

1 **Methods for cost-efficient, whole genome**
2 **sequencing surveillance for enhanced detection of**
3 **outbreaks in a hospital setting**

4
5 **1.1 Author names**

6 *Kady D. Waggle^{1,2,4}, Marissa Pacey Griffith^{1,2}, Alecia B. Rokes^{5,6}, Vatsala Rangachar*
7 *Srinivasa^{1,2,3}, Deena Ereifej^{1,2,3}, Rose Patrick^{1,2}, Hunter Coyle^{1,2}, Shurmin Chaudhary^{1,2},*
8 *Nathan J. Raabe^{1,2,3}, Alexander J. Sundermann^{1,2}, Vaughn S. Cooper^{5,6}, Lee H. Harrison^{1,2,3},*
9 *Lora Lee Pless^{1,2#}*

10

11 **ORCID iDs**

12 Kady D. Waggle  <https://orcid.org/0009-0001-9237-9964>
13 Marissa Pacey Griffith  <https://orcid.org/0000-0003-0677-8203>
14 Alecia B. Rokes  <https://orcid.org/0000-0001-9496-0296>
15 Vatsala Rangachar Srinivasa  <https://orcid.org/0000-0001-7376-8764>
16 Deena Ereifej  <https://orcid.org/0009-0002-2984-3784>
17 Rose Patrick  <https://orcid.org/0009-0003-3541-5513>
18 Hunter Coyle  <https://orcid.org/0009-0006-0463-4168>
19 Shurmin Chaudhary  <https://orcid.org/0009-0009-1640-7132>
20 Nathan J. Raabe  <https://orcid.org/0000-0002-3258-9003>
21 Alexander J. Sundermann  <https://orcid.org/0000-0002-2439-6287>
22 Vaughn S. Cooper  <https://orcid.org/0000-0001-7726-0765>
23 Lee H. Harrison  <https://orcid.org/0000-0002-3787-2705>

24 Lora Lee Pless  <https://orcid.org/0000-0003-0163-3689>

25

26

27 **1.2 Affiliation(s)**

28 ¹Microbial Genomic Epidemiology Laboratory, Center for Genomic Epidemiology,
29 University of Pittsburgh, 3507 Victoria Street, BST-10 E1000-4A, Pittsburgh,
30 Pennsylvania 15213, USA.

31 ²Division of Infectious Diseases, University of Pittsburgh School of Medicine, 3550 Terrace
32 Street, 818 Scaife Hall, Pittsburgh, Pennsylvania 15261, USA.

33 ³Department of Epidemiology, School of Public Health, University of Pittsburgh, 130 De
34 Soto Street, Pittsburgh, Pennsylvania 15261, USA

35 ⁴Department of Infectious Diseases, School of Public Health, University of Pittsburgh, 130
36 De Soto Street, Pittsburgh, Pennsylvania 15261, USA

37 ⁵Department of Microbiology and Molecular Genetics, University of Pittsburgh School of
38 Medicine, Pittsburgh, PA, USA

39 ⁶Center for Evolutionary Biology and Medicine, University of Pittsburgh School of Medicine,
40 Pittsburgh, PA, USA.

41

42 **1.3 Corresponding author and email address**

43 #Lora Lee Pless

44 lora.pless@pitt.edu

45 University of Pittsburgh

46 Starzl Biomedical Science Tower

47 200 Lothrop St

48 Pittsburgh, PA 15213

49 **1.4 Keywords**

50 Healthcare-Associated Transmission; Bacteria, Antimicrobial Resistance, Whole Genome

51 Sequencing; Laboratory Methods; Cost Analysis

52

53 2. Abstract

54 **Introduction.** Outbreaks of healthcare-associated infections (HAI) result in substantial
55 patient morbidity and mortality; mitigation efforts by infection prevention teams have the
56 potential to curb outbreaks and prevent transmission to additional patients. The incorporation
57 of whole genome sequencing (WGS) surveillance of suspected high-risk pathogens often
58 identifies outbreaks that are not detected by traditional infection prevention methods and
59 provides evidence for transmission. Our approach to real-time WGS surveillance, the
60 Enhanced Detection System for Healthcare-Associated Transmission (EDS-HAT), has 1)
61 identified serious outbreaks that were otherwise undetected and 2) shown the potential to be
62 cost saving because HAIs are expensive to treat and WGS has become relatively
63 inexpensive.

64 **Methods.** We describe a cost-efficient method to perform WGS surveillance and data
65 analysis of pathogens for hospitals that are interested in incorporating WGS surveillance.
66 We provide an overview of the weekly workflow of EDS-HAT, discussing both the laboratory
67 and bioinformatics methods utilized, as well as the costs associated with performing these
68 methods.

69 **Results.** In an average week at our tertiary healthcare system, we sequenced 48 samples at
70 a cost of less than \$100 per sample, inclusive of laboratory reagents and staff salaries. The
71 average turnaround time, from sample collection to data reporting to the infection prevention
72 and control team, was ten days.

73 **Conclusions.** Our findings demonstrate that performing EDS-HAT in real-time can be both
74 affordable and time-efficient. Providing such timely information to aid in outbreak
75 investigations can identify transmission events sooner and thus increase patient safety.

76 **3. Impact statement**

77 Whole genome sequencing (WGS) surveillance to confirm or refute suspected outbreaks of
78 potential healthcare-associated infections (HAI) is a highly effective approach for outbreak
79 detection. Since November 2021, we have been conducting WGS surveillance in real-time
80 through a program called the Enhanced Detection System for Hospital-Associated
81 Transmission (EDS-HAT), to assist our hospital infection prevention and control (IP&C) team
82 to identify and stop outbreaks. To our knowledge, our laboratory is the only group in the
83 United States that has successfully implemented real-time WGS surveillance of multiple
84 pathogens in the hospital setting. Our weekly workflow includes identifying HAI pathogens
85 and performing WGS, followed by a variety of bioinformatic analyses that include species
86 confirmation, determination of sequence type, and genetic relatedness comparisons. Based
87 on this information, transmission clusters are identified, and the electronic health record is
88 reviewed to determine probable transmission routes. Finally, IP&C implements appropriate
89 interventions to mitigate the spread of infection. We detail the laboratory and analytical
90 methods, along with the cost associated for laboratory materials and staff salary, for
91 successful implementation of WGS surveillance in real-time establishing EDS-HAT as a
92 unique and effective tool to detect HAI outbreaks.

