1  **Enabling the democratization of the genomics revolution with a fully integrated web-**

2  **based bioinformatics platform**

3

4  Po-E Li*[1], Chien-Chi Lo*[1], Joseph J. Anderson[2,3], Karen W. Davenport[1], Kimberly A. Bishop-

5  Lilly[3,4], Yan Xu[1], Sanaa Ahmed[1], Shihai Feng[1], Vishwesh P. Mokashi[3], and Patrick S. G. Chain[1]

6

7  [1]Biome Sciences, Bioscience Division, Los Alamos National Laboratory, Los Alamos, NM

8  87545; [2]Defense Threat Reduction Agency, Fort Belvoir, VA; [3]Naval Medical Research Center –

9  Frederick, Fort Detrick, MD 21702; [4]Henry M. Jackson Foundation, Bethesda, MD 20817

10

11  *Contributed equally to the work. Correspondence should be addressed to P.S.G.C.

12

13  Patrick S. G. Chain

14  Bioinformatics and Analytics Team,

15  Bioscience Division

16  Los Alamos National Laboratory

17  Los Alamos, NM 87545

18  Tel.:  505-665-4019

19  Fax:  505-665-3024

20  Email: pchain@lanl.gov

21

22  Running title: NGS bioinformatics with EDGE

23

24  Keywords: cloud compatible bioinformatics pipelines and workflows, integrated web-based

25  genome analysis platform, isolate genome and metagenome NGS data analysis and

26  visualization, next generation sequencing genomics and metagenomics, EDGE bioinformatics

27

## ABSTRACT

29  Continued advancements in sequencing technologies have fueled the development of new

30  sequencing applications and promise to flood current databases with raw data. A number of

31  factors prevent the seamless and easy use of these data, including the breadth of project goals,

32  the wide array of tools that individually perform fractions of any given analysis, the large number

33  of associated software/hardware dependencies, and the detailed expertise required to perform

34  these analyses. To address these issues, we have developed an intuitive web-based

35  environment with a wide assortment of integrated and cutting-edge bioinformatics tools. These

36  preconfigured workflows provide even novice next-generation sequencing users with the ability

37  to perform many complex analyses with only a few mouse clicks, and, within the context of the

38  same environment, to visualize and further interrogate their results. This bioinformatics platform

39  is an initial attempt at Empowering the Development of Genomics Expertise (EDGE) in a wide

40  range of applications.

41

## INTRODUCTION

43  The field of genomics has made tremendous technological leaps in recent years, and the

44  combined decrease in sequencing costs and expansion in applications (transcriptomics,

45  metagenomics, single cell genomics) have truly revolutionized the way scientists approach

46  biological questions (for a recent review, see (Buermans and den Dunnen 2014)). Now that a

47  trained technician can single-handedly produce gigabases of sequence data in essentially a

48  day's work, "next generation sequencing" (NGS) is being applied by many smaller laboratories,

49  as well as the large traditional sequencing centers, across a wide range of disciplines in order to

50  answer a variety of complex problems. For instance, NGS is being applied to the

51  characterization and attribution of outbreaks in clinical environments (Conlan et al. 2014), food

52    safety (den Bakker et al. 2014), the development of alternative energy sources (Wang et al.

53    2012; Wohlbach et al. 2014), and many other fields.

54

55    Although many advances have been made in bioinformatics methods development, the so-

56    called "democratization of genomics" (Koren et al. 2014) has not yet fully expanded to the

57    bioinformatic realm, making it difficult for investigators to adequately analyze genomic big data

58    (Daber et al. 2013; Watson-Haigh et al. 2013). While NGS no longer seems new, it has really

59    only been since 2005 that a revolutionary new technology (pyrosequencing) (Margulies et al.

60    2005) was introduced after more than twenty years of chemical degradation (Maxam and Gilbert

61    1977) and chain termination (Sanger et al. 1977) sequencing. Some of these NGS technologies

62    have already been abandoned even after strong market performance; other new technologies

63    are only now emerging, and the ones that have thus far survived continue to undergo

64    improvement. Despite reads of limited length, Illumina[®] (Bennett 2004) currently dominates the

65    market, in part due to its very high throughput and low cost.

66

67    Analysis of the massive datasets produced in NGS studies and interpretation of the results

68    requires expertise in both computer science and biology. Therefore, although the decreasing

69    cost and decreasing laboratory footprint of NGS technologies make the production of these

70    datasets a more realistic goal for many laboratories, there still remain at least three core issues

71    in bioinformatics that hamper the broader use of NGS data. First, the numerous and diverse

72    specific questions being asked of NGS data require highly specialized pipelines. While any

73    given question can sometimes make use of the same basic tool(s) with different parameters and

74    post-processing, other questions may require similar bioinformatic manipulation but are

75    optimally answered using different tools, and further questions may require entirely new

76    methods or algorithms. Second, there is the related issue of having numerous available (and

77    somewhat redundant) options for extremely complex NGS bioinformatics data analysis tools.

78     Because NGS data and their formats frequently change, the analytical tools must adapt; new

79     tools arise frequently through efforts to improve upon initially developed algorithms, or to

80     complement other methods. One can often identify dozens or even hundreds of individual tools

81     that can perform the same type of analysis, and it has been an increasing challenge to decide

82     which tools are best for which specific applications. In addition, some tools are tailored to

83     specialized hardware architectures. Lastly, few laboratories have the degree of expertise

84     required to implement robust methods, install the appropriate tools, or construct standardized

85     pipelines for processing data. The need for such expertise can delay studies and make

86     comparisons of disparate studies very difficult. While some systems have allowed the open-

87     source integration of selected tools within a single environment (e.g., Galaxy (Blankenberg et al.

88     2010)), users must often already know which tools or pipelines to select and what specific

89     parameters to use for their particular goals. A more costly approach includes commercial

90     packages that can perform similar operations and further help to visualize results, but these

91     packages use proprietary software that can be inflexible or, if one does not know the details of

92     the programs and parameters, can affect downstream interpretation.

93

94     Because we view bioinformatics as the key bottleneck in the use of NGS data, we present an

95     integrated platform toward Empowering the Development of Genomics Expertise (EDGE). This

96     bioinformatics effort is intended to truly democratize the use of NGS for exploring novel

97     genomes and metagenomes. We developed EDGE Bioinformatics as an initial suite of pre-

98     configured bioinformatics workflows that allow rapid analysis of NGS data, coupled with

99     visualization and interactive features. These features allow users to view results and explore

100     ongoing data processing within an intuitive and user-friendly web-based environment. The

101     software is freely available (https://lanl-bioinformatics.github.io/EDGE/) and a webserver is

102     provided (https://bioedge.lanl.gov/) for use with publicly available data via the NCBI Sequence

103     Read Archive.

104

**RESULTS**

105

**The EDGE Bioinformatics overview**

106

107 An overview of the EDGE Bioinformatics workflow is shown in Figure 1, with a more detailed

108 workflow shown in Supplementary Figure S1.  Because most sequencers can now output data

109 as one or more FASTQ files we opted for this format (full or compressed) as the required input

110 for raw sequencing data. EDGE can use files derived from multiple libraries, runs or lanes by

111 specifying the location of one or more FASTQ files or by retrieving them from the Sequence

112 Read Archive (SRA) at NCBI (Supplementary Figure S2). EDGE was originally designed for use

113 with Illumina$^®$ reads and performs best with these short sequence data types, but the

114 development of alternative workflows are envisioned for future versions to better handle other

115 types of data (e.g. longer reads, different error models, etc.). There are a number of additional

116 options such as specifying number of CPUs to use or allowing batch submission of many

