied, of the ODNI,
484 IARPA, ARO, or the US Government. The contributions of ALB, NHB, MJR, DDS and SR were
485 funded under Contract No. HSHQDC-15-C-00064 awarded by the Department of Homeland
486 Security (DHS) Science and Technology Directorate (S&T) for the operation and management
487 of the National Biodefense Analysis and Countermeasures Center (NBACC), a Federally
488 Funded Research and Development Center. The views and conclusions contained in this
489 document are those of the authors and should not be interpreted as necessarily representing
490 the official policies, either expressed or implied, of the DHS or S&T. In no event shall DHS,
491 NBACC, S&T or Battelle National Biodefense Institute have any responsibility or liability for any
492 use, misuse, inability to use, or reliance upon the information contained herein. DHS does not
493 endorse any products or commercial services mentioned in this publication. JGC was supported
494 in part by National Cancer Institute of the National Institutes of Health under the awards for the
495 Partnership of Native American Cancer Prevention U54CA143924 (UACC) and U54CA143925
496 (NAU), and by the National Science Foundation award 1565100. SC was supported by NIH
497 Intramural Research. JD and GU were supported in part by the NSF, grant DEB1556574 to JD.
498 BDO was supported by the Intramural Research Program of the National Human Genome
499 Research Institute, National Institutes of Health and utilized the computational resources of the
500 NIH HPC Biowulf cluster (<https://hpc.nih.gov>).

501

502 **Authors' contributions**

503

504 TT and MP organized the meeting. JK, JSM and DN reviewed abstracts, and SMR, NHB, and
505 JD selected abstracts for the sessions they chaired. All listed authors contributed to the writing
506 of the report or met at least one of the following requirements: they (1) gave an oral
507 presentation, (2) presented a poster, and/or (3) attended and contributed to the interactive
508 break-out sessions.

509

510 **Acknowledgements**

511

512 We would like to thank all those who helped make this meeting a success, especially Barbara
513 Lewis (UMD) who organized the administrative aspects of the program and CosmosID, Inc. for
514 funding the event. We would also like to thank Dr. Jayne Morrow and Dr. Robert Hanisch for
515 providing helpful feedback on the manuscript.

516

517 Opinions expressed in this paper are the authors' and do not necessarily reflect the policies and
518 views of NIST or affiliated venues. Certain commercial equipment, instruments, or materials are
519 identified in this paper in order to specify the experimental procedure adequately. Such
520 identification is not intended to imply recommendations or endorsement by NIST, nor is it
521 intended to imply that the materials or equipment identified are necessarily the best available for
522 the purpose. Official contribution of NIST; not subject to copyrights in USA.

523

524 **Figure Legends**

525

526 **Figure 1.** Different sectors and institutions represented at the January 2018 M³ Meetup

Research Gaps	Current Limitations	Community Goals
<p>Tracking microbial communities across time and topography (Key Conclusions 1 & 3)</p> <p><u>Importance:</u> Studies incorporating temporal and/or spatial sampling allow us to detect important shifts in community dynamics</p> <p><u>Application example:</u> Detecting the spread of infection in a hospital or of a pathogen contaminating crops and spreading food-borne illness</p>	<ul style="list-style-type: none"> Sequencing strategies are not able to quantify viable organisms (which is essential for biodefense applications) Lack of well-established statistical approaches for exploring longitudinal microbiome data Increased sample size makes these studies more expensive and harder to obtain sufficient statistical power for all subjects/timepoints/regions 	<ul style="list-style-type: none"> Collection, sequencing, and sharing of more time-series datasets Development of statistical methods and tools to help analyze longitudinal and/or geospatial microbiome datasets
<p>Looking beyond bacterial pathogens (Key Conclusion 2)</p> <p><u>Importance:</u> Viral and fungal components of the microbiome are often under-explored, despite their potential implications in biodefense</p> <p><u>Application example:</u> Better understanding the transmission of infectious viruses, like influenza</p>	<ul style="list-style-type: none"> Lack of a universally distributed marker gene (viruses) Difficult to obtain sufficient material from low biomass environments High levels of host contamination Incomplete databases 	<ul style="list-style-type: none"> More consistent database curation and maintenance (potentially incentivized financially or with publications) Improved gene function identification
<p>Development and application of metagenomic analysis tools (Key Conclusion 4)</p> <p><u>Importance:</u> Computational tools need to be developed to help improve the utility of high-throughput sequencing strategies for biodefense problems</p> <p><u>Application example:</u> Improved metagenome assembly methods could better delineate between different strains of a pathogen in samples</p>	<ul style="list-style-type: none"> Tools for metagenome pre-processing, assembly, and binning are not always sensitive or fast enough for detection of pathogens in a sample As sequencing technologies advance, we need new tools to handle output from long and short read technologies, as well as single cell metagenomics approaches 	<ul style="list-style-type: none"> Easy to install, open-access software with comprehensive documentation detailing best and worst use cases Defined metrics for critical assessment and validation of existing tools Software and database versions should be more consistently reported in the literature and preserved for future replication of analyses

