

1 **Title**

2 BaGPipe: an automated, reproducible, and flexible pipeline for bacterial genome-
3 wide association studies

4

5 **Authors**

6 Kuangyi Charles Wei 1,2, Beth Blane 1,3, Jacqueline Toussaint 4, Sandra Reuter 5,
7 Michelle S. Toleman 6, Mili Estee Torok 3, Sharon J. Peacock 3, Ewan M Harrison
8 1,3, Dinesh Aggarwal* 1,3,7, William Roberts-Sengier* 1

9 1. Parasites and Microbes Programme, Wellcome Sanger Institute, Hinxton, UK

10 2. Department of Genetics, University of Cambridge, Cambridge, UK

11 3. Department of Medicine, University of Cambridge, Cambridge, UK

12 4. European Bioinformatics Institute, Cambridge, UK

13 5. University of Freiburg, Freiburg, Germany

14 6. Cambridge University Hospitals NHS Foundation Trust, Cambridge, UK

15 7. Department of Infectious Diseases, Imperial College London, London, UK

16

17 Corresponding authors:

18 Kuangyi (Charles) Wei, kw524@cam.ac.uk

19 William Roberts-Sengier, wr7@sanger.ac.uk

20

21 * These authors contributed equally

22

23

24

25

26 **Abstract**

27 *Background*

28 Microbial genome-wide association studies (GWAS) are crucial for linking genetic
29 variation to phenotypic traits in bacteria. However, current tools often involve complex
30 manual processing, limited scalability, and fragmented workflows, which constrain
31 large-scale or routine bacterial GWAS.

32 *Results*

33 We developed BaGPipe, an automated and flexible bacterial GWAS pipeline built
34 using Nextflow and incorporating Pyseer for association analysis. BaGPipe integrates
35 pre-processing, statistical analysis, and downstream visualisation into a unified
36 workflow that is reproducible and easy to deploy across diverse computational
37 environments. BaGPipe was validated on a publicly available dataset of *Streptococcus*
38 *pneumoniae* whole-genome sequences, and reproduced published findings with
39 improved computational efficiency. BaGPipe was then applied to a dataset of
40 *Staphylococcus aureus* whole-genome sequences, successfully identifying known
41 and novel antibiotic resistance associations.

42 *Conclusions*

43 By offering an accessible, efficient, and reproducible platform, BaGPipe accelerates
44 bacterial GWAS and facilitates deeper exploration into the genetic underpinnings of
45 phenotypic traits. BaGPipe is freely available at [https://github.com/sanger-](https://github.com/sanger-pathogens/BaGPipe)
46 [pathogens/BaGPipe](https://github.com/sanger-pathogens/BaGPipe).

47

48 **Keyword**

49 microbial GWAS, Nextflow pipeline, Pyseer, whole genome sequencing, genome-
50 wide association study

51

52 **Background**

53 Genome-wide association studies (GWAS) have become an essential tool for
54 uncovering the genetic basis of phenotypic traits across diverse organisms,
55 particularly for investigating mechanisms of pathogenicity and resistance [1–5].
56 However, bacterial GWAS faces unique challenges due to horizontal gene transfer,
57 complex population structures, and extensive genomic diversity, which complicate
58 accurate genotype-phenotype associations [6–9]. Current bacterial GWAS tools often
59 exacerbate these difficulties by requiring fragmented workflows, extensive manual
60 intervention, and lacking scalability and reproducibility.

