1 Solu – a cloud platform for real-time genomic

² pathogen surveillance

- 3 Timo J Moilanen^{1*}, Kerkko Visuri¹, Jonatan Lehtinen¹, Irene Ortega-Sanz², Jacob L
- 4 Steenwyk³, Samuel Sihvonen¹

- 6 1 Solu Healthcare Oy, Fredrikinkatu 63 A 302, 00100 Helsinki, Finland
- 7 2 Department of Food Technology, Safety and Health, Faculty of Bioscience Engineering,
- 8 Ghent University, Coupure Links 653, 9000 Ghent, Belgium
- 9 3 Howards Hughes Medical Institute and the Department of Molecular and Cell Biology,
- 10 University of California, Berkeley, Berkeley, CA, USA
- 11
- 12 *Corresponding author: timo@solugenomics.com (Timo J Moilanen)
- 13
- 14
- 15

16 Abstract

17 Background

18	Genomic surveillance is extensively used for tracking public health outbreaks and
19	healthcare-associated pathogens. Despite advancements in bioinformatics pipelines, there
20	are still significant challenges in terms of infrastructure, expertise, and security when it
21	comes to continuous surveillance. The existing pipelines often require the user to set up and
22	manage their own infrastructure and are not designed for continuous surveillance that
23	demands integration of new and regularly generated sequencing data with previous
24	analyses. Additionally, academic projects often do not meet the privacy requirements of
25	healthcare providers.

26 Results

- 27 We present Solu, a cloud-based platform that integrates genomic data into a real-time,
- 28 privacy-focused surveillance system.

29 Evaluation

- 30 Solu's accuracy for taxonomy assignment, antimicrobial resistance genes, and
- 31 phylogenetics, was comparable to established pathogen surveillance pipelines. In some
- 32 cases, Solu identified antimicrobial resistance genes that were previously undetected.
- 33 Together, these findings demonstrate the efficacy of our platform.

34 Conclusions

- 35 By enabling reliable, user-friendly, and privacy-focused genomic surveillance, Solu has the
- 36 potential to bridge the gap between cutting-edge research and practical, widespread

- 37 application in healthcare settings. The platform is available for free academic use at
- 38 platform.solugenomics.com.
- 39 Keywords: Workflow, Genomics, Whole-genome sequencing, Infection prevention,
- 40 Outbreak, Phylogeny, Privacy

41 Background

42 Bacterial and fungal pathogens, along with their antimicrobial resistance, are causing an 43 increasing burden on healthcare and public health (1-3). Advances in microbial genomics 44 have significantly enhanced infection prevention and outbreak surveillance by providing 45 detailed information about pathogen species, antimicrobial resistance, and phylogenetics 46 (4,5). As the cost of Whole-Genome Sequencing (WGS) has decreased rapidly, continuous 47 genomic surveillance has become a cost-effective method for infection prevention and 48 control (6,7). The interest towards genomic analysis has led to the emergence of several 49 pathogen analysis tools, such as nf-core (8), TheiaProk (9), ASA3P (10), CamPype (11), 50 Nullarbor (12), Bactopia (13), and Galaxy (14), which enable genomic analysis also for users 51 without in-depth expertise in bioinformatics or computer science. 52 Despite these advancements, bioinformatics still remains a bottleneck for the widespread 53 adoption on pathogen genomic surveillance due to limitations in usability, speed, and 54 security (7). 55 Most existing pipelines (8–13) are operated using the command-line interface (CLI) and 56 require the user to manage their own data storage and computation infrastructure. While it is 57 possible to learn their usage without advanced computational knowledge (15), many

- 58 practitioners simply don't have the time or willingness for it and prefer graphical user
- 59 interfaces instead. Additionally, most existing tools are designed for single-use execution,
 - 3

60 which is a challenge for continuous surveillance where new sequencing data is often

61 generated in small batches (6,16). To facilitate ongoing analysis, users must implement their

62	own processes for integrating new and old data.
63	Fast time to results has been identified as a key component for effective genomic
64	surveillance (6,16). As new samples arrive in batches and need to be compared to all
65	previously accumulated samples, computation time can become a significant bottleneck if
66	using a single-workstation installation. Also, unless the pipelines are highly automated,
67	running the analyses often requires specially trained personnel who might not be available
68	immediately upon the arrival of new data.
69	Academia-led projects developed under FAIR (Findable, Accessible, Interoperable,
70	Reusable) principles often lack the necessary privacy focus to meet the stringent

requirements of healthcare providers (17). In contrast, healthcare providers must adhere to

stringent legal requirements, such as the U.S. HIPAA Privacy Rule (18), ruling out many

73 existing online platforms for genomic surveillance.

74 It is possible for healthcare providers to overcome these limitations by implementing their

75 own automated pipelines, but it requires significant investments in bioinformatics and

76 computational infrastructure, and the lack of these resources is a challenge in many facilities

(19). To fill this gap, we present Solu – an automated, fast, and secure web application for

78 analyzing WGS samples.