93 **4. Introduction**

94 Healthcare-associated infections (HAIs) are a growing concern in hospital settings and can
95 be associated with substantial morbidity and mortality. HAIs also impose a significant
96 economic burden on healthcare systems, costing hospitals an estimated USD\$ 9.6 billion
97 per year (1,2). Whole genome sequencing (WGS) for HAI organisms can provide insight on
98 the transmission dynamics in hospital settings (3). Historically, determining the degree of
99 genomic variation between organisms was accomplished using pulsed field gel

100 electrophoresis (PFGE; 4). Given recent declines in costs and its many advantages, WGS
101 has emerged as the leading method for determining genetic relatedness between clinical
102 isolates (5). Reactive WGS is currently the most commonly used method to confirm or refute
103 the presence of a suspected outbreak. This approach can result in a failure to detect
104 important outbreaks for a variety of reasons, including outbreaks caused by common
105 organisms, those not clustering on a single nursing unit, those consisting of a small number
106 of patients, or those caused by an unsuspected or complex transmission route (6). Reactive
107 WGS could also falsely identify an outbreak supported by epidemiological methods by
108 clustering genetically distinct isolates (7). Furthermore, using WGS to obtain information
109 about the entire genome provides the required data for determining organism phylogeny,
110 detecting the presence of antimicrobial resistance (AMR) genes and mobile genetic
111 elements, and identifying rare or novel genetic variants (7,8). The Microbial Genomic
112 Epidemiology Laboratory (MiGEL) at the University of Pittsburgh developed the Enhanced
113 Detection System for Healthcare-Associated Transmission (EDS-HAT) to identify outbreaks
114 of HAIs in real-time using WGS surveillance methods in partnership with the UPMC IP&C
115 team and the UPMC Clinical Laboratories. EDS-HAT has been operational in real-time at our
116 institution for over two years (7,9–12). The barriers for most hospital systems for
117 implementing proactive WGS are cost, lack of technical guidance, and inadequate
118 infrastructure.

119 In this paper, we describe our methods for WGS, the bioinformatics workflow, and
120 provide a cost estimate of WGS surveillance, with the goal of providing guidance to hospitals
121 who wish to implement WGS surveillance.

122 5. Theory and Implementation

123 Study Setting.

124 MiGEL is a non-Clinical Laboratory Improvement Amendments (CLIA) certified research
125 laboratory located on the University of Pittsburgh main campus, in Pittsburgh PA, USA.
126 EDS-HAT was developed and is currently implemented in real time at MiGEL in coordination
127 with the University of Pittsburgh, UPMC, the UPMC Clinical Laboratory Building (CLB), the
128 UPMC IP&C team, and Carnegie Mellon University (CMU). UPMC Presbyterian is an adult
129 tertiary acute care hospital with 758 total beds, 134 critical care beds, and over 400 annual
130 solid organ transplants. The main campus is UPMC Presbyterian Hospital, also located in
131 the Oakland neighborhood of Pittsburgh, PA, adjacent to the University of Pittsburgh main
132 campus and CMU. The University of Pittsburgh Institutional Review Board provided ethics
133 approval for EDS-HAT (Protocol: STUDY21040126).

134 Clinical Specimen Collection.

135 **Isolate Inclusion Criteria.** A list of select, high-concern bacterial pathogens was generated
136 twice per week using Theradoc (5.4.0.HF1.102, Pittsburgh, PA; Fig 1A). Pathogens of
137 interest include: extended-spectrum B-lactamase-producing (ESBL) *Escherichia coli*, ESBL
138 *Enterobacter* species, *Acinetobacter* species, *Pseudomonas* species, *Klebsiella* species,
139 *Stenotrophomonas* species, *Serratia* species, *Burkholderia* species, *Providencia* species,
140 *Proteus* species, *Citrobacter* species, vancomycin-resistant *Enterococcus* (VRE), methicillin-
141 resistant *Staphylococcus aureus* (MRSA), and *Clostridioides difficile*. EDS-HAT isolate
142 inclusion criteria included patients who have been in the hospital for three or more days
143 and/or had a previous hospital exposure during the 30-days prior to culture (7). For this
144 study, we described the samples and methods utilized during a one-year period of time of
145 performing real-time EDS-HAT (March 2022-March 2023).

146 **Isolate Collection.** Bacterial samples were collected by MiGEL twice per week at the UPMC
147 CLB from pure cultures isolated from clinical specimens that were prompted by clinician
148 suspicion of infection (Fig 1.B1). To ensure availability of the isolates for sequencing, CLB
149 technologists subcultured all gram-negative isolates from aerobic bacterial cultures onto
150 nutrient agar slants. We identified the gram-negative isolate slants of interest from the CLB,
151 and then isolates of interest were subcultured to Trypticase Soy Agar with 5% sheep blood
152 (BAP) plates (BD, Franklin Lakes, NJ), transported to MiGEL, and incubated at 37°C
153 overnight in the presence of 5% CO₂. The gram-positive isolates from the CLB were
154 transferred from one BAP to another and then transported and incubated at MiGEL following
155 the same procedure. The next day, sample information was imported into the MiGEL
156 database, and a de-identified specimen ID was generated for each sample.

157 ***Clostridioides difficile* collection and culture.** In contrast to the methods described for the
158 above organisms that were isolated as pure cultures, we collected and cultured clinical stool
159 specimens that tested positive for *C. difficile* by culture-independent diagnostic testing (13).
160 This organism is anaerobic; thus, we performed the following protocol to isolate this
161 organism directly from the clinical stool specimens. In a biosafety cabinet, each stool sample
162 was subcultured onto cycloserine-cefoxitin-mannitol-agar with taurocholate and lysozyme
163 (CCMA-TAL) plates to select for *C. difficile* growth. Plates were transferred into a Coy
164 anaerobic chamber (Coy Laboratory Products, Grass Lake, MI) and incubated at 37°C for 48
165 hours. Colonies of *C. difficile* were passaged to a second CCMA plate and incubated at
166 37°C in the anaerobic chamber for an additional 24-48 hours. Isolates were confirmed as *C.*
167 *difficile* by testing for the production of L-Proline aminopeptidase using a PRO Disc test
168 (Remel, San Diego, CA; Fig1.B2).