117 samples.
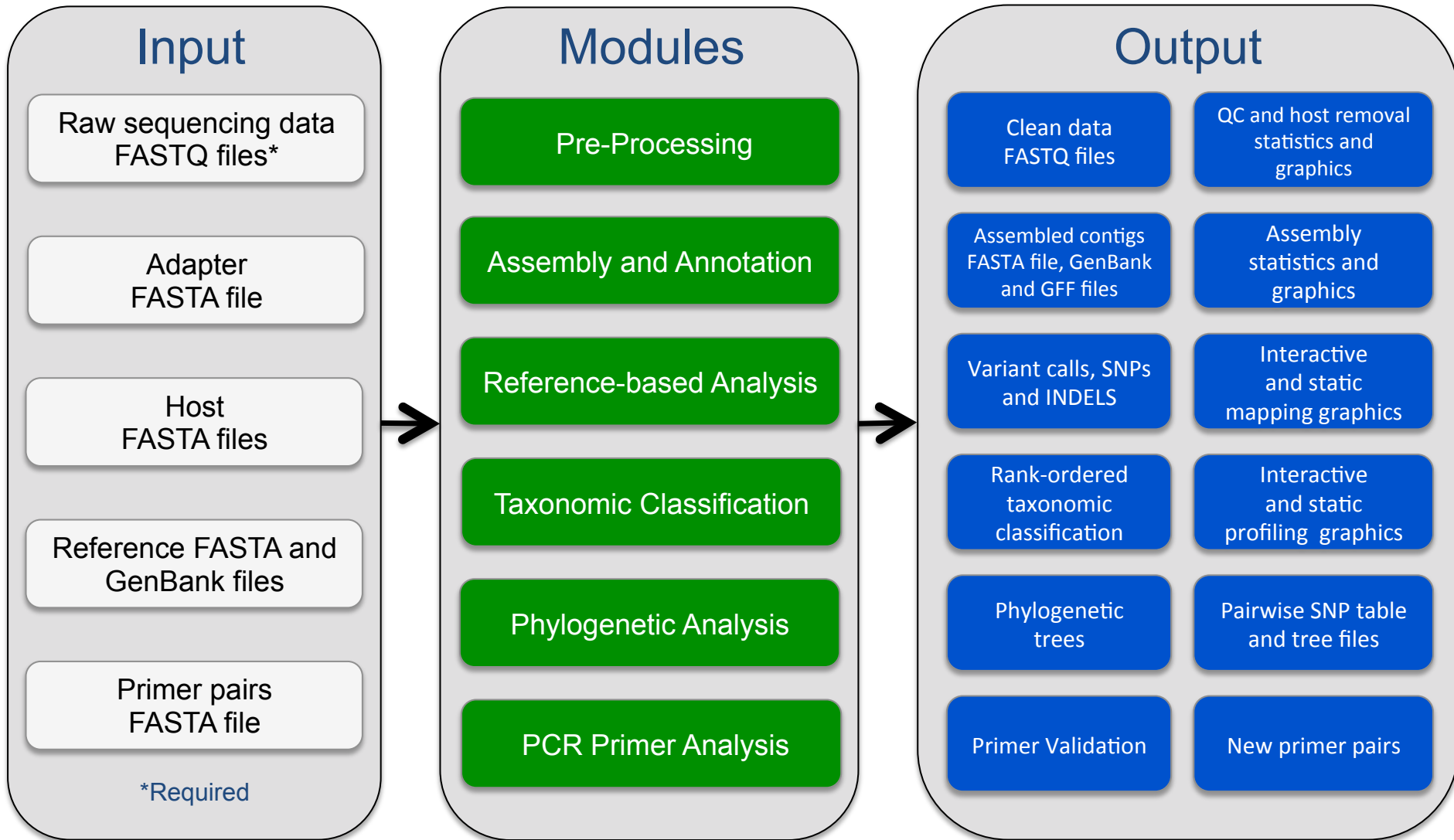
118

# The EDGE Environment



**Figure 1. An overview of the EDGE Bioinformatics Environment.** The only Inputs required from the user are raw sequencing data and a project name. The user can create specific workflows with any combination of the modules. In addition, tailored parameters dictating how each module functions can be modified by the user. EDGE outputs a variety of files, tables and graphics which can be viewed on screen or downloaded. A more detailed overview is shown in Supplementary Figure S1. All Modules are described in the Methods section.

119

120    Optional inputs depend on the selected modules (see Methods) and can include an adapter

121    FASTA file for adapter filtering, a host FASTA file for removal of host reads, PacBio/Nanopore

122    long read FASTA/FASTQ files for use with the SPAdes (Bankevich et al. 2012) assembler, one

123    or more reference genomes for comparative genomic analysis, and a primer pair(s) file in

124    FASTA format for *in silico* primer validation. While there are several optional environmental

125    parameters that can control the way EDGE runs, the users need only specify a project name,

126    select the input file(s), toggle which modules they would like to use, and click Submit. The

127    results of each project are displayed within its own project page (see Methods and

128    Supplementary Figure S3). Descriptions of all modules are in the Methods section.

129

130    **Analysis in EDGE**

131    To demonstrate the utility and versatility of EDGE, we tested this platform using a number of

132    different samples that represent varied scenarios, including examples of isolate sequencing and

133    analysis of several clinical metagenome samples with known, suspected, and unknown etiologic

134    agents (Table 1). Not all results are described in depth, but the different datasets are used to

135    highlight some of the various modules and analytic capabilities encompassed within the EDGE

136    Bioinformatics platform. All datasets and project pages with full results are publicly available on

137    our webserver (https://bioedge.lanl.gov/). There, users can view or select and run their own

138    analyses of these data or other publicly accessible SRA data.

139

**Table 1. Descriptions of samples and EDGE modules tested.**

| Sample description | Sample type (material) | # of reads (millions) | Sequence type | EDGE Modules* | | | | | | CPUs | Run Time (hours) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | 1 | 2 | 3 | 4 | 5 | 6 | | |
| *Bacillus anthracis* strain SK-102 SRR1993644 | Isolate (gDNA) | 28.6 | HiSeq 2x101 nt | X | X | X | X | X | X | 8 | 4:12:03 |
| *Bacillus anthracis* strain SK-102 SRR1993644 | Isolate (gDNA) | 28.6 | HiSeq 2x101 nt | X | X | X | X | X | X | 20 | 3:33:52 |
| *Yersinia pestis* strain Harbin 35 SRR1993645 | Isolate (gDNA) | 15 | GAII 2x110 nt | X | X | X | X | X | X | 8 | 3:35:39 |
| Human Microbiome Project (staggered mock community) SRR172903 | Metagenome (DNA) | 7.93 | GAII 75 nt | X | X | | X | | | 8 | 0:53:59 |
| Patient plasma sample 2014 *Ebola* outbreak (IDBA assembly) SRR1553609** | Metagenome (RNA) | 0.93 | HiSeq 2x100 nt | X | X | X | X | | | 12 | 0:38:07 |
| Patient plasma sample 2014 *Ebola* outbreak (SPAdes assembly) SRR1553609** | Metagenome (RNA) | 0.93 | HiSeq 2x100 nt | X | X | X | X | | | 12 | 0:47:24 |
| Patient fecal sample 2011 *E. coli* outbreak SRR2164314 | Metagenome (DNA) | 273 | HiSeq 2x100 nt | X | X | | X | | | 8 | 34:43:30 |
| Patient nasal swab acute respiratory illness SRP062772** | Metagenome (DNA) | 2.52 | MiSeq 2x300 nt | X | X | | X | | | 8 | 0:20:59 |

* EDGE Modules are described in Methods: 1. Pre-Processing; 2. Assembly and Annotation; 3. Reference-Based Analysis; 4. Taxonomic Classification; 5. Phylogenetic Analysis; 6. PCR Primer Analysis
** These samples were retrieved directly from the NCBI SRA.

140

141 **Analysis of isolate genome sequencing projects**

142 To highlight and validate some of the features and integration of utilities within EDGE, we tested

143 the various modules using two datasets (sequenced at two different institutions) from recently

144 completed isolate genome sequencing projects: *Bacillus anthracis* strain SK-102 (Johnson et al.

145 2015b) and *Yersinia pestis* strain Harbin 35 (Johnson et al. 2015a). After quality control, 96-98%

146 of the reads were retained for *B. anthracis* and *Y. pestis* (Supplementary Figure S4). Results

147 from the Assembly and Annotation module were consistent with known genome complexity

148 (repeated elements such as insertion sequences and rRNA operons), genome size, and

149 associated number of genes. The *B. anthracis* assembly was 5.5 Mb in size, consisting of 89

150 contigs with a maximum contig size of 450kb and an average contig fold coverage of 328X,

151 consistent with the amount of data sequenced (Supplementary Figure S5). The *Y. pestis*

152 assembly (4.6 Mb with 306X fold coverage) was more fragmented (329 contigs) with smaller

153 contig sizes (maximum contig size of 115kb) owing to the large number of repeat sequences

154 within the genome. However, using the reference-based analysis module, all of the *Y. pestis*

155 contigs, and all but a single contig of the *B. anthracis* assembly, could be mapped to the

156 selected reference genome (*Y. pestis* CO92 and *B. anthracis* Ames Ancestor, respectively).

157 More than 98% of the reads of either sample could also be mapped, covering >97-100% of the

158 reference chromosomes and plasmids (Supplementary Figure S6).

159

160 While the identities of the organisms sequenced in this case are not in question, the taxonomy

161 classification module can be used to identify a contaminant, or otherwise suggest similarity to

162 another taxon. The consensus for all the taxonomy classification tools encompassed in EDGE

163 confirmed the presumed identities of the organisms sequenced. With *Y. pestis*, both GOTTCHA

164 (Freitas et al. 2015) and Metaphlan (Segata et al. 2012) provided the cleanest results,

165 suggesting only *Y. pestis* reads comprise the dataset (Figure 2A), however with *B. anthracis*, a

166    number of different organisms were found by these tools (Figure 2B), even at the genus level.

167    At the species level, both GOTTCHA and Metaphlan identified *B. cereus* and *Francisella*

168    *philomiragia* in addition to the dominant *B. anthracis*. In addition, GOTTCHA found signatures of

169    *Y. pestis* and *B. weihenstephanensis*, while Metaphlan suggested *B. thuringiensis* was present.

170    Upon further investigation, we discovered that the *B. anthracis* SK-102 sample was sequenced

171    within the same Illumina lane as many other samples, including *F. philomiragia* ATCC25018,

172    two *Y. pestis* strains (771 and 790), *B. cereus* BACI291, *B. mycoides* BACI084 (a near neighbor

173    to *B. weihenstephanensis* (Soufiane and Cote 2013)), and several fecal samples from Condors

174    (found to contain dominant amounts of *Clostridia* sequences, consistent with dominance of

175    *Clostridia* in the Vulture hindgut (Roggenbuck et al. 2014)). Therefore these additional

176    identifications are likely the result of index cross contamination (or other mis-assignment) of

177    barcodes to sample, often found among samples run within the same lane (Kircher et al. 2012).