<p>Navigating the trade-off between speed and accuracy (Key Conclusion 4)</p> <p><u>Importance:</u> Metagenomic analysis used for pathogen detection and identification are time-sensitive</p> <p><u>Application example:</u> Deciding if a food product should be recalled due to contamination</p>	<ul style="list-style-type: none"> • Current algorithms vary in speed and accuracy (often sacrificing one for the other) • Large datasets, error-prone heuristics, and coarse resolution of k-mer based methods present challenges 	<ul style="list-style-type: none"> • Better documentation of available tools to help users optimize their software choice based on their available resources • Improvements in sequencing technologies and tools/algorithms to improve both speed and accuracy
<p>Storing and sharing data (Key Conclusion 5)</p> <p><u>Importance:</u> Access to publicly available datasets will help in verification of results and advance of scientific knowledge. Scientists need to be encouraged to move their data out of private silos and into shared databases</p>	<ul style="list-style-type: none"> • Not all data can be shared because it is important to protect personally identifiable information or intellectual property rights • Lack of sufficient infrastructure or manpower to upload or store datasets at scale 	<ul style="list-style-type: none"> • Defined quality standard to maintain usable, open repositories • Improved ways for secure interrogation of genomic datasets that cannot be openly shared due to privacy regulations

529

530 **Table 1.** Outline of current research gaps and future goals discussed at the January 2018 M³

531 Meeting

532

533

	Speaker	Title
Keynote	Tara O'Toole In-Q-Tel, Inc.	Bioterror, and Biodefense, 2018
Invited Seminar	Don Milton UMD SPH	College Dorms as a Laboratory for Studying Respiratory Infection
Data-driven session Chair: Stephanie Rogers B.Next (In-Q-Tel, Inc.)	Stephanie Rogers B.Next (In-Q-Tel, Inc.)	Bridging technology, venture, and national security
	Daniel Nasko UMD CBCB	Tragedy of the commons: RefSeq database growth influences the accuracy and sensitivity of species identification from metagenomic samples
	Sarah Allard UMD SPH	Comparison of sequencing and culture-based methods for the detection of foodborne pathogens in non-traditional irrigation water in the Mid-Atlantic United States: A CONSERVE study
	Sean Conlan NHGRI, NIH	Tracking Antibiotic Resistance Across Space and Time
	Greg Caporaso NAU, NCI, NIH	Longitudinal analysis of microbiomes
Methods-driven session Chair: Nicholas Bergman NBACC	Nicholas Bergman NBACC	Metagenomics in bioforensics
	Brian Ondov UMD, NHGRI, NIH	Mash Screen: Fast sequence containment estimation using MinHash
	Nathan D. Olson UMD CBCB and NIST	A Sample Mixture Experiment to Assess 16S rRNA Metagenomic Methods
	Victoria Cepeda UMD CBCB	MetaCompass: Reference-guided Assembly of Metagenomes
	Héctor Corrada-Bravo UMD CBCB	Metaviz: Interactive Statistical and Visual Analysis of Human Microbiome Project Data