169 **Sample Preparation and DNA Extraction**

170 To begin sample preparation for WGS, microcentrifuge tubes containing 750 μ L phosphate
171 buffered saline (PBS) were inoculated with a quarter-portion of a 10 μ L loop of bacteria (a
172 half-portion was used for *C. difficile*) from the BAP or CCMA plate. The tubes were
173 centrifuged at $6.0 \times g$ for 10 minutes to generate a pellet, and the supernatant was removed
174 using a P1000 pipette (Fig 1C). For samples not proceeding immediately to extractions, the
175 pellets were stored at -20°C . Isolate stocks for long-term storage for all bacterial isolates
176 (including *C. difficile*) were prepared by inoculating a 10 μ L loop of bacteria into cryovials
177 containing 1 mL of nutrient broth mixed with 20% glycerol and then stored at -80°C .

178 The bacterial pellets were re-suspended in 500 μ L PBS prior to extraction. DNA was
179 extracted using the MagMAX DNA Multi-Sample Ultra 2.0 extraction kit on the King Fisher
180 Apex (Thermo Fisher Scientific, Waltham, MA) per manufacturer's instructions (Fig 1D).
181 Briefly, this procedure isolates and purifies nucleic acids using magnetic bead-based
182 technology. DNA was eluted in 100 μ L of elution buffer supplied by the kit and then
183 quantified using a Qubit broad range dsDNA kit (Life Technologies, Carlsbad, CA). Samples
184 with a concentration $\geq 3.5\text{ng}/\mu\text{L}$ were considered for WGS. For samples that did not meet this
185 criterion, DNA was extracted again.

186 **WGS Library Preparation**

187 DNA libraries were prepared on an epMotion 5075t (Eppendorf, Hamburg, Germany) liquid
188 handler using a DNA Prep (M) Tagmentation kit (Illumina, San Diego, CA), utilizing half-
189 volume reactions for BLT/TB1 and EPM reagents (Fig 1E). A unique 10-mer index adapter
190 sequence was ligated to each sample (IDT, Coralville, IA). Briefly, the DNA Prep protocol
191 uses bead-linked transposomes to tagment and amplify the adapter-tagged DNA segments.
192 Eight individual libraries were pooled together by combining 5 μ L per library into a single
193 tube. Pooled libraries were quantified using a Qubit high sensitivity dsDNA kit. The library

194 pool was normalized to 4 nM with resuspension buffer (RSB). Additional pools were
195 combined using equimolar concentration into a single pool. The distribution of the fragment
196 sizes for the sequencing pool was assessed using an Agilent Tapestation D5000 screen
197 tape and reagents per manufacturer's protocol (Agilent Technologies, Santa Clara, CA).

198 **Whole Genome Sequencing**

199 DNA libraries were sequenced weekly using an Illumina MiSeq (≤ 32 samples on a v3, 600
200 cycle kit) or NextSeq550 (> 32 samples on a v2.5, 300 cycle kit) platform (Fig 1F). The DNA
201 library was denatured using 0.2N NaOH and spiked with 1% PhiX to increase diversity on
202 the flow cell. The DNA library was diluted, using the average library length, to the final
203 loading concentration of 16 pM for the MiSeq or 1.5-1.6 pM for the NextSeq550. A
204 commercial lab was used for sequencing in rare cases where personnel were unavailable for
205 in-house sequencing. For these occasions, DNA was extracted and sent for same-day
206 delivery using a local medical courier service, followed by library preparation and sequencing
207 at the commercial lab. For sequencing using any of the options described, DNA extraction
208 and library preparation were performed using automated methods; however, it was possible
209 to perform all steps manually.

210 **Bioinformatics and Data Analysis**

211 **Sequencing Data Quality Control (QC).** We have developed a real-time bioinformatics
212 pipeline that is executed once per week as a single command written in the programming
213 language Python. This customized pipeline is one of four commands that are executed on
214 the new samples, as well as previously sequenced genomes. These commands include: 1)
215 data download from the BaseSpace Sequence Hub v7.18.0 (Illumina); 2) sample
216 demultiplexing; 3) file transfer into individual directories; and 4) real-time bioinformatics
217 pipeline execution. Specifically, we begin by converting and demultiplexing the base call files
218 using Illumina bcl2fastq (v2.20) software. WGS reads were assembled using Unicycler

219 v0.5.0 and then annotated using Prokka v1.14 (14). Multilocus sequence types (STs) were
220 assigned using PubMLST typing schemes for all organisms with the exception of *Serratia*
221 *spp.* and *Providencia spp.*, which do not have ST schemes (mlst v2.11; 15). Reads were
222 mapped using Kraken2 with the Kraken standard database to determine the most prevalent
223 species (16). Isolates passed QC if 1) the most prevalent species by Kraken2 was the
224 expected organism, 2) the assembly length was within 20% of the expected genome length,
225 3) the assembly was ≤ 350 contigs, and 4) there was at least 35 \times depth (Fig 1G).

226 **Determining Infection Clusters and Downstream Applications.** Pairwise single
227 nucleotide polymorphisms (SNPs) between all real-time EDS-HAT isolates of the same
228 species were determined using one of two programs (Fig 1H). i) Pairwise core genome
229 SNPs (cgSNPs) were determined using Snippy v4.3.0, a reference-based method, for
230 isolates with the same ST (17). SNP distances were calculated from the core alignment
231 using 'snp-dists' (18). ii) SKA v1.0, a reference-free method, was used to calculate SNP
232 distances using the 'ska distance' command for isolates of the same species (19). We
233 selected the minimum SNP distance for each pairwise comparison quantified by Snippy or
234 SKA to determine clusters of genetically similar isolates. These genetically similar clusters
235 were defined using hierarchical clustering with average linkage and a cutoff of ≤ 15 SNPs for
236 all species except *C. difficile*, for which a cutoff of ≤ 2 SNPs was used (Fig 1H;
237 <https://scipy.org/>). The electronic health records for patients with genetically similar isolates
238 were reviewed to determine potential epidemiological links. This information was then
239 communicated to the hospital IP&C team, which implemented targeted mitigation measures
240 when possible. See Supplementary Figure 1 for real-time bioinformatics pipeline.

241 **Cost Analysis**

242 A cost estimate for EDS-HAT real-time genomic surveillance methods was determined in
243 2023 US dollars and included the cost of personnel, reagents, and supplies, and was

244 analyzed comparatively for each sequencing platform used by MiGEL (Supplementary Table
245 1). Non-fringe personnel costs (salary) were determined using the average pay scale of
246 Laboratory Technician III (90% effort) and Bioinformatics Research Analyst II (50% effort)
247 positions at UPMC, Pittsburgh, PA in 2023. Reagent and supply costs were determined
248 using manufacturer pricing (data accessed: December 1, 2023).