79 Implementation

Solu is a cloud-based platform for the analysis of bacterial and fungal WGS samples. Its
automated bioinformatics pipeline includes genomic characterization and phylogenetic

- 82 comparison. Its cloud implementation is built to match the usability, speed, and security
- 83 requirements of ongoing genomic surveillance in healthcare facilities.

84 Bioinformatics pipeline

- 85 The platform runs a fully automated pathogen analysis pipeline, which is illustrated in Figure
- 1. The pipeline includes *de novo* assembly, quality assurance (QA), species identification
- 87 and genomic characterization for each uploaded sample, and phylogenetic comparison
- 88 between all uploaded samples of the same species. It is triggered automatically after each
- 89 file upload and cannot be configured by the user. This section presents an overview of the
- 90 pipeline, and a detailed description can be found in Supplementary Material 1.



92 Figure 1. Bioinformatics pipeline

93 Input format

- 94 The supports three input types: paired-end short reads in FASTQ format, long reads in
- 95 FASTQ format, or an assembled genome in FASTA format. Analysis of long reads is still
- 96 considered an experimental feature.

97 Assembly, QA, and species detection

- 98 Short reads are quality checked using FastQC (20), quality corrected with fastp (21), and
- 99 assembled using Shovill (22). Long reads are pre-processed, assembled and polished with

100 Dragonflye (23). After assembly, all samples are standardized using any2fasta (24) and

- 101 quality assessed with Quast (25).
- 102 Species is identified with Bactinspector (26). To identify fungal species, Bactinspector's
- 103 default database was augmented with all fungal reference genomes from the NCBI
- 104 Taxonomy (27). The augmented database also includes clade-level reference genomes for
- 105 Candida auris.

106 Genomic characterization

- 107 Analysis of bacterial species includes multi-locus sequence typing (MLST) with mlst (28),
- 108 AMR annotation using AMRFinderPlus (29), and plasmid analysis with MOB-suite (30).
- 109 The pipeline also includes an experimental antifungal resistance (AFR) gene annotation for
- 110 the species Candida auris. AFR annotation is implemented using AMRFinderPlus with a
- 111 custom database of known AFR point mutations sourced from AFRBase (31).

112 Phylogenetic comparison

- 113 The pipeline's phylogeny is based on constructing a multiple sequence alignment for each
- species. Based on the species in question, multiple sequence alignment is computed by
- 115 either a reference-based or reference-free method.
 - 7

- 116 The reference-based alignment is considered more robust and has been implemented for 21
- 117 commonly analyzed species. It includes aligning each sample to the species' reference
- 118 genome using Snippy (32), creating a multiple-sequence alignment using snippy-core (32),
- and filtering out low quality SNPs using an in-house script.
- 120 The reference-free alternative is implemented to support analysis of species that are not yet
- 121 supported by the reference-based alignment. It is computed using the split-kmer analysis
- 122 tool SKA (33).
- 123 After constructing the multiple sequencing alignment, the phylogenetic comparison includes
- 124 pairwise SNP distances, clustering, and phylogenetic tree inference. SNP distances are
- 125 counted from the multiple sequence alignment using snp-sites (34) and snp-dists (35).
- 126 Samples are clustered with a 20-SNP single-linkage clustering threshold using an in-house
- 127 Python script. Phylogenetic trees are inferred using a general time reversible maximum
- 128 likelihood model from IQ-TREE 2 (36) and midpoint-rooted using TreeTime (37). Both IQ-
- 129 TREE 2 and TreeTime are run using the Augur toolkit (38).

130 Automated cloud infrastructure

131 The cloud infrastructure of Solu is built on three principles: usability, speed, and security.



132

133 Figure 2. Solu's cloud platform implementation

134 Usability

The Solu Platform is web based, enabling practitioners to use it without installing software or running command-line tools (Figure 2). New samples are uploaded using the drag-and-drop web UI, which automatically triggers a bioinformatics pipeline. The pipeline requires zero configuration from the user, which promotes repeatability and alleviates the need for in-depth bioinformatics knowledge. The analysis results are stored in the cloud, eliminating the need for a self-implemented storage system and enabling effortless result sharing betweencolleagues.

Each newly uploaded sample is also automatically compared to previously uploaded
samples of the same species, which enables detecting potential new outbreaks quickly. For
instance, when a user uploads a *Salmonella enterica* sample, the platform automatically recomputes the phylogenetics for all uploaded *Salmonella enterica* samples and highlight
possible clusters.

147 **Speed**

The platform's cloud infrastructure is optimized for speed even during peak usage. This is achieved by running each computation-intensive workload in a separate Docker container with optimized resource (CPU and memory) distribution. These containers are orchestrated by a cloud computing cluster that is automatically scaled up and down based on usage, up to a maximum of 512 CPUs and 2 TB of memory. The cluster also contains a pool of hot standby resources, which allows starting the analysis of a new sample within seconds of its upload.