61

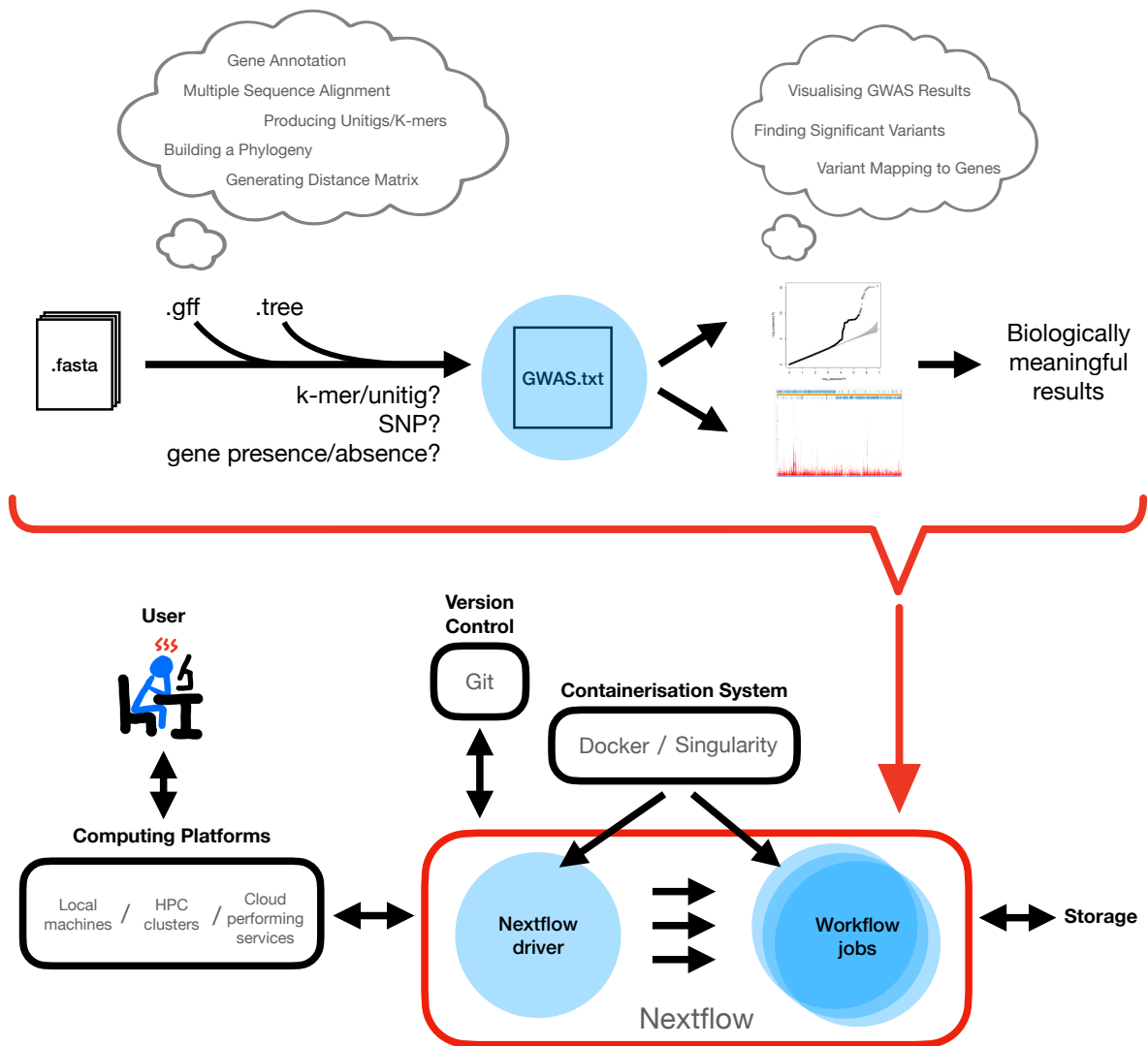
62 Executing bacterial GWAS requires analysing large genomic datasets to identify
63 associations with phenotypes such as antimicrobial resistance (AMR), while
64 accounting for factors like population structure and statistical power [10]. Current
65 bacterial GWAS tools fall into three categories [9, 11, 12]: phylogenetic, non-
66 phylogenetic, and machine learning approaches (See **Supplementary Table 1**).
67 Phylogenetic tools like Scoary [13] and TreeWAS [14] rely on well-defined
68 phylogenetic structures and are suitable for datasets where recombination can be
69 mitigated but are less practical for highly diverse species or large datasets. Non-
70 phylogenetic tools such as Bugwas [15] and SEER [16] adapt concepts from human
71 GWAS to develop phylogeny-independent methods. While computationally efficient
72 and scalable, they face challenges like elevated false-positive rates due to insufficient
73 correction for population structure. Machine learning GWASs tools are rarely used,
74 and no benchmarking has been performed.