249 **6. Results**

250 **Weekly Sequencing Runs**

251 From March 2022 to March 2023, MiGEL collected and sequenced 2,070 bacterial isolates
252 (with an average of 48 isolates per week) as part of real-time EDS-HAT. The most
253 commonly sequenced organism was *Pseudomonas aeruginosa* (617 genomes) and the
254 least sequenced was *Burkholderia sp.* (11 genomes; Table 1). To determine which platform
255 was best suited for weekly sequencing, we considered the count of organisms and the
256 average genome size. The weekly average genome size was 4.85 Mbp, roughly equating to
257 a maximum of 37 or 98 samples on the MiSeq or NextSeq flow cells, respectively, to achieve
258 a minimum target of 80x coverage. When sequencing pools of organisms with smaller
259 average genome sizes, a greater number of isolates could be appropriately accommodated
260 per flow cell without compromising run quality or per organism coverage data (Figure 2).
261 Based on the MiGEL average genome size and to maximize cost efficiency, runs containing
262 a range of 32-40 samples were sequenced on the MiSeq platform, and runs containing > 40
263 samples were sequenced on the NextSeq550 platform. During this study, 17 runs were
264 performed on the MiSeq platform, and 28 runs were performed on the NextSeq550 platform.
265 For an average run of 48 samples on the NextSeq550, MiGEL observed a maximum output
266 of 52 Gb of data and an average of 100 million reads (Supplementary Table 2). The average
267 turnaround time to complete the EDS-HAT workflow from sample collection by MiGEL to
268 bioinformatic analysis using either platform for sequencing was approximately 10 days, with

269 an average WGS instrument run time of 25 hours (Supplementary Table 2). The turnaround
270 time for using a commercial lab was approximately two weeks or less.

271 **Cost Analysis**

272 The cost to run real-time EDS-HAT weekly was categorized into sample processing, DNA
273 extraction and quantification, library preparation, and flow cell cost (Table 2). The lowest cost
274 per sample (\$48) was achieved when the maximum number of samples (n=96) were
275 sequenced using the NextSeq550 platform. Costs ranged from \$48 to \$83 per sample,
276 dependent on platform and sample counts. There was an inverse relationship between the
277 number of samples sequenced and flow cell cost as per sample costs significantly
278 decreased when a greater number of samples were multiplexed on the appropriate flow cell.
279 The gray dashed line in Figure 3 shows the cost to sequence 40 samples using all
280 sequencing options, with the MiSeq having the lowest cost and commercial lab having the
281 highest. The estimated weekly cost of personnel, based on the pre-tax salaries for one lab
282 technician and one bioinformatician based on percent efforts, totaled \$1,077. When all costs
283 were considered, the cost to run EDS-HAT on an average week totaled \$4,293 (min \$3,626
284 – max \$4,758) or \$223,236 per year (min \$188,552 – max \$247,416).

285 **7. Discussion**

286 In this study, we detailed an efficient laboratory workflow, our approach for bioinformatics
287 analyses, and estimated the cost associated with implementing real-time WGS surveillance
288 for pathogenic bacteria in a hospital system that was designed to detect otherwise
289 unrecognized hospital outbreaks. EDS-HAT began in 2016 as a retrospective study (7) and,
290 once we demonstrated the superiority of the system over traditional approaches, transitioned
291 in November 2021 to a real-time workflow, subsequent bioinformatic analyses, and reporting
292 of results to the hospital IP&C team. To our knowledge, UPMC is the only hospital system in

293 the US that is actively performing prospective WGS surveillance methods for multiple
294 pathogens in real time. By doing so, our hospital system has dramatically changed the way
295 outbreaks are being detected.

296 We provide details about our methods for a one-year timeframe, after our initial
297 optimization period, beginning in March 2022. We determined that the per sample cost for
298 WGS ranged from \$48 to \$83, with an average of \$65. Furthermore, with the addition of staff
299 salaries, the mean weekly cost for an average week of real-time sequencing was \$4,293
300 ($N=48$ samples). This cost of real-time WGS is lower compared to prior studies (20,21). Our
301 lower cost was achieved, in part, by increasing sample counts per flow cell while utilizing the
302 appropriate instrument, using half-volumes of reagents for some stages of library
303 preparation, and an overall decline in sequencing costs.

304 With our quick turnaround time from the day the sample is collected by MiGEL, we have
305 identified ongoing outbreaks that serve as a guide for the IP&C team to implement infection
306 prevention interventions. We previously showed that there was an estimated cost savings of
307 \$96,204–\$346,266 per year by implementing a real-time WGS surveillance system, which
308 was based on an average cost of \$86 for sample preparation and sequencing (adjusted for
309 inflation to 2023 USD; 22). We optimized the average per isolate cost of sample preparation
310 and sequencing from \$74 on the MiSeq platform (SD, \$3.30) to \$60 on the NextSeq platform
311 (SD, \$6.40), achieving even greater cost savings per year. 2/16/2024 1:00:00 PM
312 importantly, stopping transmission events quickly at the first sign of an outbreak cluster has
313 the potential to reduce further spread of the infection and thus reduce patient morbidity and
314 mortality.

315 The foremost concern of hospital systems with implementing programs like EDS-HAT is
316 cost, with the vast majority of interested parties assuming that there is a large expense
317 associated with real-time sequencing surveillance. While this was true years ago, the cost of

318 sequencing has decreased over time (23). In addition, the laboratory and bioinformatics
319 methods have become more streamlined, automatable, and efficient. Furthermore, the cost
320 of treating preventable hospital infections is high, and, in fact, EDS-HAT has been shown to
321 be cost saving. Taken together, these facts and the evidence that this approach can identify
322 important, otherwise-undetected outbreaks, suggest that WGS surveillance should
323 eventually become standard practice in hospitals.