This auto-scaling capability brings the user substantial speed improvements by allowing the parallelization of some of the analyses in the pipeline. In addition, it allows analyzing a whole batch of samples simultaneously, leading to a significant reduction in overall time-to-results when analyzing a batch of samples. Importantly, these speed improvements are achieved without a significant increase in computation costs.

160 Security

161 The platform's data is stored in a secure cloud storage with set read and write permissions.

162 All computations occur within a virtual private network, monitored by automated access

- 163 control checks. Solu implements strict data security protocols, including appropriate access
- 164 permissions, encryption, continuous monitoring, code reviews, staff training, and other
- 165 cybersecurity measures. Accordingly, Solu adheres to the U.S. HIPAA rule, and can sign a
- 166 Business Associate Agreement (BAA) for enterprise customers. Solu also allows enterprise
- 167 customers to choose between U.S. or EU as their data storage. Further information
- 168 regarding data security practices can be found at <u>https://solugenomics.trust.site/</u>.

169 Evaluation

170 To evaluate the Solu platform, we reproduced four outbreak investigation studies using

171 published genomic data (Table 1). Data was obtained from the European Nucleotide Archive

- as raw reads and uploaded to the Solu platform. All samples were paired-end short-reads in
- 173 FASTQ format.
- 174 We evaluated Solu's performance by computing metric scores for species identification,
- 175 MLST, clade construction, and AMR predictions against the references. Phylogenetic trees
- 176 were exported from Solu. Tree topologies were compared visually, and where raw tree data
- 177 was available, we calculated a Robinson-Foulds distance using TreeDist (39). Both
- 178 reference-based and reference-free phylogenetic pipelines were run for all datasets and
- 179 compared against each other to validate the platform's internal consistency.
- 180 We also measured the required time for analysis of each sample. Plasmids were not
- 181 evaluated in this study due to the absence of plasmid annotations in the original studies.
- 182

Species	BioProject accession	Number of samples	Raw data size (GB, zipped)
Staphylococcus aureus	PRJNA400143 (40)	135	16.47
Enterococcus faecium	PRJEB34664 (41)	99	15.18
Salmonella enterica	Multiple (42)	23	5.03

183 Table 1. Overview of evaluation datasets

Candida auris	PRJNA328792 (43)	47	50.94
---------------	------------------	----	-------

184 Results

185 Evaluation results

- 186 The Solu platform successfully completed the bioinformatics pipeline for all 304 samples. A
- 187 screenshot of the platform's home screen is shown in Figure 3. This workspace, including all
- 188 samples and results, is also accessible at a user-friendly web interface at
- 189 <u>https://platform.solugenomics.com/w/solu-publication-2</u> and
- 190 <u>https://platform.solugenomics.com/w/solu-publication-3.</u>

•• 🗆 -	< >			ii platfor	m.sok	sgen	omics.com			C			0 + 4	
b solu-publication-2		ERR3569619	Enterococcus faecium	80		15	(WE)	2	•		Not computed		Z months ago	
 Gamples All samples 	2	ERR3569664	Enterococcus faecium	117	W	33	WE	54	11	14	Not computed		2 months ago	
E. faecium		ERR3569683	Enterococcus faocium	117	19	35	VHE .		10		Not computed		2 months ago	
 S. enterica S. aureus 		ERR3569662	Enterococcus faecium	292	190) 1900	43	(VHE)	~	11	390	Not computed		2 months ago	
Clusters	>	ERR3569661	Enterococcus faecium	80		15	990	2	10	~	Not computed		2 months ago	
		ERR3569660	Enterpococcus faecium	117	*	14	VHE	19		199	Not complated		2 months aga	
		ERR3569659	Enterococcus faecium	80		15	sec.	10	9		Not computed	**	2 months ago	
		ERR3569658	Enterococcus faecium	80	290	15	VIE		0	СФ.	Not computed	**	2 months ago	
		ERR3569657	Enterococcus faecium	117	•	35	(VRC)	2	9	~	Not computed		2 months ago	
			ERR2569656	Enterprocess faecium	117	1	19	(Int)		10	1	Not computed		2 months ago
		ERR3569655	Enterococcus lancium	80	Y	35	SHE	i i i	11	*	Not asimputed		2 months ago	
		ERR3569653	Enterococcus faecium	78	1	16	VIE	3	9		Net computed		2 months ago	
🕼 Help & feedback		ERR3569652	Enterococcus faecium	\$17	~	19	SHE	9	10	1	Not computed		2 months ago	
<u>Den un</u> her ynur own Solur		ERR3569651	Enterococcus faecium	117	-	35	VIE	÷	10	÷	Not computed		Z months ago	

191 192

Figure 3. Summary of the samples shown in Solu

193 Species, MLST and clade assignment

- 194 Solu accurately identified the species of all 304 samples and assigned the correct clade for
- 195 all 47 Candida auris samples.
- 196 Exact MLST matches were observed in 210 out of 230 (91.3%) isolates with known
- 197 sequence types (Table 2). However, 18 of 20 non-exact MSLT matches were single-locus

198 variants.