178    In addition, and consistent with the bacteria in this sample, GOTTCHA viral analysis suggested

179    three *Bacillus* phages as well as *Staphylococcus* phage SpaA1, which is similar to *Bacillus*

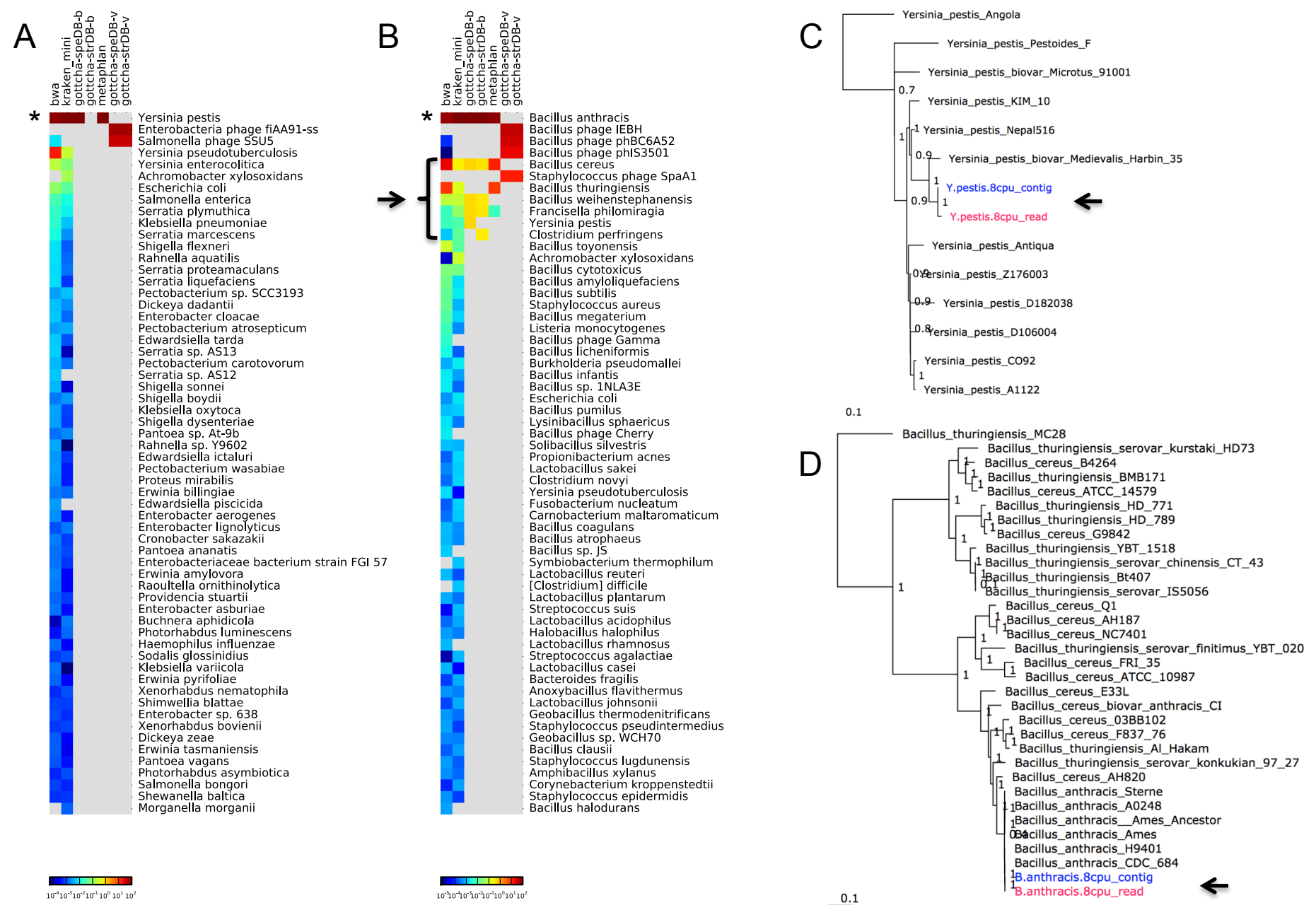180    prophages and can infect *Bacillus* spp. (Swanson et al. 2012).

181

**Figure 2. Taxonomy and phylogenetic evaluations of bacterial isolates.** Panels A and B show taxonomic classification of reads for A) the *Y. pestis* Harbin35 sample and B) the *B. anthracis* SK-102 sample. The stars indicate the consistent dominant taxonomic calls for all tools, while the black arrow and bracket indicate identified contamination in the *B. anthracis* sample. Panels C and D indicate the inferred phylogenetic trees for the C) *Y. pestis* and D) *B. anthracis*; black arrows point to the read dataset (pink) and contigs (blue) that were placed in these trees.

182

183 Phylogenetic analysis was performed for each dataset, selecting all available NCBI RefSeq

184 genomes for either *Y. pestis*, or for *B. anthracis*, *B. cereus*, and *B. thuringiensis*. This

185 phylogenetic module, based on PhaME (Ahmed et al. 2015),  independently treats the input

186 reads and resulting contigs (when assembly is selected) for whole genome SNP analysis, and

187 consistently placed the datasets within their respective phylogenetic trees (Figure 2C-D). The *Y.*

188 *pestis* tree was inferred from a 4.0 Mb core genome with 2,077 SNPs and the *Y. pestis* sample

189 was placed nearest a previously sequenced *Y. pestis* Harbin35. The *Bacillus* tree was based on

190 a core genome of 3.1 Mb with 384,568 SNPs, is fully consistent with known *Bacillus*

191 relationships (Soufiane and Cote 2013), and placed the reads and the resulting contigs of the *B.*

192 *anthracis* SK-102 closest to *B. anthracis* CDC684.

193 Using the PCR Primer Tools module, published primers that have been used to detect either *Y.*

194 *pestis* (Hinnebusch and Schwan 1993; Begier et al. 2006) or *B. anthracis* (Fasanella et al. 2003;

195 Francy et al. 2009; 2012) were input for validation against these isolates and confirmed the

196 appropriate amplicon sizes using electronic PCR against the respective assemblies. For *B.*

197 *anthracis*, two novel PCR primer pairs were suggested by the primer design software that would

198 specifically amplify only this strain compared with all other NCBI genomes (Supplementary

199 Figure S7).

200

201 **Analysis of a mock human microbiome sample of known complexity.**

202 The Human Microbiome Project's (HMP) staggered mock community (Human Microbiome

203 Project 2012) was used to evaluate the metagenome analysis potential of EDGE. This dataset,

204 consisting of sequencing reads derived from a mixture of 21 known bacterial strains and one

205 eukaryotic strain, was analyzed using the Pre-processing, Assembly, and Taxonomy

206 classification modules with default parameters. The FaQCs (Lo and Chain 2014) quality control

207 pipeline retained 81.2% of the reads and 76.7% of the data from the 7.9M read dataset, while

9

208    the subsequent assembly produced 13,097 contigs totaling 14.8 Mb. Read mapping validation

209    suggested that the assembly represents 77.6% of the reads with a contig average fold coverage

210    of 24X (Supplementary Figure S8). Both the read- (Figure 3A), and contig-based (Figure 3B)

211    taxonomy classification tools accurately identified most of the known community members of

212    this sample with the exception of the eukaryote since these tools were implemented to identify

213    bacteria and viruses. The contig plot of average G+C (%) versus average fold coverage can

214    also help distinguish groups of contigs that belong to different organisms (Figure 3C). Similar

215    graphics and results can be found at various taxonomic levels.
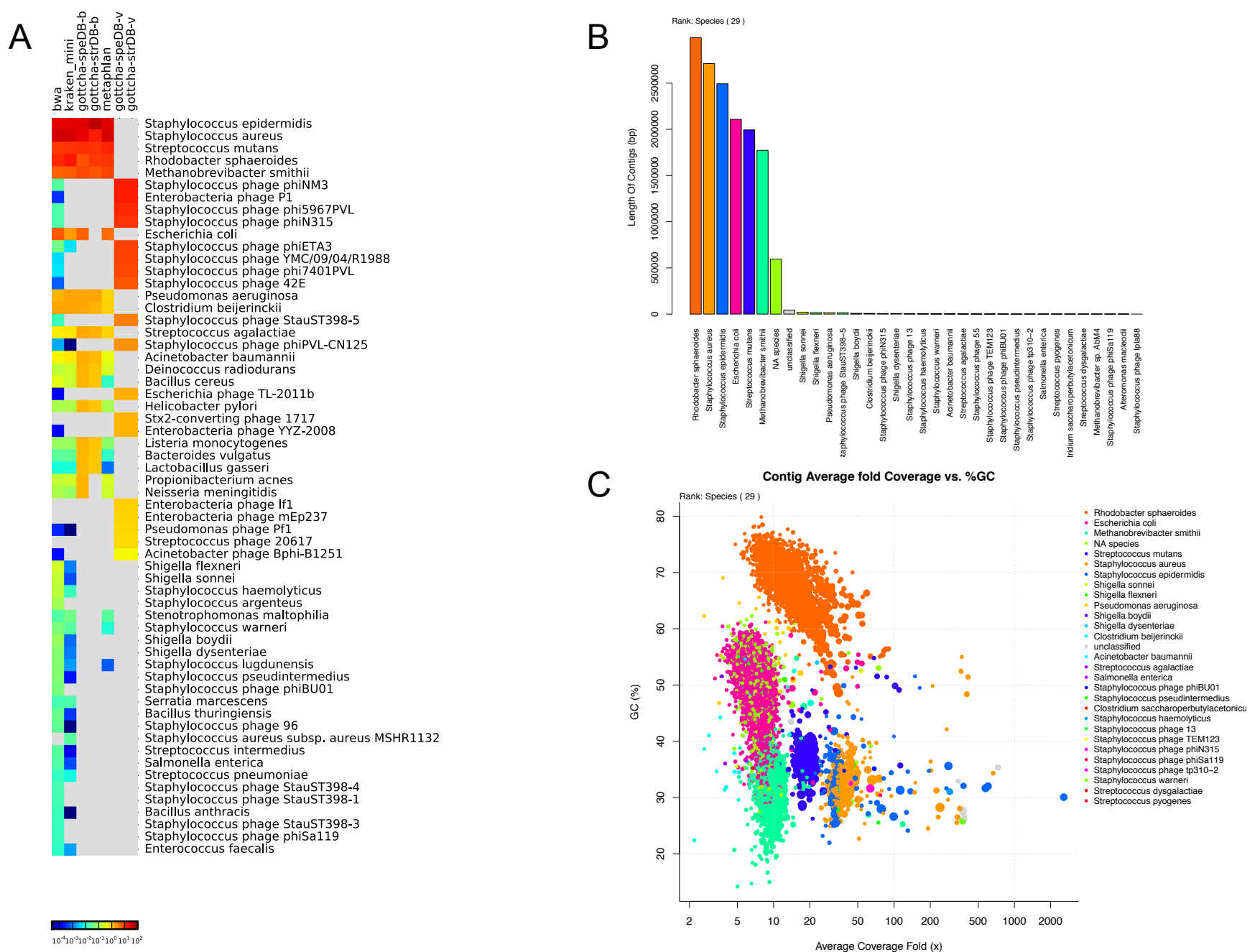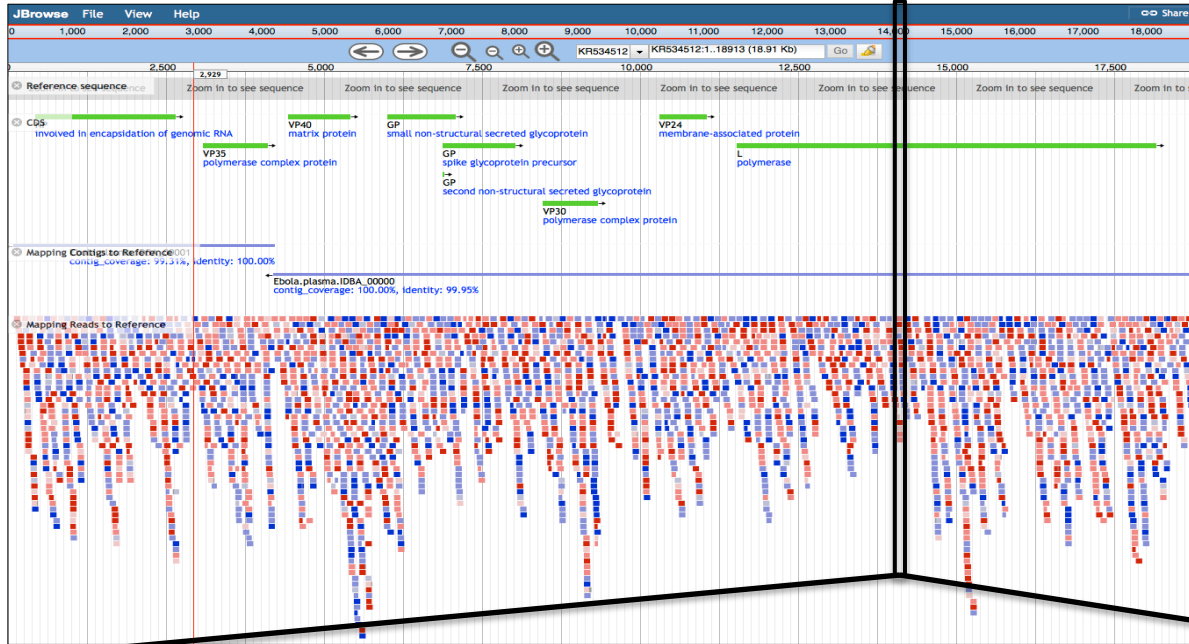
216