75

76 Among non-phylogenetic tools, Pyseer [17] effectively addresses population structure
77 and mitigates false positives by incorporating a linear mixed model [18]. By leveraging
78 dimensionality reduction techniques like multidimensional scaling (MDS) and
79 employing k-mer based association studies, Pyseer has emerged as a standard tool
80 in the field due to its robustness and speed, especially in identifying genetic variants
81 with small effect sizes [12, 19–22].

82

83 Despite these advancements, executing bacterial GWAS analyses remains hindered
84 by significant technical challenges. Pyseer and similar tools [23–25] lack integrated
85 solutions for crucial pre-processing steps — such as genome annotation, phylogenetic
86 inference, or distance matrix generation — and downstream analyses like result
87 visualisation and annotation (**Figure 1**). These steps demand expertise across
88 multiple bioinformatics tools, each with its own dependencies and idiosyncrasies,
89 leading to an increased risk of error, wasted computational resources, and ultimately,
90 a diversion from biological interpretation towards technical troubleshooting. Existing
91 pipelines, such as bacterialGWAS [26], DBGWAS [23] and GEMMA kmer_pipeline
92 [27], attempt to address specific aspects of bacterial GWAS but remain limited by their
93 narrow scope, fragmented workflows, or lack of scalability. A recent Snakemake
94 pipeline, microGWAS [28], offers a more comprehensive wrapper around Pyseer, but
95 it does not fully automate pre- and post-processing steps or offer the modular
96 scalability enabled by the popular modern workflow manager, Nextflow [29].



97

98 **Figure 1: Schematic overview of the problem addressed and the solution by BaGPipe.**

99 Top: A comprehensive bacterial GWAS workflow involves miscellaneous pre-processing steps
 100 and various post-processing analyses, yet most bacterial GWAS tools only achieve the
 101 association study (shown by a blue circle). Bottom: The objective is to compile all relevant
 102 tools into a Nextflow pipeline, which can be easily deployed by users across computing
 103 platforms. The quantile-quantile plot (Q-Q plot) and the Manhattan plot are taken from the
 104 Pyseer tutorial [30].

105

106 To address these unmet needs, we present BaGPipe, a highly scalable and
 107 reproducible pipeline built on Nextflow [29] (**Figure 1**). BaGPipe fully leverages Pyseer
 108 for bacterial GWAS analysis while encompassing the entire GWAS workflow. We

109 validate BaGPipe by replicating previously published findings in a well-characterised
110 *Streptococcus pneumoniae* dataset [30–32] and further examine its utility by exploring
111 novel genetic associations in a clinically relevant *Staphylococcus aureus* genomic
112 dataset [33]. We then benchmark BaGPipe with AURORA [34], a Machine Learning
113 tool, and microGWAS [28], a Snakemake pipeline, on the same *S. pneumoniae*
114 dataset to compare their outputs and computational efficiencies.