- 199
- 200
- 201 Table 2. Species identification and MLST concordance. Solu's results compared to the original
- 202 publications when available there

Species	Species	MLST exact match	MLST accuracy		
	identification	accuracy	including single-		
	accuracy		locus variants		
S. aureus	100.0% (135/135)	87.1% (115/132)	98.4% (130/132)		
E. faecium	100.0% (99/99)	96.9% (95/98)	100.0% (98/98)		
S. enterica	100.0% (23/23)	MLST not reported in	the original article		
C. auris	100.0% (47/47)	MLST scheme not defined			

203

204 Antimicrobial resistance

- 205 The antimicrobial resistance (AMR) gene detection results from Solu were compared with
- those of the references. Concordance varied by species, ranging from 99.6% for *E. faecium*

- 207 to 93.1% for *S. aureus*. Table 3 summarizes some commonly studied AMR loci, while the full
- 208 results can be viewed online.

209	Table 3. Sensitivity of Solu's AMR prediction
-----	---

Species	Overall AMR locus sensitivity	vanA	vanB	mecA	Other key genes			
S. aureus	93.1% (309/332)	N/A	N/A	98.8% (81/82)	<i>mupA</i> : 100% (3/3), <i>blaZ</i> : 90.9% (100/110)			
E. faecium	99.6% (515/517)	100.0%	100.0% (69/69)	N/A	<i>ermB</i> :100.0% (95/95) <i>tet</i> : 100.0% (4/4)			
S. enterica	S. enterica AMR results not reported in the original article							
C. auris	90.9% (40/44)	N/A	N/A	N/A	<i>ERG11_K143R</i> : 100.0% (8/8)			

210 Note: N/A indicates that the gene is not typically relevant for the species.

211 Key AMR genes, such as the *vanA/vanB* type and *mecA*, were detected with 100% and

212 98.8% sensitivity, respectively. Antifungal resistance mutation detection for Candida auris

showed a 90.9% sensitivity. The results matched the original findings for 43 isolates.

However, Solu identified the ERG11_K143R mutation in 2 Clade I isolates, which were

215 originally reported as having the *Y123F* mutation, and detected 2 isolates lacking ERG

216 mutations.

- 217 The main reason for lower agreement in the S. aureus dataset was Solu's inability to find
- any *dfrA* matches, whereas the original article reported 10 isolates with *dfrA*. (44)

219 Phylogenetics

- 220 Solu automatically generated phylogenetic trees for all four datasets, which can be viewed
- and downloaded in the published workspace in "Tree view" (Figure 4).
- 222 The E. faecium phylogenetic tree generated by Solu demonstrated a high degree of
- 223 concordance with the reported SNP subclusters, complex types (CTs) and sequence types
- (STs) of the reference (Figure 5). For the Salmonella enterica dataset, Solu produced a
- similar topology to the reference tree (Figure 6) where the outbreak samples are separate
- from the outgroup. Robinson-Foulds distance to the *S. enterica* reference tree was 2.



227

228 Figure 4. Screenshot of the Salmonella enterica tree in the graphical user interface



- 231
- 232 Figure 5. Solu's *E. faecium* phylogenetic tree shown as a cladogram (inner) vs. SNP
- 233 clusters, complex types (CT), and sequence types (ST) from the original publication (outer
- rings). The cladogram visualization was created using Dendroscope (45).



Figure 6. Tanglegram of Solu's *Salmonella enterica* phylogenetic tree (left) vs. the reference
tree (right).

238

239 For the S. aureus and C. auris datasets, Solu generated phylogenetic trees in which isolates

240 with identical sequence types or clades consistently clustered together. Further detailed

comparisons were not possible due to the lack of raw tree data and subtyping information.

242 The resulting trees are provided in Supplementary Material 2.

243

244 In comparing the reference-free and reference-based pipelines, Solu's reference-free

pipeline generated highly concordant phylogenetic trees with the reference-based pipeline.

246 The Robinson-Foulds distances ranged from 0.08 to 0.46, as computed using TreeDist (46),

247 indicating a high level of similarity (see Supplementary Material 2)

248

249 Time-to-results

- 250 Time-to-results for the four datasets are presented in Table 4 and Figure 7. For bacterial
- samples, Solu completed de novo assembly in an average of 7.2 minutes and variant calling

- in 9.5 minutes from upload. For *C. auris* samples, de novo assembly averaged 17 minutes,
- and variant calling took 23.6 minutes from upload.
- Table 4. Average time-to-results per dataset. Shortest and longest recorded times in
- 255 parentheses. Genome size (27) and sample count of each dataset included for additional
- 256 context.