**Figure 3. Taxonomic Classification of the HMP staggered mock sample.** A) Read-based classification using various taxonomy profiling tools; B) contig-based classification displaying length of all classified contigs per taxon; and C) a scatterplot of contig % GC vs. fold coverage of the contigs, colored by taxon.

217

218    **Analysis of complex clinical samples.**

219    We also used EDGE to evaluate datasets from several clinical samples with suspected

220    pathogens. In the first example we used EDGE to characterize one of the recent 2014 Ebola

221    outbreak samples. Using the Sierra Leone human plasma RNA sequencing sample

222    SRR1553609 retrieved directly from the SRA, we ran all EDGE modules with the exception of

223    phylogenetic and primer analyses. Pre-processing removed ~25% of the data, and human host

224    removal only identified 605 reads that matched the human reference. IDBA (Peng et al. 2012)

225    assembly of the remaining reads resulted in 1588 contigs, a total assembly size of 665kb and a

226    largest contig of 14.6kb. Due to the complexity of the sample, only 15% of the data assembled.

227    We examined the use of the alternate assembler, SPAdes (Bankevich et al. 2012), with this

228    sample and found an increased run time (Table 1) balanced by an improved 36% read

229    incorporation (vs. 15%) into the assembly, resulting in 12,105 contigs, a total assembly size of

230    >3.8Mb and a largest contig of 18.6kb. Using as reference the *H.sapiens*-wt/GIN/2014/Makona-

231    Gueckedou-633 Zaire ebolavirus (a sequence from Guinea, 2014), we found that only 3,228

232    reads (0.43% of the input reads) could be mapped to the genome, covering 98.9% of the length

233    with 10 potential single nucleotide variants. Two of the IDBA contigs overlapped and together

234    covered 99.2% of the genome, while a single SPAdes contig covered 97.8% of the reference.

235    Both assemblies identified the same 8 SNPs with respect to the reference genome. The

236    genome browser in EDGE helped resolve the disparate variant analysis found between the

237    reads and the contigs (Figure 4). While almost all of the reads confirmed all 8 SNPs found within

238    the contigs, the two additional variants identified with read-based analysis likely reflected the

239    quasispecies nature of the virus, with strong support but fewer than 50% of the reads at those

240    positions carrying the additional point mutations. This shows the utility of a multi-pronged

241    approach when performing such comparisons. The taxonomy classification module showed that

242    Ebola could indeed be found within the reads, though only with the GOTTCHA and BWA

243    pipelines, and also provided a list of bacteria that may have also been present within the

244    sequenced sample, including *Ralstonia*, *Bradyrhizobium*, *Propionibacterium*, and *Pseudomonas*

245    (Supplementary Figure S9). The contig-based taxonomy analyses also clearly showed Ebola

246    virus to be present, and confirmed that many contigs belonged to the same bacterial groups

247    identified by read-based analyses.

248

**Figure 4. Interactive genome browsing view of a reference-based analysis in EDGE with a human clinical sample containing Ebolavirus.** A) An Ebola reference genome and its genes (green lines) are displayed together with contig-based (using IDBA) and read-based comparisons. The two contigs (blue lines) from IDBA are shown aligned along the length of the reference as well as the reads (red and blue). B) A zoomed-in view of one section of the genome where SNPs were identified. The SNP and coding difference is outlined under the contig alignment, while the variants are indicated under the read alignments.

249

250    In the second clinical example, we analyzed data derived from a fecal sample of a patient

251    returning from Germany during the 2011 enterohemorrhagic *Escherichia coli* outbreak, and who

252    was suspected of harboring *E. coli* O104:H4. Trimming and filtering removed 13.3% of the

253    bases while host removal identified only 0.15% of the reads as human and 0.02% as PhiX (a

254    spike-in control commonly used in Illumina sequencing). Assembling the remaining 253M reads

255    resulted in 2,957 contigs totaling 10.5 Mb, comprising 23.9% of the reads. The single

256    chromosome and three plasmids of *E. coli* O104:H4 2011C-3493, were used as reference for

257    both read- and contig-based comparisons. Using reads, 99.99% of the reference chromosome

258    was covered at 115X, while the three plasmids were covered 100% at fold-coverages ranging

259    from 250X for the largest plasmid to 7.6 million fold coverage for the smallest plasmid. Using

260    contigs, all replicons were covered >99.7% with the exception of the small plasmid which was

261    absent from the assembly (this absence is likely due to the excessive fold coverage known to

262    create assembly issues). All taxonomy profiling tools clearly showed that *E. coli* (or *Shigella*)

263    was the dominant organism and that the Shiga-toxin phage was also present (Supplementary

264    Figure S10). Whole genome SNPs were identified and phylogenetic analysis was performed

265    with both reads and contigs, easily done within EDGE using the drop down menu to select 68 *E.*

266    *coli* and *Shigella* genomes. Both the predominantly *E. coli* metagenome reads and the

267    assembled contigs were placed within the same clade as the other *E. coli* O104 strains,

268    reaffirming the initial suspicion of *E. coli* O104:H4 as the etiologic agent (Figure 5A).

269

270    A nasal swab sample from a patient with acute respiratory illness of unknown etiology was used

271    as a final test of EDGE's utility for analysis of clinically derived metagenomic datasets. In this

272    case, while >99% of the data passed FaQCs quality control, the majority of sequence reads

273    (78.9%) were human-derived and removed (data not shown). The remaining reads were

274    submitted to SRA and used for assembly and taxonomy classification. A number of expected

275    organisms (Rawlings et al. 2013; Bassis et al. 2015) ranked among the most abundant genera

276    identified, including *Prevotella, Veillonella* and *Streptococcus*. Unexpectedly, *E. coli* was

277    identified by GOTTCHA, and also detected (at a substantially lower level) by BWA and Kraken

278    mini (Figure 5B). Upon closer inspection, the mapping results demonstrated that all of the *E. coli*

279    hits were to the plasmid (with no matches to the chromosome) in *E. coli* strain ABU83972,

280    covering approximately 80% of this replicon. Interestingly, this plasmid is very similar (>90%

281    identity) to a number of enteric plasmids, as well as to the *Corynebacterium renale* plasmid

282    pCR1, suggesting that the presence of this plasmid might be the result of colonization or

283    infection by a *Cornyebacterium* species, which are common in nasal cavities (Bassis et al.