115

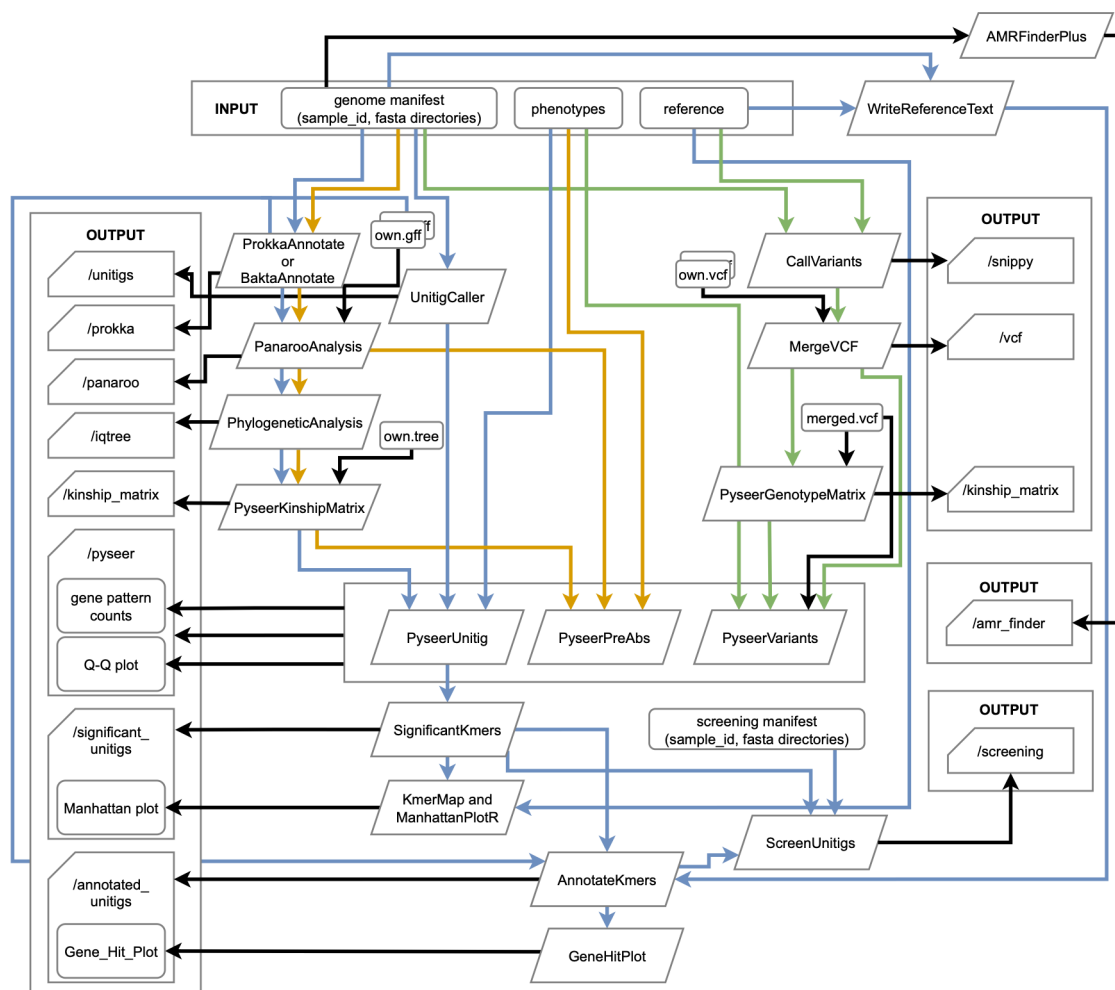
116 **Methods**

117 **Pipeline Overview**

118 BaGPipe is a bacterial GWAS pipeline designed to integrate a wide array of
119 bioinformatics tools into a unified workflow that can efficiently handle pre-processing
120 steps to post-processing analysis and visualisation. The pipeline was developed to
121 address computational challenges and facilitate genome-wide association studies on
122 bacterial datasets in a reproducible and user-friendly manner. The pipeline utilises
123 Nextflow [29] for workflow management, ensuring compatibility across diverse
124 computational environments while optimising resource usage. A step-by-step tutorial
125 and comprehensive documentation are available online ([https://github.com/sanger-](https://github.com/sanger-pathogens/BaGPipe)
126 [pathogens/BaGPipe](https://github.com/sanger-pathogens/BaGPipe)).

127

128 BaGPipe supports different entry points to provide maximum flexibility for users
129 (**Figure 2**), allowing for the integration of their own input data, such as phylogenetic
130 trees or variant call formats (VCFs). This feature allows BaGPipe to accommodate
131 diverse user requirements, reducing the burden of having to repeat computationally
132 intensive pre-processing tasks.



133

134 **Figure 2: Overview of the BaGPipe pipeline.** Input files are displayed on the top and output
 135 files on the sides. Alternative input files (GFF file, phylogenetic tree file, or VCF file) can be
 136 provided using appropriate option parameters. In a comprehensive analysis using the unitig
 137 mode, there will be eight output folders. The colour defines the genotype mode of analysis
 138 (blue: unitig/k-mer; yellow: gene presence and absence; green: short variations, including
 139 SNPs and Indels).