324 To accompany our methods, we computed the cost per sample, which accounts for staff
325 salaries, to be \$91 on average (range \$62-\$119), and is specific for the greater Pittsburgh
326 region in Pennsylvania, USA and is likely to be different at other locations. This fact is
327 summarized by Price and colleagues, who find the cost to perform WGS varies by country
328 and city (20). For example, Price (20) converted the cost per sample from prior studies to
329 2023 USD and showed the per sample cost of sequencing ranged from approximately \$72-
330 \$470 for the US and Italy, respectively. In this study, we determined our average per sample
331 cost (without considering staff salaries, for comparison) was \$65 per sample. The primary
332 factor in determining this cost estimate was sample count per run and average organism
333 genome size. For reference, we provide the maximum number of samples that can be
334 sequenced on either MiSeq or NextSeq platforms by organism, considering genome size,
335 along with the average genome size sequenced over one year by MiGEL (Figure 2). In
336 addition, we show in Figure 3 that a sample count of 40 is an appropriate cutoff to decide
337 which machine to use for sample sequencing, while maintaining sufficient genome coverage.
338 Generally, we find sequencing more samples at a time reduced the cost of sequencing per
339 sample, with the exception of utilizing a commercial lab. While the commercial lab offered a
340 discounted price once the sample count reached 48, we find the fixed price was overall more
341 costly than performing in-house sequencing. Furthermore, we find a decrease in sequencing
342 costs over time. MiGEL estimated a \$72 average per sample cost in 2021, which we show is

343 lower in cost by \$7 during our study period (March 8, 2022 to March 9, 2023; average cost is
344 \$65).

345 We note limitations with this study. First, the costs for reagents and supplies presented
346 in this manuscript represent discounted pricing provided to our university from some
347 manufacturers. Other institutions may have different discounted pricing or pay manufacturers
348 rates, which will alter the costs described in our methods. Second, MiGEL benefits from the
349 use of robotic instruments for nucleic acid extractions and library preparation, which can help
350 save time and decrease pipetting errors on the bench. Some institutions may not have such
351 instruments available and will need to accommodate the laboratory methods we described
352 accordingly; however, we do not think this represents a significant detriment to the process.
353 Third, we only considered Illumina-based technology for this study. Other short-read
354 sequencing technologies or long-read sequencing were not assessed. Fourth, we have
355 demonstrated the cost-efficiency at an academic, tertiary hospital system. These estimates
356 are likely not reflective of a healthcare system located at a smaller locale.

357 In conclusion, we have shown that a real-time WGS surveillance program is both
358 feasible and affordable. Healthcare institutions wishing to do the same could potentially
359 discover outbreaks that would otherwise be missed. Further adoption of this approach has
360 the potential to significantly enhance patient safety.

361

362 8. Tables

363 Table 1.

364 EDS-HAT Bacterial genome sizes and number of each organism that has been sequenced
365 by MiGEL (March 8, 2022 to March 9, 2023).

Organism	Genome Size (Mb)	Count Sequenced by MiGEL
<i>Acinetobacter species</i>	3.9	58
<i>Burkholderia species</i>	7.2	11
<i>Clostridioides difficile</i>	4.2	67
<i>Citrobacter species</i>	4.7	46
<i>Enterobacter cloacae</i>	5.3	51
<i>Escherichia coli</i>	5.3	152
<i>Klebsiella oxytoca</i>	5.9	28
<i>Klebsiella pneumoniae</i>	5.3	115
MRSA	2.9	306
<i>Providencia species</i>	4.5	28
<i>Pseudomonas species</i>	6.6	643
<i>Serratia species</i>	5.2	154
<i>Stenotrophomonas species</i>	4.8	112
VRE	2.9	127
<i>Proteus species</i>	4.06	269

366

367 Table 2. Cost Estimates

Method	MiSeq cost per sample	NextSeq cost per sample	Commercial Sequencing Cost
--------	-----------------------	-------------------------	----------------------------

	(32 samples)	(48 samples)	(96 samples)	(any # samples)
Sample Processing	\$2.06	\$2.06	\$2.06	\$2.06
DNA Extraction & Quantification	\$7.15	\$7.15	\$7.15	\$7.15
Library Preparation	\$25.90	\$24.96	\$24.02	-
Sequencing Kit (Flow Cell)	\$49.16	\$37.29	\$18.65	\$75
TOTAL COST PER SAMPLE	\$84.27	\$71.46	\$51.88	\$84.21

368

369

370 **9. Author statements**

371 **9.1 Conflicts of interest**

372 The authors declare that there are no conflicts of interest, including financial interests,
373 activities, relationships, and affiliations.

374 **9.2 Funding information**

375 This work was supported by the National Institutes of Health (grant numbers R01AI127472
376 and R21AI109459).

377 **9.3 Ethical approval**

378 The University of Pittsburgh institutional review board provided ethics approval for this study.

379 **9.4 Acknowledgements**

380 The authors would like to thank SeqCenter for their assistance with WGS. We thank the
381 leaders and staff of the UPMC Clinical Laboratories, especially Tung Phan, MD, PhD,
382 D(ABMM), and Hannah Creager PhD, D(ABMM) and all members of the UPMC

383 Presbyterian/Shadyside Infection Prevention & Control Team, especially Graham Snyder,
384 MD, Ashley Ayres, MBA, CIC for their continued support. This publication made use of the
385 PubMLST website (<https://pubmlst.org/>) developed by Keith Jolley (Jolley & Maiden 2010,
386 BMC Bioinformatics, 11:595) and sited at the University of Oxford. The development of that
387 website was funded by the Wellcome Trust.