Ecology-driven session Chair: Jocelyne DiRuggerio JHU	Gherman Uritskiy JHU	Dynamic Response of Atacama Desert Extremophiles to Weather Perturbations
	Sarah Preheim JHU	Frequency and impact of transitions between microbial populations mediating biogeochemical cycling in a freshwater lake
	Nur A. Hasan CosmosID	Cloud based bioinformatics platform for cross-disciplinary microbiome research

535

536 **Supplementary Table 1.** Outline of oral presentations at the January 2018 M³ Meeting

537

Break-out Session	Chairs	Topics
Viral and fungal pathogens: off the beaten path	Jacquelyn Meisel UMD CBCB and Daniel Nasko UMD CBCB	<ul style="list-style-type: none"> • Barriers to studying viruses and fungi and actionable goals to address • Detection of viral and fungal pathogens
Crowdsourcing biodefense: data sharing, data standards, data security, and data visualization	Todd Treangen UMD CBCB and Brian Ondov UMD, NHGRI, NIH	<ul style="list-style-type: none"> • Homomorphic encryption techniques for secure computations • Data visualization as incentive for data sharing • Data standards: raising the bar to ensure quality vs lowering the bar to bring data out of silos • Cloud based data sharing: practical or impractical?
Need for speed: navigating the trade-off between analysis accuracy and speed	Adam Bazinet NBACC and Nathan D. Olson UMD CBCB and NIST	<ul style="list-style-type: none"> • Evaluation of time-sensitive metagenomic analyses • Strategies for speeding up such analyses without compromising accuracy

538

539 **Supplementary Table 2.** Outline of interactive break-out sessions at the January 2018 M³

540 Meeting

541

542 **References**

- 543 1. Drew TW, Mueller-Doblies UU. Dual use issues in research - A subject of increasing
544 concern? *Vaccine*. 2017;35:5990–4.
- 545 2. Gardy JL, Loman NJ. Towards a genomics-informed, real-time, global pathogen surveillance
546 system. *Nat. Rev. Genet.* Nature Publishing Group; 2018;19:9–20.
- 547 3. Robinson ER, Walker TM, Pallen MJ. Genomics and outbreak investigation: from sequence
548 to consequence. *Genome Med.* BioMed Central; 2013;5:36.
- 549 4. Miller RR, Montoya V, Gardy JL, Patrick DM, Tang P. Metagenomics for pathogen detection
550 in public health. *Genome Med.* BioMed Central; 2013;5:81.
- 551 5. Lipkin WI. The changing face of pathogen discovery and surveillance. Nature Publishing
552 Group. Nature Publishing Group; 2013;11:133–41.
- 553 6. Forbes JD, Knox NC, Ronholm J, Pagotto F, Reimer A. Metagenomics: The Next Culture-
554 Independent Game Changer. *Front Microbiol.* Frontiers; 2017;8:1069.
- 555 7. main groups.io Group. Available from: <https://m3.groups.io/g/main>
- 556 8. Winter 2018 Mid-Atlantic Microbiome Meetup Biodefense and Pathogen Detection Agenda
557 [Internet]. [cited 2018 May 4]. Available from: [https://cpb-us-
558 e1.wpmucdn.com/blog.umd.edu/dist/d/418/files/2017/10/WinterM3_agenda_final-27afpqx.pdf](https://cpb-us-e1.wpmucdn.com/blog.umd.edu/dist/d/418/files/2017/10/WinterM3_agenda_final-27afpqx.pdf)
- 559 9. CONSERVE: A Center of Excellence at the Nexus of Sustainable Water Reuse, Food, and
560 Health, Year 1 Achievements (March 2016-February 2017) [Internet]. Available from:
561 [https://static1.squarespace.com/static/578101761b631b1a87aa0a3c/t/59f8f8e8e31d19ae52831
562 0e9/1509488877173/CONSERVE_annual_report.pdf](https://static1.squarespace.com/static/578101761b631b1a87aa0a3c/t/59f8f8e8e31d19ae528310e9/1509488877173/CONSERVE_annual_report.pdf)