284    2015). This hypothesis is partially supported by BWA and Kraken, which identified a different

285    *Cornyebacterium* at low levels, as well as by 16S sequence data in which *E. coli* is not detected

286    but the genus *Cornyebacterium* is found (Supplementary Table S1). As a result of these findings

287    a new feature now present in EDGE separates plasmid from chromosomal hits for GOTTCHA,

288    thereby allowing for greater specificity in evaluating taxonomic profiling results (Figure 5C). The

289    differences in bacterial species found by Metaphlan compared with all other tools can be

290    explained by the additional draft genome references included within the Metaphlan database

291    (Segata et al. 2012) and which are not yet available in RefSeq.
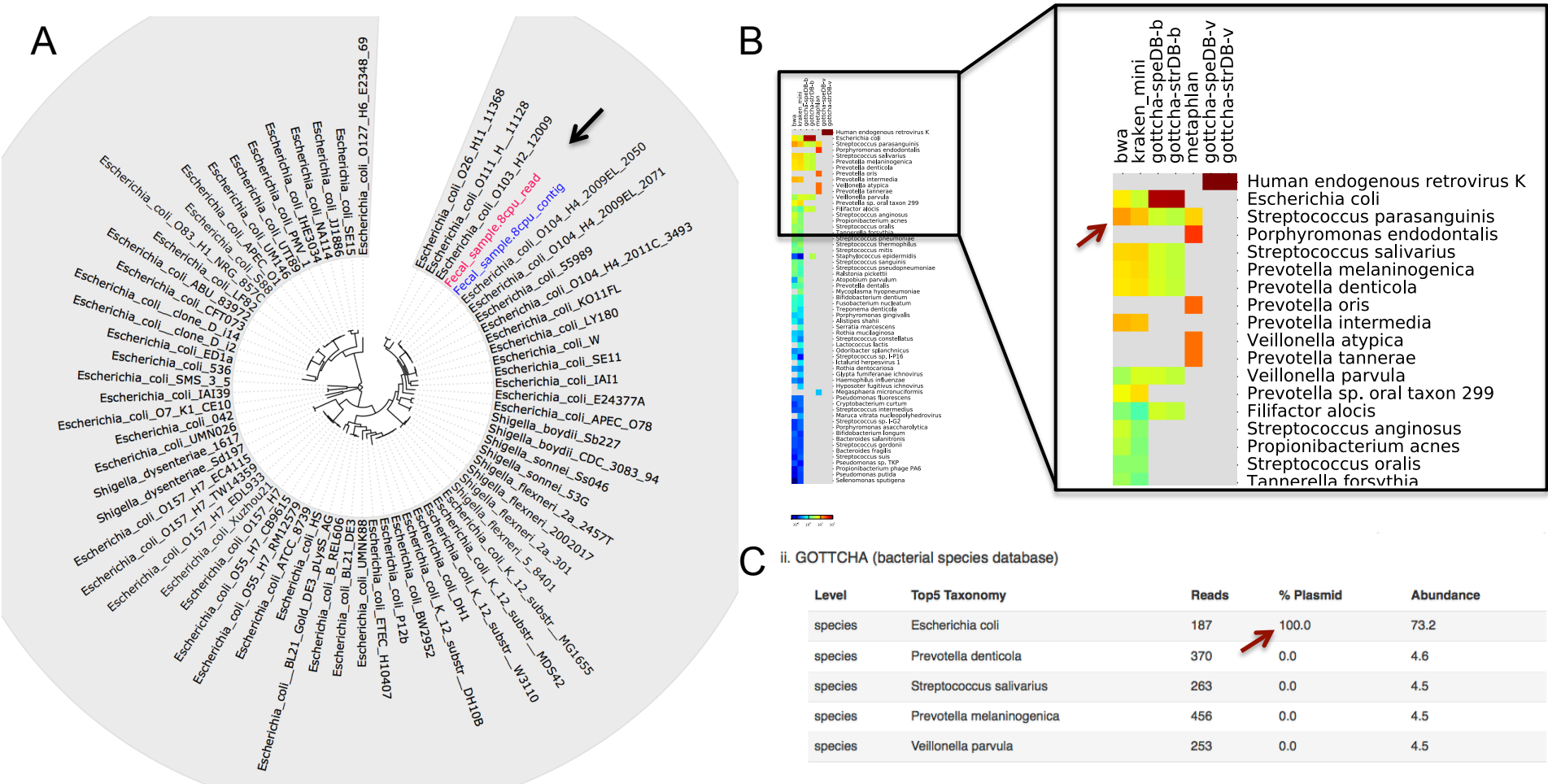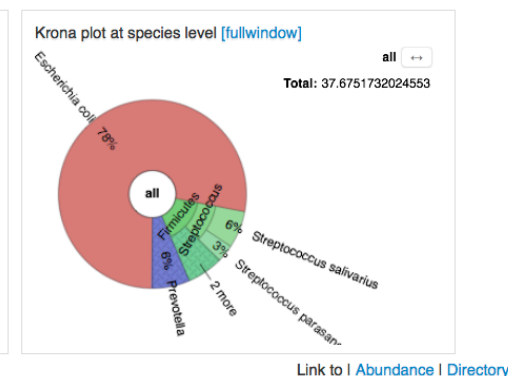
292

**Figure 5. Phylogenetic and taxonomic analysis of human clinical samples with suspected and unknown causative agents.** A) Circular phylogenetic tree clearly places within the *E. coli* O104 group both the raw reads and the contigs obtained from a clinical fecal sample. B) A comparative heatmap view of identified taxa from a nasal swab sample demonstrates the abundance of typical nasal cavity organisms. C) The *E. coli* identified with GOTTCHA in the nasal swab sample (in B) is described in greater detail under the tool-specific EDGE view (red arrow), showing the percent of hits to plasmids for each identified taxon; below are a taxonomic dendrogram featuring the taxa detected with circles representing relative abundance, and a Krona plot view of the same data.

293

294    **DISCUSSION**

295    As the number of investigations that apply sequencing continues to climb, the wider genomics

296    community will greatly benefit from a user-friendly bioinformatics environment of integrated tools

297    and pipelines designed to address a large number of scenarios and scientific end-goals. The

298    initial system and the tools we developed and used in EDGE are available as open source

299    software, and we encourage other developers to contribute best-practice tools and pipelines, as

300    there are yet a number of use cases not addressed within this initial platform. For the tools in

301    current use, the focus was on accuracy, speed, flexibility, and ability to run within a modest

302    computational environment. In some cases, like with read-based taxonomy profiling, given that

303    this is a still emerging field of exploration, we provide a suite of tools based on different

304    algorithms, and present a comparative view of the results for further scrutiny by researchers. In

305    other cases, tools were selected that perform well under a diverse set of circumstances, and are

306    computationally friendly with respect to speed and memory considerations. While novel tools

307    continue to be developed and databases continue to grow, future focus will be on the systematic

308    incorporation of better tools and updating of databases alongside the development of new

309    modules.

310

311    Collectively, our results and experiences suggest that EDGE provides significant advantages

312    over the current status quo. For example, significant expertise is generally required to determine

313    what tools (and parameters) are ideal for different scenarios, and in many cases, to run these

314    tools and manipulate the results. EDGE assists non-expert users by providing pre-defined

315    pipelines to run cutting-edge tools and a web interface that makes inspection of results quick

316    and easy. Comparative views of results output by complex metagenome taxonomy profiling

317    tools distinguish this system from all others along with the ability to easily perform whole

318    genome SNP phylogenies with user-selected genomes. The ability to integrate read-based with

15

319  assembly-based analyses provides complimentary views of genomic data. Real time tracking of

320  projects and system resources allows for better monitoring and job queuing. With embedded log

321  files detailing the specifics of each run, a wide adoption of systems like EDGE can also provide

322  a form of standardized data analysis which would allow for more robust comparisons to be

323  made across different independent projects and laboratories.

324

325  EDGE is a unique bioinformatic software package both for the variety of open-source tools that

326  are encompassed and for its ease of use. To our knowledge, there is no other freely available

327  bioinformatic software package that incorporates these types of analyses and tools within a

328  framework of intuitive pipelines and interactive graphical and tabular results. This software

329  package is designed to enable scientists with limited experience in bioinformatics to perform a

330  variety of genomic analyses with resources that are available in smaller laboratories, rather than

331  requiring extensive computational and personnel infrastructure. The EDGE Bioinformatics

332  software therefore represents a critical step forward in democratizing genomics analyses.

333

334  **METHODS**

335  **EDGE Bioinformatics computational design**

336  EDGE Bioinformatics is built around a collection of publicly available, open-source software

337  packaged in six modules. The main wrapper script is written in Perl, while the various tools

338  currently include BLAST, version 2.2.26 (Altschul et al. 1990), BowTie2, version 2.1.0

339  (Langmead and Salzberg 2012), BWA, version 0.7.9 (Li and Durbin 2010), FaQCs, version 1.33

340  (Lo and Chain 2014), FastTree, version 2.1 (Price et al. 2010), GOTTCHA, version 1.0b (Freitas

341  et al. 2015), IDBA_UD, version 1.1.1 (Peng et al. 2012), SPAdes, version 3.5.0 (Bankevich et al.

342  2012), JBrowse, version 1.11.6 (Skinner et al. 2009), jsPhyloSVG, version 1.55 (Smits and

343  Ouverney 2010), Kraken, version 0.10.4-beta (Wood and Salzberg 2014), KronaTools, version

344  2.4 (Ondov et al. 2011), MetaPhlAn, version 1.7.7 (Segata et al. 2012), MUMmer3, version 3.23

345 (Kurtz et al. 2004), Phage_Finder, version 2.1 (Fouts 2006), PhaME (Ahmed et al. 2015),

346 Primer3, version 2.3.5 (Untergasser et al. 2012), Prokka, version 1.11 (Seemann 2014), RATT,

347 version 08-Oct-2010 (Otto et al. 2011), RAxML, version 8.0.26 (Stamatakis 2014), and

348 SAMtools, version 0.1.19 (Li et al. 2009).

349

350 All tools and modules can be run on Unix command line, however we provide a user-friendly

351 web-based graphic user interface (GUI). The GUI is primarily implemented using the JQuery

352 Mobile javascript framework and HTML5 on the client-side, and implements perl CGI using

353 Apache or Python on the server-side. This implementation makes EDGE accessible on any

354 platform, including all smartphones, tablets, and desktop devices. The EDGE software tools

355 were selected or developed based on the desire (and need) for both accuracy and speed, with

356 the assumption of moderate computational hardware resources. More detail regarding the

357 installation, implementation, and the tools encompassed within EDGE can be found at

358 http://edge.readthedocs.org/.