			Average total time from sample upload (minutes						
	Genome size (Mb)	Sample count	Read quality analysis and correction	De novo assembly	Variant calling				
Enterococcus faecium	2.9	99	1.5 (1.0-2.8)	5.2 (1.9-11.8)	7.4 (3.7–14.4)				
Salmonella enterica	5.0	23	1.4 (1.1–1.8)	5.0 (2.9–10.1)	7.6 (5.4–12.7)				
Staphylococcus aureus	2.8	135	1.9 (1.0–6.1)	9.0 (2.4–21.9)	11.4 (4.2–25.6)				
Candida auris	12.2	47	3.5 (2.4–5.9)	17.0 (10.4–32.3)	23.6 (15.1–37.9)				





259 Discussion

260 Bioinformatics remains a bottleneck to widespread use of genomic pathogen surveillance.

261 Usability, speed, and security are additional requirements for practical outbreak analysis in a

healthcare setting.

263 Our evaluation demonstrates that the Solu platform produces outputs that are largely

264 consistent with prior outbreak studies, using raw sequencing reads and requiring zero

- configuration, with a runtime of approximately 10 minutes for bacterial samples and 20
- 266 minutes for fungal samples. Solu's phylogenetic pipelines produced results that were
- 267 internally and externally consistent.

The largest discrepancies were observed in the *Staphylococcus aureus* dataset, where the original study used PCR for MLST assignment and applied 90% identity and 75% coverage thresholds for AMR gene detection. We hypothesize that the different pipeline parameters allow for higher sensitivity, at the cost of potential misidentification.

272 Compared to some other pipelines, Solu platform's zero-configuration design prevents users 273 from customizing pipeline parameters, which may result in some variation in the results. This 274 approach was chosen to promote usability and prevent users from inadvertently selecting 275 unsuitable parameters. Despite this limitation, default tool configurations provide sufficient 276 accuracy for a wide variety of research applications, including AMR gene characterization 277 and clonality assessment (47). Future studies leveraging more in-depth datasets and 278 epidemiologically validated outbreaks hold great potential to further strengthen and expand 279 the applicability of our findings.

- 280 We aim to improve the analytical capacity of the platform in future iterations, featuring
- additional tooling, modifications to the analytical workflow, broader support for species and
- databases, and improved runtimes among other features. We encourage users to contact
- the authors to request any additional analyses or databases of interest.
- 284 In conclusion, by focusing on a robust, privacy-focused infrastructure, Solu facilitates
- 285 broader adoption of genomic pathogen surveillance, potentially bridging the gap between
- research and practice.

287 Availability and requirements

- 288 Project name: Solu Platform
- 289 Project home page: <u>https://platform.solugenomics.com</u>

- 290 Operating system(s): Platform independent
- 291 Programming language: Typescript
- 292 Other requirements: Modern internet browser such as Google Chrome, Mozilla Firefox or
- 293 Safari
- 294 Pricing, academic free plan, and non-academic license descriptions:
- 295 www.solugenomics.com/pricing

296 List of abbreviations

- 297 AFR: antifungal resistance; AMR: antimicrobial resistance; CPU: Central Processing Unit;
- 298 CT: complex type; FAIR: Findable, Accessible, Interoperable, Reusable; HIPAA: Health
- 299 Insurance Portability and Accountability Act; MLST: Multi-locus Sequence Typing; QA:
- 300 quality assurance; SNP: single nucleotide polymorphism; ST: sequence type; WGS: Whole-
- 301 Genome Sequencing.

302 **Declarations**

³⁰³ Ethics approval and consent to participate

- 304 Not applicable.
- 305 Consent for publication
- 306 Not applicable.

307 Availability of data and materials

- 308 The data analyzed in the study are available from the European Nucleotide Archive at
- 309 https://www.ebi.ac.uk/ena using the accessions provided in Table 1. The Solu workspace,
- 310 including all samples and results, is accessible at https://platform.solugenomics.com/w/solu-
- 311 publication-2 and https://platform.solugenomics.com/w/solu-publication-3

312 Competing interests

- 313 TJM, JL, KV and SS have been employed by Solu Healthcare Oy. JLS is an advisor for
- 314 ForensisGroup Inc.

315 Funding

- 316 JLS is a Howard Hughes Medical Institute Awardee of the Life Sciences Research
- 317 Foundation.

318

319 Authors' contributions

- 320 TM: Conceptualization, Software, Formal Analysis, Investigation, Project Administration,
- 321 Writing. KV: Software, Visualization, Methodology, Validation, Writing. JL: Software,
- 322 Methodology, Validation, Writing. IOS: Methodology, Review & Editing, Visualization. JLS:
- 323 Writing -- Review & Editing, Funding Acquisition, and Visualization. SS: Conceptualization,
- 324 Writing. All authors read and approved the final manuscript.