140

141 **Pipeline Implementation and Reproducibility**

142 To ensure ease of use and reproducibility, BaGPipe was developed using Docker and
 143 Nextflow [29]. Docker standardises the software environment, encapsulating all
 144 essential tools and specific versions into a single container for consistency across

145 platforms. Nextflow manages the modular pipeline, allowing users to execute the
146 workflow with a single command while specifying input files and analysis settings:

```
147     # install pipeline
148     nextflow pull sanger-pathogens/bagpipe
149     # run pipeline
150     nextflow run sanger-pathogens/bagpipe --help
```

151

152 For high-performance computing (HPC) environments, BaGPipe supports the Load
153 Sharing Facility (LSF) executor, enabling efficient parallel processing by dispatching
154 each task as an independent job. It dynamically allocates memory and CPU based on
155 task needs and can escalate resource requests if a job fails due to underestimation.

156

157 **User Accessibility and Configuration**

158 The pipeline includes a default configuration file with all the key options and
159 parameters necessary for running BaGPipe. This configuration is easily modifiable:
160 non-experienced users only need to adjust key parameters, while advanced users
161 can modify any parameter by editing the source code or profile configuration. As
162 BaGPipe is implemented in Nextflow, it can be deployed on a wide range of HPC
163 infrastructures and schedulers, including LSF and SLURM, by supplying an
164 appropriate institutional profile. This design ensures that BaGPipe remains both
165 accessible and highly adaptable to diverse compute environments and research
166 workflows.

167

168 **Pre-processing**

169 In its most comprehensive mode, BaGPipe starts with genome assemblies and
170 automates pre-processing, including the generation of k-mers or unitigs, genome
171 annotation, pangenome analysis, and phylogenetic analysis. This pre-processing
172 stage is supported by Prokka (v.1.14.6) [35] and Bakta (v.1.11.3) [36] for annotation,
173 Panaroo (v.1.5.1) [37] for constructing pangenomes and creating a multiple sequence
174 alignment of core genes, and IQ-TREE (v.2.2.6) [38] for generating a core genome
175 phylogeny. The output from these steps provides the essential elements for
176 subsequent analysis, such as annotated assemblies, a core phylogenetic tree, and a
177 distance matrix that reflects the genetic relationships within the dataset.

178

179 **Association Analysis**

180 BaGPipe employs Pyseer (v.1.3.11) [17], a Python implementation of SEER [16], to
181 conduct the association study, leveraging a linear mixed model to identify genotype-
182 phenotype relationships. The typical bacterial GWAS requires three core types of
183 inputs: genotypes, phenotypes, and an interaction or population structure correction.
184 In BaGPipe, users can conduct association analyses using different genotyping
185 approaches, including unitigs, k-mers, SNPs and indels, which provide flexibility in how
186 genomic variation is represented.

187

188 The linear mixed model in Pyseer controls for population structure using kinship
189 matrices derived from phylogenetic analysis. Alternatively, a pairwise distance matrix
190 can be generated from genome assemblies using Mash [39], which is then employed
191 to manage population stratification. BaGPipe supports Pyseer's k-mer/unitig-based
192 association study, which is recognised as a best practice for bacterial GWAS due to
193 its capacity to handle high genetic diversity while controlling for population structure

194 [30]. This k-mer/unitig approach aids in controlling population structure by providing a
195 comprehensive genetic representation of each isolate, capturing all types of genetic
196 variations — including SNPs, indels, and structural variants. By using k-mers/unitigs
197 to quantify genetic similarities among isolates, a more accurate kinship matrix can be
198 constructed, reflecting the true genetic relationships based on the presence or
199 absence of k-mers across genomes.

200

201 Q-Q plots are generated to visualise the p -value distribution against the null hypothesis,
202 ensuring that significant associations are not artefacts of confounding factors, such as
203 population structure, phylogenetic relatedness, horizontal gene transfer, or technical
204 biases. Additionally, as a good practice inherited from using Pyseer, BaGPipe provides
205 a count of genotype patterns to guide users on whether corrections for multiple testing
206 are necessary, if k-mers or unitigs are used as the genotyping approach.

207

208 **Post-processing and Visualisation**

209 BaGPipe facilitates intuitive analysis by automatically processing significant unitigs or
210 other genomic elements identified during the association analysis. It generates
211 Manhattan plots to visualise the genomic locations of significant associations and
212 annotates these markers using both reference and draft assemblies. The output
213 includes a "gene_hits.tsv" file summarising significant unitigs and their gene
214 annotations, providing insights into potential biological functions.