388 10. References

- 389 1. Scott RD. The Direct medical costs of healthcare-associated infections in U.S. hospitals
390 and the benefits of prevention [Internet]. 2009. 16 p. Available from:
391 <https://stacks.cdc.gov/view/cdc/11550>
- 392 2. Alioth Finance. \$6,500,000,000 in 2007 → 2023 | Inflation Calculator [Internet]. 2023.
393 Available from: <https://www.officialdata.org/us/inflation/2007?amount=6500000000>
- 394 3. Mustapha MM, Srinivasa VR, Griffith MP, Cho ST, Evans DR, Waggle K, et al. Genomic
395 Diversity of Hospital-Acquired Infections Revealed through Prospective Whole-Genome
396 Sequencing-Based Surveillance. *mSystems*. 2022 Jun 28;7(3):e0138421.
- 397 4. Neoh HM, Tan XE, Sapri HF, Tan TL. Pulsed-field gel electrophoresis (PFGE): A review
398 of the “gold standard” for bacteria typing and current alternatives. *Infect Genet Evol J Mol
399 Epidemiol Evol Genet Infect Dis*. 2019 Oct;74:103935.
- 400 5. Quainoo S, Coolen JPM, van Hijum SAFT, Huynen MA, Melchers WJG, van Schaik W, et
401 al. Whole-Genome Sequencing of Bacterial Pathogens: the Future of Nosocomial
402 Outbreak Analysis. *Clin Microbiol Rev*. 2017 Oct;30(4):1015–63.
- 403 6. Sundermann AJ, Babiker A, Marsh JW, Shutt KA, Mustapha MM, Pasculle AW, et al.
404 Outbreak of Vancomycin-resistant *Enterococcus faecium* in Interventional Radiology:
405 Detection Through Whole-genome Sequencing-based Surveillance. *Clin Infect Dis Off
406 Publ Infect Dis Soc Am*. 2020 May 23;70(11):2336–43.
- 407 7. Sundermann AJ, Chen J, Kumar P, Ayres AM, Cho ST, Ezeonwuka C, et al. Whole-
408 Genome Sequencing Surveillance and Machine Learning of the Electronic Health Record
409 for Enhanced Healthcare Outbreak Detection. *Clin Infect Dis*. 2022 Aug 31;75(3):476–82.
- 410 8. Didelot X, Bowden R, Wilson DJ, Peto TEA, Crook DW. Transforming clinical
411 microbiology with bacterial genome sequencing. *Nat Rev Genet*. 2012 Sep;13(9):601–12.
- 412 9. Sundermann AJ, Rangachar Srinivasa V, Mills EG, Griffith MP, Waggle KD, Ayres AM, et
413 al. Two Artificial Tears Outbreak-Associated Cases of Extensively Drug-Resistant
414 *Pseudomonas aeruginosa* Detected Through Whole Genome Sequencing-Based
415 Surveillance. *J Infect Dis*. 2023 Sep 13;jiad318.

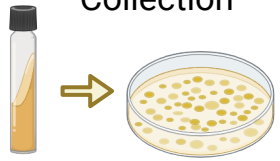
- 416 10. Raabe NJ, Valek AL, Griffith MP, Mills E, Waggle K, Srinivasa VR, et al. Genomic
417 Epidemiologic Investigation of a Multispecies Hospital Outbreak of NDM-5-Producing
418 Enterobacteriales Infections. *MedRxiv Prepr Serv Health Sci*. 2023 Sep
419 1;2023.08.31.23294545.
- 420 11. Sundermann AJ, Griffith M, Rangachar Srinivasa V, Ereifej D, Waggle K, Van Tyne
421 D, et al. Environmental contamination of postmortem blood cultures detected by whole-
422 genome sequencing surveillance. *Infect Control Hosp Epidemiol*. 2023 Aug 24;1–2.
- 423 12. Sundermann AJ, Chen J, Miller JK, Saul MI, Shutt KA, Griffith MP, et al. Outbreak of
424 *Pseudomonas aeruginosa* Infections from a Contaminated Gastroscope Detected by
425 Whole Genome Sequencing Surveillance. *Clin Infect Dis Off Publ Infect Dis Soc Am*.
426 2021 Aug 2;73(3):e638–42.
- 427 13. Crobach MJT, Planche T, Eckert C, Barbut F, Terveer EM, Dekkers OM, et al.
428 European Society of Clinical Microbiology and Infectious Diseases: update of the
429 diagnostic guidance document for *Clostridium difficile* infection. *Clin Microbiol Infect Off*
430 *Publ Eur Soc Clin Microbiol Infect Dis*. 2016 Aug;22 Suppl 4:S63-81.
- 431 14. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinforma Oxf Engl*.
432 2014 Jul 15;30(14):2068–9.
- 433 15. Seemann T. mlst [Internet]. Available from: <https://github.com/tseemann/mlst>
- 434 16. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2.
435 *Genome Biol*. 2019 Nov 28;20(1):257.
- 436 17. Seemann T. Snippy [Internet]. 2023. Available from:
437 <https://github.com/tseemann/snippy>
- 438 18. Seemann T. snp-dists [Internet]. 2023. Available from:
439 <https://github.com/tseemann/snp-dists>
- 440 19. Harris SR. SKA: Split Kmer Analysis Toolkit for Bacterial Genomic Epidemiology
441 [Internet]. *Genomics*; 2018 Oct [cited 2023 Oct 4]. Available from:
442 <http://biorxiv.org/lookup/doi/10.1101/453142>
- 443 20. Price V, Ngwira LG, Lewis JM, Baker KS, Peacock SJ, Jauneikaite E, et al. A
444 systematic review of economic evaluations of whole-genome sequencing for the
445 surveillance of bacterial pathogens. *Microb Genomics*. 2023 Feb;9(2):mgen000947.
- 446 21. Havelaar AH, Kirk MD, Torgerson PR, Gibb HJ, Hald T, Lake RJ, et al. World Health
447 Organization Global Estimates and Regional Comparisons of the Burden of Foodborne
448 Disease in 2010. *PLoS Med*. 2015 Dec;12(12):e1001923.
- 449 22. Kumar P, Sundermann AJ, Martin EM, Snyder GM, Marsh JW, Harrison LH, et al.
450 Method for Economic Evaluation of Bacterial Whole Genome Sequencing Surveillance

- 451 Compared to Standard of Care in Detecting Hospital Outbreaks. Clin Infect Dis. 2021 Jul
452 1;73(1):e9–18.
- 453 23. Wetterstrand K. DNA Sequencing Costs: Data from the NHGRI Genome Sequencing
454 Program (GSP) Available at: www.genome.gov/sequencingcostsdata.
- 455