- 563 10. Singer E, Wagner M, Woyke T. Capturing the genetic makeup of the active microbiome in
564 situ. *The ISME journal*. Nature Publishing Group; 2017;11:1949–63.
- 565 11. Teytelman L, Stoliartchouk A, Kindler L, Hurwitz BL. *Protocols.io: Virtual Communities for*
566 *Protocol Development and Discussion*. *Plos Biol. Public Library of Science*; 2016;14:e1002538.
- 567 12. VERVE Net [Internet]. *protocols.io*. Available from: protocols.io/g/verve-net
- 568 13. Paez-Espino D, Eloe-Fadrosch EA, Pavlopoulos GA, Thomas AD, Huntemann M, Mikhailova
569 N, et al. Uncovering Earth's virome. *Nature*. *Nature Research*; 2016;536:425–30.
- 570 14. Yan J, Grantham M, Pantelic J, Bueno de Mesquita PJ, Albert B, Liu F, et al. Infectious virus
571 in exhaled breath of symptomatic seasonal influenza cases from a college community.
572 *Proceedings of the National Academy of Sciences of the United States of America*.
573 2018;115:1081–6.
- 574 15. Killingley B, Nguyen-Van-Tam J. Routes of influenza transmission. *Influenza Other Respir*
575 *Viruses*. Wiley/Blackwell (10.1111); 2013;7 Suppl 2:42–51.
- 576 16. Tellier R. Aerosol transmission of influenza A virus: a review of new studies. *J R Soc*
577 *Interface*. The Royal Society; 2009;6 Suppl 6:S783–90.
- 578 17. Bridges CB, Kuehnert MJ, Hall CB. Transmission of influenza: implications for control in
579 health care settings. *Clin. Infect. Dis*. 2003;37:1094–101.
- 580 18. Tellier R. Review of Aerosol Transmission of Influenza A Virus. *Emerging Infect. Dis*.
581 *Centers for Disease Control and Prevention*; 2006;12:1657–62.
- 582 19. Lemieux C, Brankston G, Gitterman L, Hirji Z, Gardam M. Questioning aerosol transmission
583 of influenza. *Emerging Infect. Dis*. 2007;13:173–4–authorreply174–5.

- 584 20. Brankston G, Gitterman L, Hirji Z, Lemieux C, Gardam M. Transmission of influenza A in
585 human beings. *The Lancet Infectious Diseases*. Elsevier; 2007;7:257–65.
- 586 21. Conlan S, Park M, Deming C, Thomas PJ, Young AC, Coleman H, et al. Plasmid Dynamics
587 in KPC-Positive *Klebsiella pneumoniae* during Long-Term Patient Colonization. *mBio*.
588 2016;7:e00742–16.
- 589 22. Weingarten RA, Johnson RC, Conlan S, Ramsburg AM, Dekker JP, Lau AF, et al. Genomic
590 Analysis of Hospital Plumbing Reveals Diverse Reservoir of Bacterial Plasmids Conferring
591 Carbapenem Resistance. Bonomo RA, editor. *mBio*. American Society for Microbiology;
592 2018;9:e02011–7.
- 593 23. Faust K, Lahti L, Gonze D, de Vos WM, Raes J. Metagenomics meets time series analysis:
594 unraveling microbial community dynamics. *Curr. Opin. Microbiol.* Elsevier Current Trends;
595 2015;25:56–66.
- 596 24. Gerber GK. The dynamic microbiome. *FEBS Lett.* Wiley-Blackwell; 2014;588:4131–9.
- 597 25. Bokulich N, Zhang Y, Dillon M, Rideout JR, Bolyen E, Li H, et al. q2-longitudinal: a QIIME 2
598 plugin for longitudinal and paired-sample analyses of microbiome data. *bioRxiv*. Cold Spring
599 Harbor Laboratory; 2017;:223974.
- 600 26. Cepeda V, Liu B, Almeida M, Hill CM, Koren S, Treangen TJ, et al. MetaCompass:
601 Reference-guided Assembly of Metagenomes. 2017.
- 602 27. Uritskiy GV, DiRuggiero J, Taylor J. MetaWRAP - a flexible pipeline for genome-resolved
603 metagenomic data analysis. *bioRxiv*. 2018.