359

360 The modular design and open source license also allow other researchers to expand the

361 available capabilities beyond our initial implementation. For expert bioinformaticians, another

362 benefit is that EDGE can also be integrated into other workflows and be used via command line

363 to submit jobs on a cluster. More information can be found at the EDGE homepage (https://lanl-

364 bioinformatics.github.io/EDGE/), and the software is available at https://github.com/LANL-

365 Bioinformatics/edge. We also offer a partially modified, public EDGE webserver, available at

366 https://bioedge.lanl.gov/, which can be used to analyze publicly available data deposited in the

367 NCBI SRA or EMBL ENA. All project datasets and results discussed in this manuscript are

368 provided on this site.

369

17

370     One of the key features of the EDGE Bioinformatics platform is that visualization of the results is

371     fully integrated with, and accessible directly on, the webpage in real time.  Many graphics are

372     displayed on each project page as thumbnails that link to either a full-page view or a lightbox

373     (quick zoom) view, including quality control graphics, assembly summary charts, heat maps,

374     phylogenetic trees, etc. In addition, there are links to the interactive genome browser JBrowse

375     and to interactive classification results via Krona, as well as links to output directories where all

376     resulting data for each pipeline are stored.

377

378     Because some of the most challenging aspects of genomics involve the exponentially

379     increasing size of datasets and the resources required to move large datasets, a key benefit of

380     the EDGE Bioinformatics software is that it can be implemented on a stand-alone server that

381     can access datasets in local storage or in network-mounted space. We have tested EDGE

382     Bioinformatics with datasets of up to hundreds of millions of reads, on a variety of servers (e.g.

383     12-64 core servers with 64GB-512GB of RAM), with run times ranging from minutes to hours.

384     Using more CPUs will decrease runtime (see Table 1). All projects run for this manuscript were

385     performed on the publicly available server (https://bioedge.lanl.gov/), which is a Dell PowerEdge

386     R720 with 24 cores, 512GB RAM, and 7 TB disk space. On this particular server, we have

387     restricted use to a maximum of 20 CPUs that can be specified by any given analysis.

388

389     A user management system has been implemented to provide a level of privacy/security for a

390     user's submitted projects. When this system is activated, any user can view projects that have

391     been made public, but other projects can only be accessed by logging into the system using a

392     registered local EDGE account or via an existing social media account (Facebook, Google+,

393     Windows, or LinkedIn). The users can then run new jobs and view their own previously run

394     projects or those that have been shared with them.

395

**The project page layout**

A left navigation menu on the EDGE website provides access to the Home page, the Run EDGE page (to initiate a new project) and the Projects list, allowing users to navigate to any desired project page (Supplementary Figure S3). A page for each project is produced as soon as it is launched within EDGE and allows the user to monitor the progress of the run and access the output summaries of each pipeline as they complete in real time. Each project page provides a summary of the project, and under a 'General' tab, a description of the input(s) provided, the modules selected for the run along with their run time statistics, and access to log files, the output directory, and a final PDF report.

A link in the upper right corner provides access to a sliding panel that contains a job progress widget, a resource monitoring widget, and an action widget. Once the job is submitted, the job progress widget reports the status for each analysis step in real time. The resource monitoring widget provides a real time view of the computational system running EDGE, and allows the user to anticipate whether there are sufficient resources to simultaneously run additional jobs, or if some projects should be moved to a different storage location. For example, projects will fail to complete one or more of the modules if there is insufficient storage for the outputs. The action widget provides the user some flexibility over the project, including allowing a user to interrupt, rerun, delete, and move his or her submitted jobs. The user can also share the project with other users, publish the project such that any user can access the results, or make the project private again (unpublish). In addition, there is a command line 'live log' view, which displays the real time actions and the Unix commands launched by EDGE.

**The EDGE modules and their outputs**

All of the six main modules within the EDGE Bioinformatics environment are optional and can be selectively run as individual modules or in any combination, thus affording the user maximum

422    flexibility in customizing each analysis to particular specifications. These consist of: 1) a pre-

423    processing module that performs quality control, trimming, and removal of sequences matching

424    an unwanted target (e.g. host removal); 2) a *de novo* assembly module which assembles the

425    data, validates the assembly, and annotates the resulting contigs; 3) a reference-based analysis

426    module, which allows users to select one or more references to which reads (and contigs) are

427    compared; 4) a taxonomy classification module, which classifies reads (and contigs); 5) a

428    phylogenetics module, which calculates a core genome, determines all SNPs, and infers a

429    phylogenetic tree from a number of input genomes; and 6) a primer and assay module which

430    allows users to validate *in silico* known primers against the *de novo* assembly, or to design new

431    primers that uniquely amplify short sequences within the *de novo* assembly. The latter module

432    does require an assembly for primer analysis.

433

434    Each module comprises a Perl wrapper with one or more bioinformatics tools tailored to handle

435    NGS reads and/or contigs, as well as several scripts to parse and post-process the results. The

436    users can also adjust a limited set of parameters or toggle options within each module. EDGE

437    produces a web page for each project with many different summaries of the results for each

438    module, including the statistics of the run (each module and time to completion), summary log

439    files and a PDF summary of all results, along with more detailed results of each individual

440    module. Each module outputs a number of files, which are accessible via a directory link and

441    are summarized with both text and figures along with some interactive graphics all within the

442    context of the website.

443

444    **Pre-Processing (Supplementary Figure S1, module 1).** This module consists of two

445    independent, selectable pipelines. For data quality control, the FaQCs software is used to

446    analyze all reads for quality and to trim or filter out reads using default parameters, unless these

447    are changed by the user (optional). Using an input reference FASTA, EDGE can also filter

448    unwanted reads that align to a selected reference. While this 'Host Removal' function was

449    originally envisioned to exclude host reads when inputting clinical samples or those derived from

450    known animals, this component can remove any data that aligns to the input reference, allowing

451    users to selectively remove any other target genome(s). Some built-in references include the

452    most recently updated GRCh38 Human reference and PhiX, which is often used as a control

453    within Illumina runs. This module aims to provide high-quality, clean reads for any subsequent

454    analysis by EDGE. If this module is not selected, the raw data will be used for all downstream

455    process modules.

456

457    Statistics and graphical outputs of the data, prior to and after processing, are provided for user

458    interpretation, along with access to the cleaned data files. The major outputs of this module are

459    shown in Supplementary Figure S1 (A, B and C) and example screen shots of output from the

460    EDGE webpage can be found in Supplementary Figure S4.

461

462    **Assembly and Annotation (Supplementary Figure S1, module 2).** EDGE performs *de novo*

463    assembly with the input reads using either IDBA-UD or SPAdes.  Because each of these

464    assemblers performs and combines multiple assemblies, both tools are capable of providing

465    reasonable assemblies from a wide variety of sample types, including isolate genomes, single

466    cell projects, and metagenomes. IDBA-UD is used by default (due to time and memory

467    considerations), and the assembly parameter option for kmer sizes begins with k=31 with a step

468    size of 20, until a maximum kmer size is reached (dependent on the read lengths). When this

469    module is selected, assembly validation is performed by mapping the short read input data to

470    the assembled contigs using Bowtie2. Additionally, the user can select to have the assembly

471    annotated (default behavior) using a modified Prokka tool (for the rapid annotation of prokaryotic

472    genomes), and prophages within microbial genomes are detected using Phage_Finder. If there

473    is an available reference that is sufficiently similar to the target genome assembly, EDGE can

21

474    also use a modified version of the Rapid Annotation Transfer Tool (RATT) to transfer the

475    annotation from the reference GenBank file (a required input for this step) to the assembly.

476    When SPAdes is selected as the assembler, there exists an additional option to input long read

477    data (PacBio or Nanopore) which can help in gap closure and repeat resolution.

478

479    The results of this module include the assembled contigs FASTA file, assembly and assembly

480    validation statistics and graphics, the annotation files (gbk and gff), and an interactive JBrowse

481    implementation, which provides visualization of the contigs and their annotation. The major

482    outputs of this module are displayed in Supplementary Figure S1 (D, E, F and G) and an

483    example screenshot can be found in Supplementary Figure S5.

484

485    **Reference-based Analysis (Supplementary Figure S1, module 3).** When this module is

486    selected, the user must choose one or more reference genomes (FASTA or Genbank formats)

487    to which the reads (and contigs, if assembly was performed) are compared. Reads are aligned

488    to the input reference using BowTie2 and variants are identified using SAMtools. Any regions

489    left uncovered by reads are also identified and reported in text files. Similarly, contigs are

490    aligned to the same reference(s) using MUMmer and the results parsed using Perl scripts to

491    catalogue SNPs and small insertions or deletions (indels), as well as regions within the contigs

492    that may be novel and do not align to the reference. If Genbank reference files are provided, the

493    variants, SNPs, and uncovered regions of the reference are further analyzed to output any

494    affected genes and reports are generated to display whether the changes also contribute to

495    synonymous or non-synonymous substitutions within coding regions. Reads and contigs that do

496    not map to the reference are parsed into separate FASTA/Q files and an option is available to

497    align these reads and contigs to RefSeq for taxonomic identification.