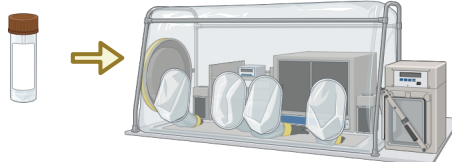
A. Collect List Generation



B1. Bacterial Isolate Collection



B2. *C. difficile* Collection and Isolation



C. Sample Processing



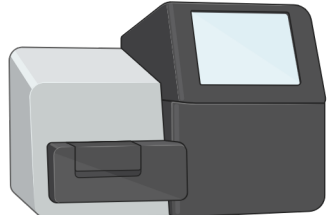
D. DNA Extraction



E. Library Preparation



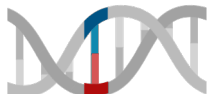
F. Whole Genome Sequencing



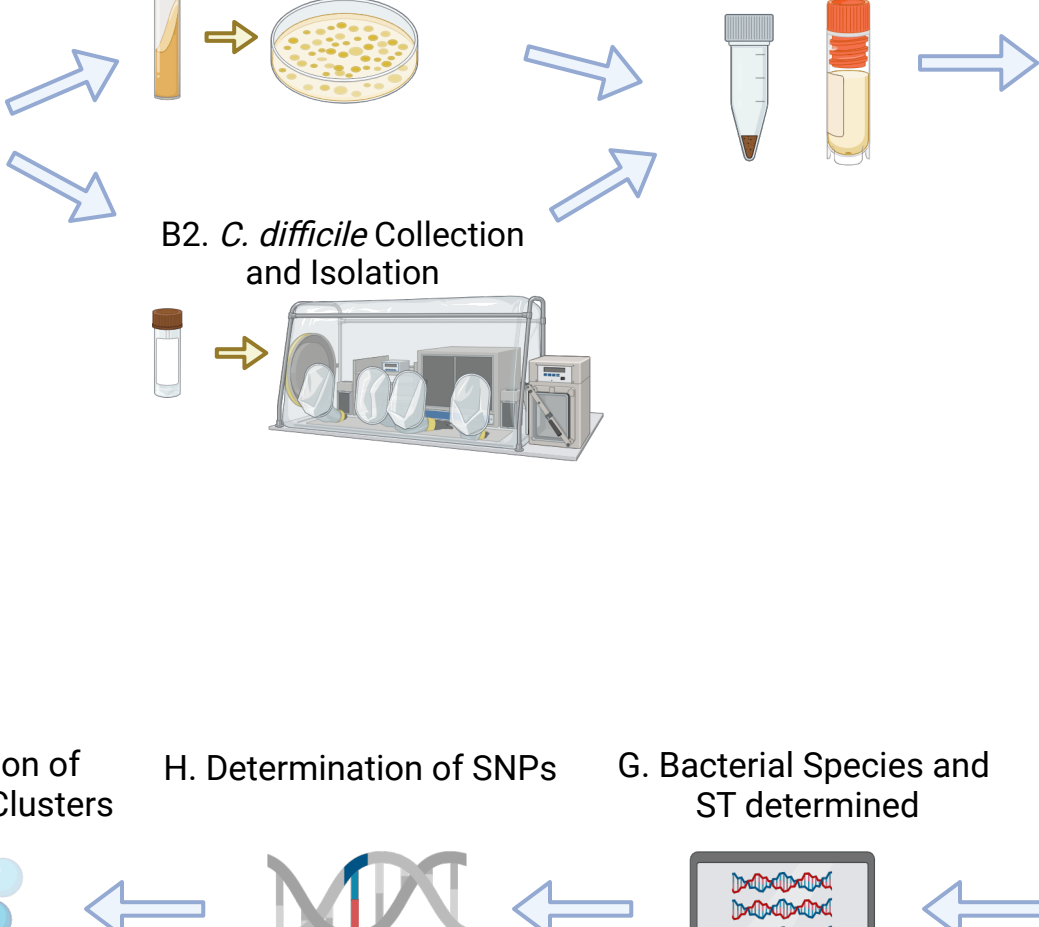
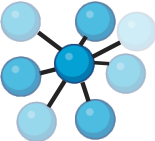
G. Bacterial Species and ST determined