604 28. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast
605 genome and metagenome distance estimation using MinHash. *Genome biology. BioMed*
606 *Central*; 2016;17:132.

607 29. Wagner J, Chelaru F, Kancherla J, Paulson JN, Zhang A, Felix V, et al. Metaviz: interactive
608 statistical and visual analysis of metagenomic data. *Nucleic acids research*. 2018;514:59.

609 30. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, et al. Critical
610 Assessment of Metagenome Interpretation—a benchmark of metagenomics software. *Nature*
611 *methods*. Nature Publishing Group; 2017;14:1063–71.

612 31. Bremges A, McHardy AC. Critical Assessment of Metagenome Interpretation Enters the
613 Second Round. *mSystems*. American Society for Microbiology Journals; 2018;3:537.

614 32. Sinha R, Abu-Ali G, Vogtmann E, Fodor AA, Ren B, Amir A, et al. Assessment of variation in
615 microbial community amplicon sequencing by the Microbiome Quality Control (MBQC) project
616 consortium. *Nat. Biotechnol*. Nature Publishing Group; 2017;35:1077.

617 33. Altman RB, Prabhu S, Sidow A, Zook JM, Goldfeder R, Litwack D, et al. A research
618 roadmap for next-generation sequencing informatics. *Sci Transl Med*. American Association for
619 the Advancement of Science; 2016;8:335ps10–0.

620 34. Nasko DJ, Koren S, Phillippy AM, Treangen TJ. RefSeq database growth influences the
621 accuracy of k-mer-based species identification. *bioRxiv*. 2018.

622 35. Snitkin ES, Zelazny AM, Thomas PJ, Stock F, NISC Comparative Sequencing Program
623 Group, Henderson DK, et al. Tracking a hospital outbreak of carbapenem-resistant *Klebsiella*
624 *pneumoniae* with whole-genome sequencing. *Sci Transl Med*. American Association for the
625 Advancement of Science; 2012;4:148ra116–6.

- 626 36. Venkateswaran K, Singh NK, Checinska Sielaff A, Pope RK, Bergman NH, van Tongeren
627 SP, et al. Non-Toxin-Producing *Bacillus cereus* Strains Belonging to the *B. anthracis* Clade
628 Isolated from the International Space Station. Bik H, editor. *mSystems*. 2017;2:e00021–17.
- 629 37. van Tongeren SP, Roest HIJ, Degener JE, Harmsen HJM. *Bacillus anthracis*-Like Bacteria
630 and Other *B. cereus* Group Members in a Microbial Community Within the International Space
631 Station: A Challenge for Rapid and Easy Molecular Detection of Virulent *B. anthracis*. Schuch R,
632 editor. *PLoS ONE*. Public Library of Science; 2014;9:e98871.
- 633 38. Venkateswaran K, Checinska Sielaff A, Ratnayake S, Pope RK, Blank TE, Stepanov VG, et
634 al. Draft Genome Sequences from a Novel Clade of *Bacillus cereus* *Sensu Lato* Strains,
635 Isolated from the International Space Station. *Genome Announc*. American Society for
636 Microbiology Journals; 2017;5:e00680–17.
- 637 39. Afshinnekoo E, Meydan C, Chowdhury S, Jaroudi D, Boyer C, Bernstein N, et al. *Geospatial*
638 *Resolution of Human and Bacterial Diversity with City-Scale Metagenomics*. *CELS*. Elsevier;
639 2015;1:1–16.
- 640 40. Ackelsberg J, Rakeman J, Hughes S, Petersen J, Mead P, Schriefer M, et al. *Lack of*
641 *Evidence for Plague or Anthrax on the New York City Subway*. *CELS*. 2015;1:4–5.
- 642 41. Afshinnekoo E, Meydan C, Chowdhury S, Jaroudi D, Boyer C, Bernstein N, et al. *Modern*
643 *Methods for Delineating Metagenomic Complexity*. *CELS*. 2015;1:6–7.
- 644 42. Bazinet AL, Ondov BD, Sommer DD, Ratnayake S. BLAST-based validation of
645 metagenomic sequence assignments. *bioRxiv*. Cold Spring Harbor Laboratory; 2017;:181636.
- 646 43. Langille MGI, Ravel J, Fricke WF. “Available upon request”: not good enough for
647 microbiome data! *Microbiome*. *BioMed Central*; 2018;6:8.