498

499    In addition to the output text files, several graphics along with statistics are provided that outline

500    linear coverage of the reference, depth of coverage along the reference, number of variants, as

501    well as percentages of input reads and contigs mapped to the reference. Interactive JBrowse

502    views allow for the display of the reference and associated annotation (genes, rRNAs, etc.),

503    along with detailed views of the aligned reads and contigs, as well as any SNPs or small indels

504    that have been discovered. The major outputs of this module are displayed in Supplementary

505    Figure S1 (G, H, and I), while example outputs can be found in Figure 4 and Supplementary

506    Figure S6.

507

508    **Taxonomy Classification (Supplementary Figure S1, module 4).** Envisioned primarily for

509    use with metagenomic datasets or with novel genomes, this module allows both read-based and

510    contig-based classification (the latter performed if assembly was also selected). For taxonomic

511    classification of the reads, the user can select one or more of several available metagenome

512    tools (currently GOTTCHA, Kraken, and MetaPhlAn) along with BWA, a read mapper used

513    against RefSeq. The default is to run all tools to take advantage of their different strengths, and

514    to provide users with additional information to help interpret their data. Each of these classifiers

515    has its own algorithm and database, parameters for the search, and required input format, all of

516    which are automatically managed within the EDGE platform. The specific output formats of each

517    tool are unified into a common framework to generate the reports/graphs displayed by EDGE.

518    There is also an option to classify only unassembled reads, if assembly is selected and the user

519    desires to only classify unassembled data.

520

521    The results of each read-based taxonomy profiling method are summarized in comparative

522    views (heatmap plots and radar charts summarize the top hits of each tool) at the user-selected

523    level of taxonomy (genus, species, strain). Results are also presented in more detail in

524     individual tool-based views with taxonomy tree dendrograms and Krona charts (e.g. Figure 5C)

525     while more detailed outputs can be found within the directory links.

526

527     For contig classification, EDGE aligns contigs to NCBI's RefSeq database using BWA-mem.

528     While contigs can match multiple taxa, each segment within a contig is assigned to a unique

529     taxon based on best hit score. While the total length within all contigs is calculated per taxon,

530     each contig is also assigned to a unique taxon based on linear coverage. Both the total length

531     per taxon (Length barplot) and the number of contigs (Count barplot) assigned to a taxon are

532     reported, along with a scatterplot showing the identity of the contig, its fold coverage by reads,

533     and its G+C content (Supplementary Figure S6). These results are reported at all levels of

534     taxonomy using the last common ancestor algorithm.

535

536     The major outputs of this module are displayed in Supplementary Figure S1 (J and K), while

537     example outputs can be found in Figures 2, 3, 5 and Supplementary Figures S9 and S10.

538

539     **Phylogenetic Analysis (Supplementary Figure S1, module 5).**  Because phylogenetic

540     analysis is a highly desired feature for many genomic investigations, we utilize a portion of a

541     newly developed tool, PhaME, which provides the ability to infer a whole genome SNP-based

542     tree from completed genomes, genome assemblies, and even from reads. Briefly, contigs and

543     selected genomes are compared with one another to identify conserved segments while

544     ignoring repeated regions, and reads are mapped to one of these genomes to continue the

545     identification of a conserved core genome. The core genome alignment is used to identify all

546     SNPs and FastTree (default, for speed considerations) or RAxML can be used to generate a

547     phylogenetic tree. This module was envisioned for use primarily with isolate genome projects

548     (however metagenomes have also been successfully used), where a target genome comprises

549     the majority of the sequencing data and the user desires to accurately place this target genome

24

550    within the context of near neighbor genomes. The user is required to select a minimum of two

551    reference genomes to which the reads and contigs (if available) will be added to infer a

552    phylogeny.

553

554    The Newick format tree files, core genome FASTA, and SNP statistics are available in the

555    directory link and the phylogenetic trees, generated using jsPhyloSVG, are provided for easy

556    viewing in either rectangular or circular tree formats (Outputs L and M in Supplementary Figure

557    S1). The input sample (reads and/or contigs) is highlighted within the trees. Output examples

558    can be found in Figures 2 and 5.

559

560    **PCR Primer Analysis (Supplementary Figure S1, module 6).** EDGE also supports both the

561    design and validation of PCR primers based on the assembly. In the validation pipeline, known

562    primers within a user-specified input file are mapped to the assembly using BWA, given a user-

563    defined number of mismatches (default of 1) to determine if an amplicon would be generated.

564    The user can also select a pipeline to design new primers based on the assembly, that will

565    differentiate the input sequenced sample from all other bacteria and viruses in NCBI's RefSeq

566    database. In this design component, unique regions are identified using BWA, and Primer3 is

567    used to select primer pairs. All primers are further filtered by melting temperature ($T_m$) difference

568    to the nearest neighbor background, within a user-specified value (5°C by default).

569

570    For primer validation, the primer binding location(s) and product sizes are reported for any

571    submitted primers (output N in Supplementary Figure S1). For primer design, a full list of

572    primers that uniquely amplify a product within the assembled contigs is reported (only five are

573    displayed by default on the project page), along with information on the nearest neighbor

574    amplicon (output O in Supplementary Figure S1). Examples of output for both primer validation

575    and primer design can be found in Supplementary Figure S7.

576

**DATA ACCESS**

All data have been deposited to NCBI and accession numbers are shown below.

*Bacillus anthracis* strain SK-102, SRR1993644

*Yersinia pestis* strain Harbin 35, SRR1993645

Patient fecal sample, 2011 *E. coli* outbreak, SRR2164314

Patient nasal swab, acute respiratory illness, SRP062772

583

**ACKNOWLEDGEMENTS AND DISCLAIMERS**

600

601    **DISCLOSURE DECLARATION**

602    There are no disclosures to declare.

603

604    **FIGURE LEGENDS**

605    **Figure 1. An overview of the EDGE Bioinformatics Environment.** The only Inputs required from the

606    user are raw sequencing data and a project name. The user can create specific workflows with any

607    combination of the modules. In addition, tailored parameters dictating how each module functions can be

608    modified by the user. EDGE outputs a variety of files, tables and graphics which can be viewed on screen

609    or downloaded. A more detailed overview is shown in Supplementary Figure S1. All Modules are

610    described in the Methods section.

611

612    **Figure 2. Taxonomy and phylogenetic evaluations of bacterial isolates.** Panels A and B show

613    taxonomic classification of reads for A) the *Y. pestis* Harbin35 sample and B) the *B. anthracis* SK-102

614    sample. The stars indicate the consistent dominant taxonomic calls for all tools, while the black arrow and

615    bracket indicate identified contamination in the *B. anthracis* sample. Panels C and D indicate the inferred

616    phylogenetic trees for the C) *Y. pestis* and D) *B. anthracis*; black arrows point to the read dataset (pink)

617    and contigs (blue) that were placed in these trees.

618

619    **Figure 3. Taxonomic Classification of the HMP staggered mock sample.** A) Read-based

620    classification using various taxonomy profiling tools; B) contig-based classification displaying length of all

621    classified contigs per taxon; and C) a scatterplot of contig % GC vs. fold coverage of the contigs, colored

622    by taxon.

623

624    **Figure 4. Interactive genome browsing view of a reference-based analysis in EDGE with a human**

625    **clinical sample containing Ebolavirus.** A) An Ebola reference genome and its genes (green lines) are

626    displayed together with contig-based (using IDBA) and read-based comparisons. The two contigs (blue

627    lines) from IDBA are shown aligned along the length of the reference as well as the reads (red and blue).

628    B) A zoomed-in view of one section of the genome where SNPs were identified. The SNP and coding

629   difference is outlined under the contig alignment, while the variants are indicated under the read

630   alignments.

631

632   **Figure 5. Phylogenetic and taxonomic analysis of human clinical samples with suspected and**

633   **unknown causative agents.** A) Circular phylogenetic tree clearly places within the *E. coli* O104 group

634   both the raw reads and the contigs obtained from a clinical fecal sample. B) A comparative heatmap view

635   of identified taxa from a nasal swab sample demonstrates the abundance of typical nasal cavity

636   organisms. C) The *E. coli* identified with GOTTCHA in the nasal swab sample (in B) is described in

637   greater detail under the tool-specific EDGE view (red arrow), showing the percent of hits to plasmids for

638   each identified taxon; below are a taxonomic dendrogram featuring the taxa detected with circles

639   representing relative abundance, and a Krona plot view of the same data.

640

641   **Table 1. Descriptions of samples and EDGE modules tested.**

642

643   **\*** EDGE Modules are described in Methods: 1. Pre-Processing; 2. Assembly and Annotation;

644   3. Reference-Based Analysis; 4. Taxonomic Classification; 5. Phylogenetic Analysis; 6. PCR Primer

645   Analysis

646   \*\* These samples were retrieved directly from the NCBI SRA.

647

648   **REFERENCES:**

649   2012. Protocol for Detection of Bacillus anthracis in Environmenal Samples During the

650         Remediation Phase of an Anthrax Event. United States Environmental Protection

651         Agency, National Homeland Security Research Center (NHSRC), Threat and

652         Consequence Assessment Division.

653   Ahmed SA, Lo C-C, Li P-E, Davenport KW, Chain PSG. 2015. From raw reads to trees: Whole

654         genome SNP phylogenetics across the tree of life. *bioRxiv*.

655    Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool.

656        *Journal of molecular biology* **215**(3): 403-410.

657    Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI,

658        Pham S, Prjibelski AD et al. 2012. SPAdes: a new genome assembly algorithm and its

659        applications to single-cell sequencing. *Journal of computational biology : a journal of*

660        *computational molecular cell biology* **19**(5): 455-477.

661    Bassis CM, Erb-Downward JR, Dickson RP, Freeman CM, Schmidt TM, Young VB, Beck JM,

662        Curtis JL, Huffnagle GB. 2015. Analysis of the upper respiratory tract microbiotas as the

663        source of the lung and gastric microbiotas in healthy individuals. *mBio* **6**(2): e00037.