325 Acknowledgements

326 The authors thank the anonymous reviewers for their constructive comments.

328 References

- Murray CJL, Ikuta KS, Sharara F, Swetschinski L, Robles Aguilar G, Gray A, et al.
 Global burden of bacterial antimicrobial resistance in 2019: a systematic analysis. The Lancet. 2022 Feb;399(10325):629–55.
- Vallabhaneni S, Mody RK, Walker T, Chiller T. The Global Burden of Fungal Diseases.
 Infect Dis Clin. 2016 Mar 1;30(1):1–11.
- Van Rhijn N, Arikan-Akdagli S, Beardsley J, Bongomin F, Chakrabarti A, Chen SCA, et
 al. Beyond bacteria: the growing threat of antifungal resistance. The Lancet. 2024
 Sep;404(10457):1017–8.
- 4. Eyre DW. Infection prevention and control insights from a decade of pathogen whole genome sequencing. J Hosp Infect. 2022 Apr;122:180–6.
- Steenwyk JL, Rokas A, Goldman GH. Know the enemy and know yourself: Addressing cryptic fungal pathogens of humans and beyond. Chowdhary A, editor. PLOS Pathog.
 2023 Oct 19;19(10):e1011704.
- Forde BM, Bergh, Haakon, Cuddihy, Thom, Hajkowich, Kristin, Hurst, Trish, Playford, E, et al. Clinical Implementation of Routine Whole-genome Sequencing for Hospital Infection Control of Multi-drug Resistant Pathogens. Clin Infect Dis. 2023 Feb 1;76(3):e1277–84.
- Fox JM, Saunders NJ, Jerwood SH. Economic and health impact modelling of a whole genome sequencing-led intervention strategy for bacterial healthcare-associated infections for England and for the USA. Microb Genomics [Internet]. 2023 Aug 9 [cited 2024 Sep 17];9(8). Available from:
- 350 https://www.microbiologyresearch.org/content/journal/mgen/10.1099/mgen.0.001087
- Ewels PA, Peltzer A, Fillinger S, Patel H, Alneberg J, Wilm A, et al. The nf-core framework for community-curated bioinformatics pipelines. Nat Biotechnol. 2020 Mar;38(3):276–8.
- Libuit KG, Doughty EL, Otieno JR, Ambrosio F, Kapsak CJ, Smith EA, et al. Accelerating bioinformatics implementation in public health. Microb Genomics. 2023 Jul 10;9(7).
- 356
 357
 358
 358
 359
 359
 350
 350
 350
 351
 352
 353
 354
 355
 355
 355
 356
 357
 357
 358
 359
 359
 359
 359
 359
 359
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
 350
- 360 11. Ortega-Sanz I, Barbero-Aparicio JA, Canepa-Oneto A, Rovira J, Melero B. CamPype: an
 361 open-source workflow for automated bacterial whole-genome sequencing analysis
 362 focused on Campylobacter. BMC Bioinformatics. 2023 Jul 20;24(1).
- 363 12. Seemann T. nullarbor [Internet]. [cited 2024 Sep 17]. Available from:
 364 https://github.com/tseemann/nullarbor