215

216 BaGPipe enriches the final output by annotating genes in the neighbourhood of
217 significant markers, which aids in biological interpretation. Additionally, a default R
218 script (adapted from the documentation of Pyseer [30]) is provided to allow users to

219 create custom visualisations of the data, further simplifying downstream exploration
220 and interpretation.

221

222 **Screening of GWAS-Identified Loci on External Datasets**

223 BaGPipe includes an optional module for screening GWAS-identified k-mers or unitigs
224 in independent external cohorts. The module accepts as input the list of statistically
225 significant loci identified during the discovery GWAS and a user-provided manifest of
226 genome assemblies from the screening dataset. Each of these genomes is screened
227 for the presence or absence of these loci using a unitig-based matching approach.
228 The module produces summary outputs including per-sample counts of detected loci,
229 per-locus prevalence across the screening cohort, and optional gene-level summaries
230 when annotations from the discovery phase are available. This functionality enables
231 screening of further genomes for potential virulence or resistance traits without
232 additional manual processing.

233

234 **Optional Annotation of Antimicrobial Resistance Genes**

235 To allow for comparison with known antimicrobial resistance (AMR) genes and allow
236 users to explore novelty within their analyses, BaGPipe provides an optional functional
237 annotation module that integrates AMRFinderPlus for the identification of known AMR
238 genes from genome assemblies. When enabled, AMRFinderPlus (v.4.0.23) [40] is run
239 automatically on all samples, generating per-sample AMR profiles and a cohort-level
240 summary of detected resistance determinants.

241

242 **Implementation of BaGPipe on the *Streptococcus pneumoniae* dataset**

243 To assess BaGPipe's ability to identify resistance-associated variants, we applied it to
244 a well-characterised *S. pneumoniae* dataset ($n = 616$) on beta-lactam resistance. This
245 dataset was chosen as a benchmark due to its prior use in bacterial GWAS [30–32],
246 allowing for comparisons with existing methods like Pyseer. Different to the Pyseer
247 tutorial [30], BaGPipe was run with unitigs (maximal non-branching paths in a
248 compacted de Bruijn graph) instead of k-mers to reduce sequence redundancy while
249 maintaining the integrity of the analysis. For inputting to BaGPipe, we made the
250 required genome manifest CSV file indicating the paths to all assemblies. To conduct
251 significant unitig analysis, we also made a reference manifest CSV file. BaGPipe was
252 executed with the LSF executor.

253

254 **Implementation of BaGPipe on the *Staphylococcus aureus* dataset**

255 To evaluate BaGPipe's performance in a different genomic context, we applied it to a
256 methicillin-resistant *S. aureus* (MRSA) dataset [33] (metadata see **Supplementary**
257 **Table 2**), aiming to identify genetic determinants of antibiotic resistance in a real-world
258 genomic surveillance setting. To quality-control the assemblies, we ran the FASTQ
259 files ($n = 520$) through bacQC (v.1.2, <https://github.com/avantonder/bacQC>), which
260 trimmed sequences using fastp [41]. Inspection on the output QC file confirms a slight
261 improvement in the quality of the data, with fewer adapter sequences and PCR
262 duplicate artifacts. From the Kraken2 [42] and Bracken [43] reports, two obviously
263 contaminated samples were identified and excluded from the data for further analysis
264 (see **Supplementary Table 3**). Next, we forwarded the FASTQ files ($n = 518$) to
265 assembleBAC (v.1.2, <https://github.com/avantonder/assembleBAC>). It was observed
266 that the smallest N50 is 72.2 Kbp and all 518 assemblies have length 2.8-3.0 Mbp.
267 We discarded samples with less than 90% reads matching to *S. aureus* and those with

268 <30x coverage from onward analyses, as well as removing assemblies with an N50
269 value <10 Kbp, length of less than 2.6 Mbp or greater than 3.0 Mbp, or with a
270 spuriously high number of contigs. No other samples were excluded based on these
271 criteria (518 samples retained).