664    Begier EM, Asiki G, Anywaine Z, Yockey B, Schriefer ME, Aleti P, Ogden-Odoi A, Staples JE,

665        Sexton C, Bearden SW et al. 2006. Pneumonic plague cluster, Uganda, 2004. *Emerging*

666        *infectious diseases* **12**(3): 460-467.

667    Bennett S. 2004. Solexa Ltd. *Pharmacogenomics* **5**(4): 433-438.

668    Blankenberg D, Von Kuster G, Coraor N, Ananda G, Lazarus R, Mangan M, Nekrutenko A,

669        Taylor J. 2010. Galaxy: a web-based genome analysis tool for experimentalists. *Current*

670        *protocols in molecular biology / edited by Frederick M Ausubel  [et al]* **Chapter 19**: Unit

671        19 10 11-21.

672    Buermans HP, den Dunnen JT. 2014. Next generation sequencing technology: Advances and

673        applications. *Biochimica et biophysica acta* **1842**(10): 1932-1941.

674    Conlan S, Thomas PJ, Deming C, Park M, Lau AF, Dekker JP, Snitkin ES, Clark TA, Luong K,

675        Song Y et al. 2014. Single-molecule sequencing to track plasmid diversity of hospital-

676        associated carbapenemase-producing Enterobacteriaceae. *Science translational*

677        *medicine* **6**(254): 254ra126.

678    Daber R, Sukhadia S, Morrissette JJ. 2013. Understanding the limitations of next generation

679        sequencing informatics, an approach to clinical pipeline validation using artificial data

680        sets. *Cancer genetics* **206**(12): 441-448.

681    den Bakker HC, Allard MW, Bopp D, Brown EW, Fontana J, Iqbal Z, Kinney A, Limberger R,

682        Musser KA, Shudt M et al. 2014. Rapid whole-genome sequencing for surveillance of

683        Salmonella enterica serovar enteritidis. *Emerging infectious diseases* **20**(8): 1306-1314.

684    Fasanella A, Losito S, Adone R, Ciuchini F, Trotta T, Altamura SA, Chiocco D, Ippolito G. 2003.

685        PCR assay to detect Bacillus anthracis spores in heat-treated specimens. *Journal of*

686        *clinical microbiology* **41**(2): 896-899.

687    Fouts DE. 2006. Phage_Finder: automated identification and classification of prophage regions

688        in complete bacterial genome sequences. *Nucleic acids research* **34**(20): 5839-5851.

689    Francy DS, Bushon RN, Grady AMG, Bertke EE, Kephart CM, Likirdopulos CA, Mailot BE, III

690        FWS, Lindquist HDA. 2009. Performance of Traditional and Molecular Methods for

691        Detecting Biological Agents in Drinking Water. In *In: US Department of the Interior, US*

692        *Geological Survey*.

693    Freitas TA, Li PE, Scholz MB, Chain PS. 2015. Accurate read-based metagenome

694        characterization using a hierarchical suite of unique signatures. *Nucleic acids research*.

695    Hinnebusch J, Schwan TG. 1993. New method for plague surveillance using polymerase chain

696        reaction to detect Yersinia pestis in fleas. *Journal of clinical microbiology* **31**(6): 1511-

697        1514.

698    Human Microbiome Project C. 2012. A framework for human microbiome research. *Nature*

699        **486**(7402): 215-221.

700    Johnson SL, Daligault HE, Davenport KW, Jaissle J, Frey KG, Ladner JT, Broomall SM, Bishop-

701        Lilly KA, Bruce DC, Coyne SR et al. 2015a. Thirty-Two Complete Genome Assemblies

702        of Nine Yersinia Species, Including Y. pestis, Y. pseudotuberculosis, and Y.

703        enterocolitica. *Genome announcements* **3**(2).

704    Johnson SL, Daligault HE, Davenport KW, Jaissle J, Frey KG, Ladner JT, Broomall SM, Bishop-

705        Lilly KA, Bruce DC, Gibbons HS et al. 2015b. Complete genome sequences for 35

706        biothreat assay-relevant bacillus species. *Genome announcements* **3**(2).

707     Kircher M, Sawyer S, Meyer M. 2012. Double indexing overcomes inaccuracies in multiplex

708         sequencing on the Illumina platform. *Nucleic acids research* **40**(1): e3.

709     Koren S, Treangen TJ, Hill CM, Pop M, Phillippy AM. 2014. Automated ensemble assembly and

710         validation of microbial genomes. *BMC bioinformatics* **15**: 126.

711     Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg SL. 2004.

712         Versatile and open software for comparing large genomes. *Genome biology* **5**(2): R12.

713     Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature methods*

714         **9**(4): 357-359.

715     Li H, Durbin R. 2010. Fast and accurate long-read alignment with Burrows-Wheeler transform.

716         *Bioinformatics* **26**(5): 589-595.

717     Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,

718         Genome Project Data Processing S. 2009. The Sequence Alignment/Map format and

719         SAMtools. *Bioinformatics* **25**(16): 2078-2079.

720     Lo CC, Chain PS. 2014. Rapid evaluation and quality control of next generation sequencing

721         data with FaQCs. *BMC bioinformatics* **15**(1): 366.

722     Margulies M, Egholm M, Altman WE, Attiya S, Bader JS, Bemben LA, Berka J, Braverman MS,

723         Chen YJ, Chen Z et al. 2005. Genome sequencing in microfabricated high-density

724         picolitre reactors. *Nature* **437**(7057): 376-380.

725     Maxam AM, Gilbert W. 1977. A new method for sequencing DNA. *Proceedings of the National*

726         *Academy of Sciences of the United States of America* **74**(2): 560-564.

727     Ondov BD, Bergman NH, Phillippy AM. 2011. Interactive metagenomic visualization in a Web

728         browser. *BMC bioinformatics* **12**: 385.

729     Otto TD, Dillon GP, Degrave WS, Berriman M. 2011. RATT: Rapid Annotation Transfer Tool.

730         *Nucleic acids research* **39**(9): e57.

731   Peng Y, Leung HC, Yiu SM, Chin FY. 2012. IDBA-UD: a de novo assembler for single-cell and

732       metagenomic sequencing data with highly uneven depth. *Bioinformatics* **28**(11): 1420-

733       1428.

734   Price MN, Dehal PS, Arkin AP. 2010. FastTree 2--approximately maximum-likelihood trees for

735       large alignments. *PloS one* **5**(3): e9490.

736   Rawlings BA, Higgins TS, Han JK. 2013. Bacterial pathogens in the nasopharynx, nasal cavity,

737       and osteomeatal complex during wellness and viral infection. *American journal of*

738       *rhinology & allergy* **27**(1): 39-42.

739   Roggenbuck M, Baerholm Schnell I, Blom N, Baelum J, Bertelsen MF, Ponten TS, Sorensen

740       SJ, Gilbert MT, Graves GR, Hansen LH. 2014. The microbiome of New World vultures.

741       *Nat Commun* **5**: 5498.

742   Sanger F, Nicklen S, Coulson AR. 1977. DNA sequencing with chain-terminating inhibitors.

743       *Proceedings of the National Academy of Sciences of the United States of America*

744       **74**(12): 5463-5467.

745   Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**(14): 2068-

746       2069.

747   Segata N, Waldron L, Ballarini A, Narasimhan V, Jousson O, Huttenhower C. 2012.

748       Metagenomic microbial community profiling using unique clade-specific marker genes.

749       *Nature methods* **9**(8): 811-814.

750   Skinner ME, Uzilov AV, Stein LD, Mungall CJ, Holmes IH. 2009. JBrowse: a next-generation

751       genome browser. *Genome research* **19**(9): 1630-1638.

752   Smits SA, Ouverney CC. 2010. jsPhyloSVG: a javascript library for visualizing interactive and

753       vector-based phylogenetic trees on the web. *PloS one* **5**(8): e12267.

754   Soufiane B, Cote JC. 2013. Bacillus weihenstephanensis characteristics are present in Bacillus

755       cereus and Bacillus mycoides strains. *FEMS microbiology letters* **341**(2): 127-137.

756    Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of

757         large phylogenies. *Bioinformatics* **30**(9): 1312-1313.

758    Swanson MM, Reavy B, Makarova KS, Cock PJ, Hopkins DW, Torrance L, Koonin EV,

759         Taliansky M. 2012. Novel bacteriophages containing a genome of another

760         bacteriophage within their genomes. *PloS one* **7**(7): e40683.

761    Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, Rozen SG. 2012.

762         Primer3--new capabilities and interfaces. *Nucleic acids research* **40**(15): e115.

763    Wang J, Chen L, Huang S, Liu J, Ren X, Tian X, Qiao J, Zhang W. 2012. RNA-seq based

764         identification and mutant validation of gene targets related to ethanol resistance in

765         cyanobacterial Synechocystis sp. PCC 6803. *Biotechnology for biofuels* **5**(1): 89.

766    Watson-Haigh NS, Shang CA, Haimel M, Kostadima M, Loos R, Deshpande N, Duesing K, Li X,

767         McGrath A, McWilliam S et al. 2013. Next-generation sequencing: a challenge to meet

768         the increasing demand for training workshops in Australia. *Briefings in bioinformatics*

769         **14**(5): 563-574.

770    Wohlbach DJ, Rovinskiy N, Lewis JA, Sardi M, Schackwitz WS, Martin JA, Deshpande S, Daum

771         CG, Lipzen A, Sato TK et al. 2014. Comparative genomics of Saccharomyces cerevisiae

772         natural isolates for bioenergy production. *Genome biology and evolution* **6**(9): 2557-

773         2566.

774    Wood DE, Salzberg SL. 2014. Kraken: ultrafast metagenomic sequence classification using

775         exact alignments. *Genome biology* **15**(3): R46.

776