- 365 13. Petit RA, Read TD. Bactopia: a Flexible Pipeline for Complete Analysis of Bacterial
 366 Genomes. Segata N, editor. mSystems. 2020 Aug 25;5(4).
- 367 14. Jalili V, Afgan E, Gu Q, Clements D, Blankenberg D, Goecks J, et al. The Galaxy
 368 platform for accessible, reproducible and collaborative biomedical analyses: 2020
 369 update. Nucleic Acids Res. 2020 Jul 2;48(W1):W395–402.
- 15. Carriço JA, Rossi M, Moran-Gilad J, Van Domselaar G, Ramirez M. A primer on
 microbial bioinformatics for nonbioinformaticians. Clin Microbiol Infect. 2018
 Apr;24(4):342–9.
- 373 16. Roberts LW, Forde BM, Hurst T, Ling W, Nimmo GR, Bergh H, et al. Genomic
 374 surveillance, characterization and intervention of a polymicrobial multidrug-resistant
 375 outbreak in critical care. Microb Genomics. 2021 Mar 1;7(3).
- 376 17. Stacey D, Wulff K, Chikhalla N, Bernardo T. From FAIR to FAIRS: Data security by
 377 design for the global burden of animal diseases. Agron J. 2022 Sep;114(5):2693–9.
- 18. U.S. Department of Health & Human Services, Rights (OCR) O for C. The HIPAA
 Privacy Rule [Internet]. 2008 [cited 2024 Sep 17]. Available from: https://www.hhs.gov/hipaa/for-professionals/privacy/index.html
- 381 19. Afolayan AO, Bernal JF, Gayeta JM, Masim ML, Shamanna V, Abrudan M, et al.
 382 Overcoming Data Bottlenecks in Genomic Pathogen Surveillance. Clin Infect Dis. 2021
 383 Dec 1;73(Supplement_4):S267–74.
- 384 20. Andrews S. FastQC: a quality control tool for high throughput sequence data.
 385 Cambridge, United Kingdom; 2010.
- 21. Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor.
 Bioinformatics. 2018 Sep 1;34(17):i884–90.
- 388 22. Seemann T. Shovill [Internet]. 2024 [cited 2024 Sep 17]. Available from:
 389 https://github.com/tseemann/shovill
- Craddock HA, Motro Y, Zilberman B, Khalfin B, Bardenstein S, Moran-Gilad J. Long Read Sequencing and Hybrid Assembly for Genomic Analysis of Clinical Brucella
 melitensis Isolates. Microorganisms. 2022 Mar 14;10(3):619.
- 393 24. Seemann T. any2fasta [Internet]. 2024 [cited 2024 Sep 17]. Available from:
 394 https://github.com/tseemann/any2fasta
- 395 25. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for
 396 genome assemblies. Bioinformatics. 2013 Apr 15;29(8):1072–5.
- 397 26. Underwood A. BactInspector [Internet]. 2020 [cited 2024 Sep 17]. Available from:
 398 https://gitlab.com/antunderwood/bactinspector
- 399 27. Schoch CL, Ciufo S, Domrachev M, Hotton CL, Kannan S, Khovanskaya R, et al. NCBI
 400 Taxonomy: a comprehensive update on curation, resources and tools. Database

- 401 [Internet]. 2020 Jan 1 [cited 2024 Sep 17];2020. Available from:
- 402 https://academic.oup.com/database/article/doi/10.1093/database/baaa062/5881509
- 403 28. Seemann T. mlst [Internet]. 2024 [cited 2024 Sep 17]. Available from:
 404 https://github.com/tseemann/mlst
- 405
 406
 406
 406
 407
 408
 408
 409
 409
 409
 409
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
 400
- 30. Robertson J, Nash JHE. MOB-suite: software tools for clustering, reconstruction and
 typing of plasmids from draft assemblies. Microb Genomics. 2018 Aug 1;4(8).
- 31. Jain A, Singhal N, Kumar M. AFRbase: a database of protein mutations responsible for
 antifungal resistance. Martelli PL, editor. Bioinformatics. 2023 Nov 1;39(11).
- 413 32. Seemann T. Snippy [Internet]. 2024 [cited 2024 Sep 17]. Available from:
 414 https://github.com/tseemann/snippy
- 415 33. Harris SR. SKA: Split Kmer Analysis Toolkit for Bacterial Genomic Epidemiology
 416 [Internet]. Cold Spring Harbor Laboratory; 2018 [cited 2024 Sep 17]. Available from:
 417 http://biorxiv.org/lookup/doi/10.1101/453142
- 418 34. Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane A, et al. SNP-sites: rapid
 419 efficient extraction of SNPs from multi- FASTA alignments. Microb Genomics.
- 35. Seemann T. snp-dists [Internet]. 2024 [cited 2024 Sep 17]. Available from:
 https://github.com/tseemann/snp-dists
- 36. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, Von Haeseler A, et
 al. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the
 Genomic Era. Teeling E, editor. Mol Biol Evol. 2020 May 1;37(5):1530–4.
- 37. Sagulenko P, Puller V, Neher RA. TreeTime: Maximum-likelihood phylodynamic
 analysis. Virus Evol. 2018 Jan 1;4(1).
- 38. Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al. Nextstrain: realtime tracking of pathogen evolution. Kelso J, editor. Bioinformatics. 2018 Dec
 1;34(23):4121–3.
- 39. Smith MR. Robust Analysis of Phylogenetic Tree Space. Hoehna S, editor. Syst Biol.
 2022 Aug 10;71(5):1255–70.
- 432 40. Manara S, Pasolli E, Dolce D, Ravenni N, Campana S, Armanini F, et al. Whole-genome
 433 epidemiology, characterisation, and phylogenetic reconstruction of Staphylococcus
 434 aureus strains in a paediatric hospital. Genome Med [Internet]. 2018 Dec [cited 2024
- 435 Sep 17];10(1). Available from:
- 436 https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-018-0593-7