272

273 We prepared the appropriate manifest files and other input files and executed
274 BaGPipe using the LSF executor. Analyses were conducted to identify the genetic
275 basis of resistance to eight antibiotics: oxacillin, ciprofloxacin, erythromycin, fusidic
276 acid, clindamycin, tetracycline, gentamicin, and mupirocin. Antibiotic resistance in
277 the *S. aureus* dataset was recorded as binary. Details of isolate collection, whole-
278 genome sequencing and antimicrobial susceptibility testing are available in the
279 published study [33].

280

281 To compare our results, we ran AMRFinderPlus (v.3.11.18) [40] for prediction of
282 antimicrobial resistance genes from known databases. For each antibiotic, we
283 compiled AMRFinder results and extracted significant gene-level associations at FDR
284 ≤ 0.05 from BaGPipe runs. Gene symbols were harmonised by lower-casing, removing
285 parentheses, stripping variant suffixes (e.g., *parC_S80F* changed to *parC*). We
286 computed a set of agreement metrics between BaGPipe significant genes and
287 AMRFinder catalogue genes, including precision (fraction of BaGPipe genes that are
288 in the catalogue), recall (fraction of catalogue genes recovered by BaGPipe), Jaccard
289 (intersection over union), and F1 score (harmonic mean of the precision and recall).
290 AMRFinder per-gene counts across the 518 isolates were used to estimate prevalence
291 ($= \text{count} / 518$). BaGPipe's *avg_beta* was used as the effect size proxy (log-odds).

292

293 **Benchmarking BaGPipe against other GWAS tools on the on the *Streptococcus***
294 ***pneumoniae* dataset**

295 To benchmark BaGPipe against a machine learning approach, we applied AURORA
296 (v.1.0.0) [34] to the *S. pneumoniae* dataset. Pangenome presence/absence was
297 generated with Panaroo (v.1.5.1) [37] in Roary-compatible mode and the resulting
298 gene_presence_absence_roary.csv was used as the feature matrix for both AURORA
299 and cross-method mapping. For gene-level benchmarking, AURORA features
300 (group_####) were mapped to gene symbols using the Roary “Non-unique Gene
301 name” and “Annotation” fields; when multiple aliases were present, all were retained
302 and the maximum score per gene was used. Features that could not be mapped
303 beyond “group_####” were excluded to ensure comparability with BaGPipe gene
304 identifiers. A maximum likelihood phylogeny was inferred from the core genome
305 alignment and pruned to genomes with phenotype labels prior to AURORA analysis.
306 AURORA’s phenotype module was used to identify non-typical or mislabelled strains
307 via a bagged random forest with outlier detection, with per-bag proximities and
308 summary plots inspected to assess class balance and outliers. GWAS was performed
309 using the Panaroo Roary-format matrix (type_bin_mat = "panaroo") and the pruned
310 phylogeny, yielding per-feature performance against the resistant class. Gene-level
311 AURORA evidence was summarised by the maximum F1 score per gene (primary
312 metric), with additional metrics (precision, recall, standardised residuals) retained for
313 sensitivity analyses and figure annotation. Cross-method comparison used AURORA
314 F1 scores and BaGPipe gene hits as ranking metrics, with overlap assessed by set
315 intersection, Jaccard distance, and Spearman rank correlation within the overlap.

316

317 To benchmark BaGPipe against a Snakemake-based GWAS pipeline, we ran
318 microGWAS (v.0.6.2) [28] on the same *S. pneumoniae* dataset using Bakta (v.1.11.3)
319 [36] generated GFF annotations produced from BaGPipe's annotation step.
320 microGWAS was executed end-to-end with its unitig mode, producing per-feature and
321 gene-level outputs. To compare gene hits, unitig-level microGWAS results were
322 mapped to genes via a Panaroo (Roary-format) dictionary with stringent tokenisation,
323 excluding "bad-chisq" rows and unmapped "group_#####" labels, scoring genes by the
324 maximum $-\log_{10}(p\text{-value})$, ranking both methods by descending score, and quantifying
325 agreement via intersection size, Jaccard index, and Spearman rank correlation on the
326 overlap. As both pipelines were run on an HPC, Slurm accounting and workflow logs
327 were collected for all runs on the same HPC profile to record wall-time, peak memory
328 and disk usage. Furthermore, BaGPipe was run on one published *Escherichia coli*
329 virulence dataset ($n = 370$) [44] originally reported by the microGWAS authors [28]
330 and the results were compared.