- 648 44. National Research Council. Science Needs for Microbial Forensics. Developing Initial
649 International Research Priorities. Washington, D.C.: National Academies Press; 2014.
- 650 45. Bierer BE, Crosas M, Pierce HH. Data Authorship as an Incentive to Data Sharing. N. Engl.
651 J. Med. 2017;376:1684–7.
- 652 46. Credit for Data Sharing [Internet]. 2018 [cited 2018 Aug 6]. Available from:
653 <https://www.aamc.org/initiatives/research/485818/datasharing.html>
- 654 47. Meyer F, Paarmann D, D'Souza M, Olson R, Glass EM, Kubal M, et al. The metagenomics
655 RAST server – a public resource for the automatic phylogenetic and functional analysis of
656 metagenomes. BMC Bioinformatics. BioMed Central; 2008;9:386.
- 657 48. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al.
658 QIIME allows analysis of high-throughput community sequencing data. Nature methods.
659 2010;7:335–6.
- 660 49. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al. Introducing
661 mothur: open-source, platform-independent, community-supported software for describing and
662 comparing microbial communities. Applied and Environmental Microbiology. American Society
663 for Microbiology; 2009;75:7537–41.
- 664 50. ACSL.net [Internet]. [cited 2018 Aug 6]. Available from: <http://ascl.net>
- 665 51. Materials Resource Registry [Internet]. [cited 2018 Aug 6]. Available from:
666 <https://materials.registry.nist.gov>
- 667 52. Hanisch RJ, Berriman GB, Lazio TJW, Emery Bunn S, Evans J, McGlynn TA, et al. The
668 Virtual Astronomical Observatory: Re-engineering access to astronomical data. Astronomy and
669 Computing. 2015;11:190–209.

- 670 53. Field D, Sterk P, Kottmann R, De Smet JW, Amaral-Zettler L, Cochrane G, et al. Genomic
671 standards consortium projects. *Stand Genomic Sci.* Michigan State University; 2014;9:599–601.
- 672 54. Field D, Amaral-Zettler L, Cochrane G, Cole JR, Dawyndt P, Garrity GM, et al. The Genomic
673 Standards Consortium. *Plos Biol.* Public Library of Science; 2011;9:e1001088.
- 674 55. Yilmaz P, Kottmann R, Field D, Knight R, Cole JR, Amaral-Zettler L, et al. Minimum
675 information about a marker gene sequence (MIMARKS) and minimum information about any (x)
676 sequence (MlxS) specifications. *Nat. Biotechnol.* Nature Publishing Group; 2011;29:415–20.
- 677 56. Titus AJ, Flower A, Hagerty P, Gamble P, Lewis C, Stavish T, et al. SIG-DB: leveraging
678 homomorphic encryption to Securely Interrogate privately held Genomic DataBases. *arXiv*
679 *Quantitative Methods.* 2018.
- 680 57. Ellington MJ, Ekelund O, Aarestrup FM, Canton R, Doumith M, Giske C, et al. The role of
681 whole genome sequencing in antimicrobial susceptibility testing of bacteria: report from the
682 EUCAST Subcommittee. *Clinical Microbiology and Infection.* Elsevier; 2017;23:2–22.