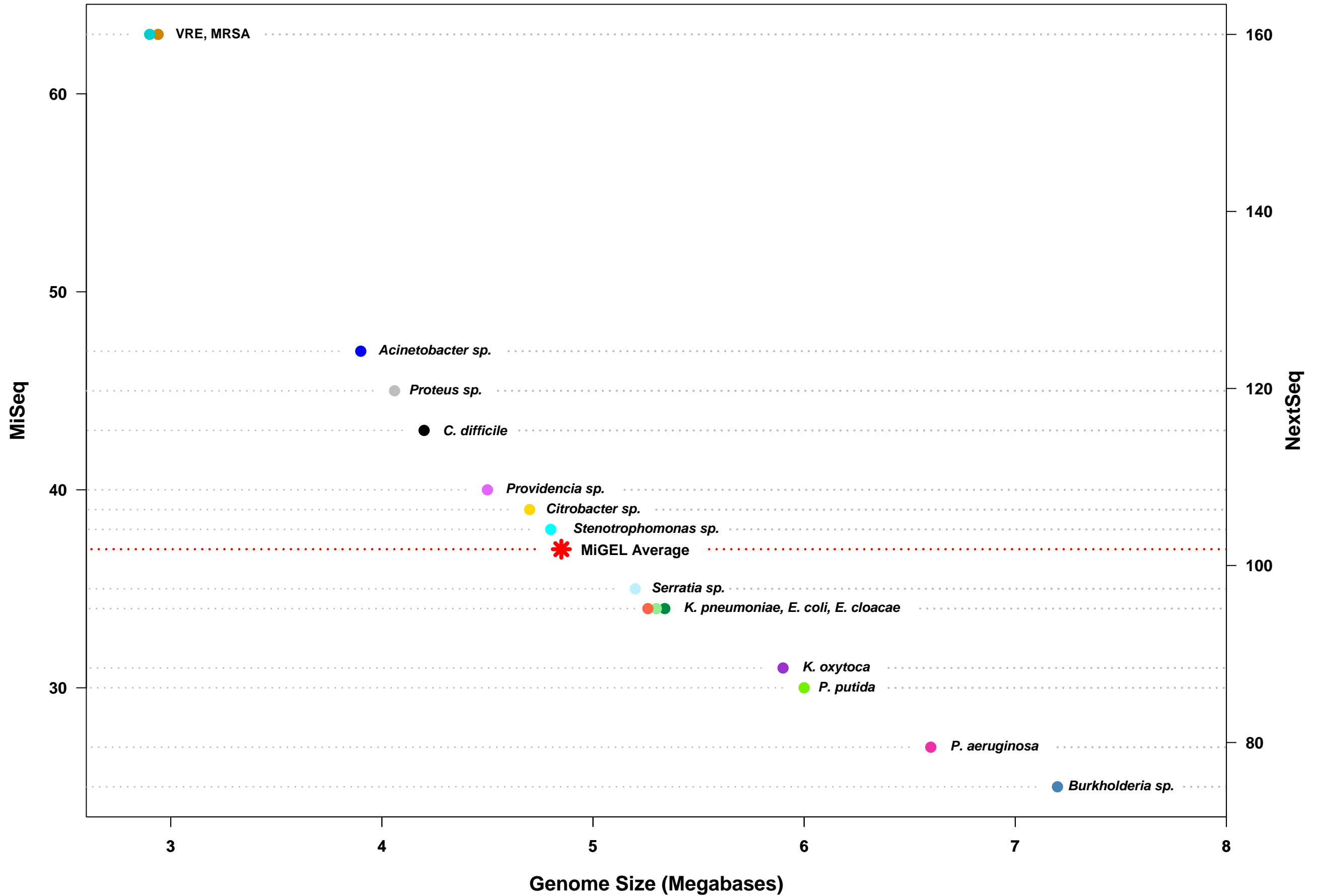
H. Determination of SNPs



I. Determination of Transmission Clusters



Maximum Sample Counts



Cost per Sample

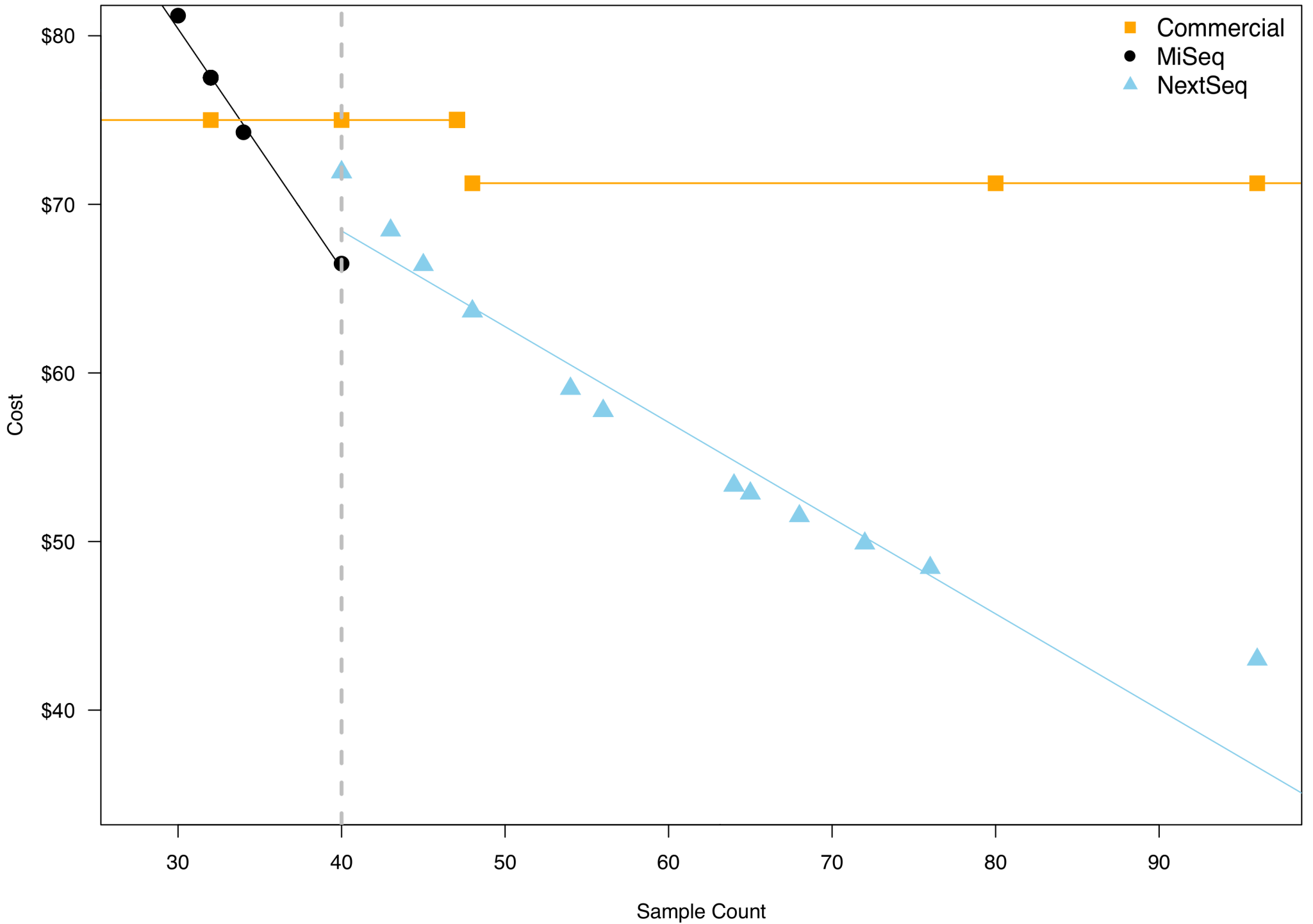


Figure 1. EDS-HAT real-time genomic surveillance methods. A) Biweekly collect list generation for all EDS-HAT organisms of interest obtained from patients admitted to the hospital for ≥ 3 days or had a previous hospital exposure in the prior 30 days, B1) Bacterial isolate collection from the UPMC Clinical Microbiology Laboratory, B2) Clinical stool specimen collection and *C. difficile* isolation using a Coy anaerobic chamber, C) Processing samples into pellets and glycerol stocks, D) DNA extraction, E) Library preparation using Eppendorf epMotion 5075, F) WGS using MiSeq or NextSeq550, G) Bioinformatic analysis to determine bacterial species and sequence type (ST), H) Determination of SNPs between isolates, and I) Determination of transmission clusters. The average turnaround time from MiGEL sample collection to determination of transmission clusters was 10 days.

Figure 2. Maximum number of genomes that can be sequenced on the MiSeq v3 600 cycle flow cell and NextSeq550 v2.5 300 cycle flow cell based on size (Mb). Sample counts were calculated using Illumina coverage calculator based on 80 \times coverage criteria. Genome size of an average run by MiGEL is shown in comparison to individual organism sizes (red star).

Figure 3. Whole genome sequencing cost per sample comparison between MiSeq, NextSeq550, and commercial laboratory. Throughput cutoff between platforms is shown at $N=40$ samples (gray dashed line). The commercial lab used by MiGEL offers a discount of 5% for orders ≥ 48 samples (as of December 2023). Data points represent instances of cost by sample count and lines of best fit are shown for each sequencing method.

Supplementary Figure 1. EDS-HAT bioinformatics pipeline.

Note, ST = sequence type; SNP = single nucleotide polymorphism; EHR = electronic health record).