- 437 41. Eisenberger D, Tuschak C, Werner M, Bogdan C, Bollinger T, Hossain H, et al. Whole438 genome analysis of vancomycin-resistant Enterococcus faecium causing nosocomial
 439 outbreaks suggests the occurrence of few endemic clonal lineages in Bavaria, Germany.
 440 J Antimicrob Chemother. 2020 Jun 1;75(6):1398–404.
- 42. Timme RE, Rand H, Shumway M, Trees EK, Simmons M, Agarwala R, et al. Benchmark
 42. datasets for phylogenomic pipeline validation, applications for foodborne pathogen
 443 surveillance. PeerJ. 2017 Oct 6;5:e3893.
- 444 43. Lockhart SR, Etienne KA, Vallabhaneni S, Farooqi J, Chowdhary A, Govender NP, et al.
 445 Simultaneous Emergence of Multidrug-Resistant*Candida auris*on 3 Continents
 446 Confirmed by Whole-Genome Sequencing and Epidemiological Analyses. Clin Infect
 447 Dis. 2017 Jan 15;64(2):134–40.
- 448 45. Huson DH, Scornavacca C. Dendroscope 3: An Interactive Tool for Rooted Phylogenetic
 449 Trees and Networks. Syst Biol. 2012 Dec 1;61(6):1061–7.
- 46. Smith MR. Information theoretic generalized Robinson–Foulds metrics for comparing
 phylogenetic trees. 2020;36(20):5007–13.
- 47. Harris PNA, Ben Zakour NL, Roberts LW, Wailan AM, Zowawi HM, Tambyah PA, et al.
 Whole genome analysis of cephalosporin-resistant Escherichia coli from bloodstream infections in Australia, New Zealand and Singapore: high prevalence of CMY-2 producers and ST131 carrying blaCTX-M-15 and blaCTX-M-27. J Antimicrob

456 Chemother. 2018 Mar 1;73(3):634–42.







••• (>			latforn	n.solu	igeno	mics.com			C			ů + ©
solu-publication-2		ERR3569619	Enterococcus faecium	80	~	15	VRE	~	8	~	Not computed	ĸv	2 months ago
☆ Samples	~	ERR3569664	Enterococcus faecium	117	~	33	VRE	~	11	~	Not computed	ĸv	2 months ago
All samples E. faecium	99	ERR3569663	Enterococcus faecium	117	~	35	VRE	~	10	Ŷ	Not computed	ĸv	2 months ago
S. entericaS. aureus	23 35	ERR3569662	Enterococcus faecium		~	43	VRE	~	11	~	Not computed	ĸv	2 months ago
Olusters	>	ERR3569661	Enterococcus faecium	80	×	15	VRE	~	10	~	Not computed	ĸv	2 months ago
		ERR3569660	Enterococcus faecium	117	×	14	VRE	~	8	~	Not computed	ĸv	2 months ago
		ERR3569659	Enterococcus faecium	80	~	15	VRE	~	9	×	Not computed	ĸv	2 months ago
		ERR3569658	Enterococcus faecium	80	~	15	VRE	~	8	~	Not computed	ĸv	2 months ago
		ERR3569657	Enterococcus faecium	117	~	35	VRE	~	9	~	Not computed	ĸv	2 months ago
		ERR3569656	Enterococcus faecium	117	~	19	VRE	~	10	~	Not computed	ĸv	2 months ago
		ERR3569655	Enterococcus faecium	80	~	35	VRE	~	11	\sim	Not computed	ĸv	2 months ago
		ERR3569653	Enterococcus faecium	78	~	16	VRE	~	9	~	Not computed	ĸv	2 months ago
餟 Help & feedback		ERR3569652	Enterococcus faecium	117	×	19	VRE	~	10	~	Not computed	KV	2 months ago
Sign up for your own Solu workspace!		ERR3569651	Enterococcus faecium	117	~	35	VRE	~	10	~	Not computed	ĸv	2 months ago



- SRR498369

SRR498373 SRR500494





SNP cluster (reference)

ST	(reference)	
----	-------------	--

ST-	ST612
ST117	ST78
ST552	ST80

CT (reference)



	Solu	Reference	
	CFSAN000212	CFSAN000212	_
	CFSAN000211	CFSAN000211	_
	CFSAN000191	CFSAN000191	_
	CFSAN000228	CFSAN000228	_
	CFSAN000189	CFSAN000189	_
	CFSAN001140	CFSAN001140	_
	CFSAN000952	CFSAN000952	_
	CFSAN000951	CFSAN000951	_
	CFSAN000753	CFSAN000753	_
	CFSAN000752	CFSAN000752	_
	CFSAN000700	CFSAN000700	_
	CFSAN000661	CFSAN000661	_
	CFSAN000669	CFSAN000669	_
	CFSAN000954	CFSAN000954	_
	CFSAN000968	CFSAN000968	_
	CFSAN000963	CFSAN000963	
	CFSAN000958	CFSAN000958	-
	CFSAN000960	CFSAN000960	_
	CFSAN000970	CFSAN000970	_
	CFSAN001112	CFSAN001112	_
	CFSAN001115	CFSAN001115	_
	CFSAN001118	CFSAN001118	_
	CFSAN000961	CFSAN000961	_