331

332 **Statistical comparison of gene sets and rankings**

333 Set overlap was quantified using the Jaccard index, defined as the number of shared
334 genes divided by the number of genes present in either set. For comparisons with
335 the curated AMR catalogue, precision was the proportion of BaGPipe genes found in
336 the catalogue, recall was the proportion of catalogue genes recovered by BaGPipe,
337 and the F1 score was the harmonic mean of precision and recall. For
338 antibiotic-specific analyses, precision, recall, and F1 were computed per antibiotic;
339 precision and recall were summarised across antibiotics using the median and
340 range, and macro-averaged F1 was summarised as the mean and range of the
341 per-antibiotic F1 values. If a denominator was zero, the metric was treated as

342 missing and excluded from summaries; if a numerator was zero, the metric was set
343 to zero. Ranking concordance was assessed in two ways. First, top-list overlap was
344 computed by comparing the number of shared genes between the two top lists for
345 sizes from 10 to 100 in steps of 10; for comparability across sizes, a top-list Jaccard
346 was also reported using the same definition of shared over either list. Ties at the
347 boundary were included and truncated deterministically with a fixed seed to ensure
348 reproducibility. Second, Spearman rank correlation was computed on the
349 overlapping genes only, using average ranks for ties and a two-sided P value for the
350 null of no monotonic association. Where significance of set overlap was reported, it
351 was evaluated with a two-sided Fisher exact test on a two-by-two table of set
352 membership, with false discovery rate control by the Benjamini-Hochberg procedure.
353 All analyses were performed in R, with two-sided tests and fixed random seeds for
354 any tie-breaking.

355

356 **Results**

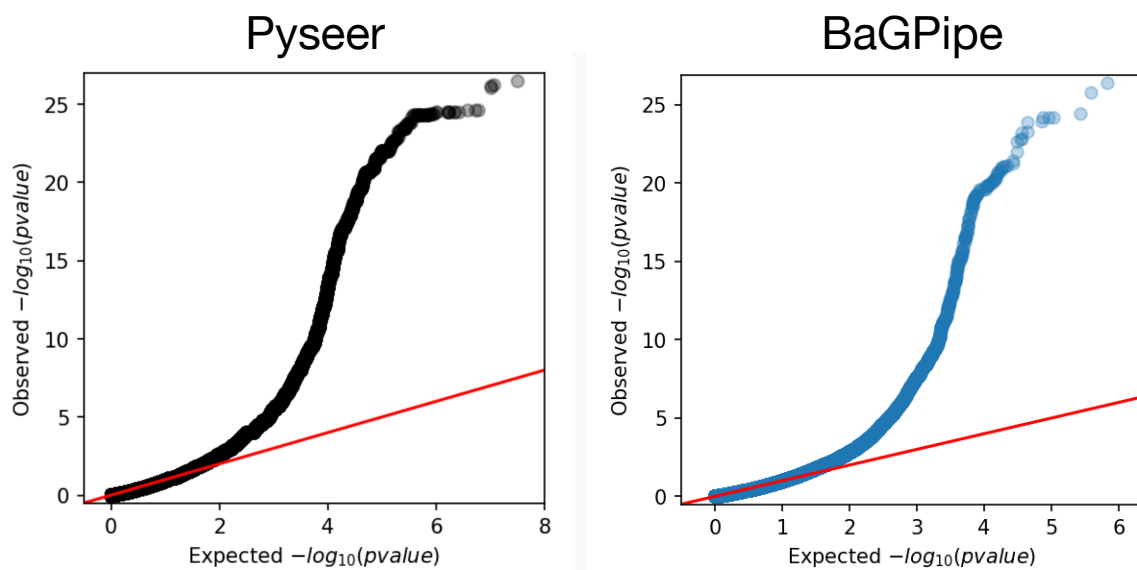
357 **Validation of BaGPipe by Reproducing Analyses from Pyseer**

358 We validated BaGPipe on 616 *S. pneumoniae* isolates, the dataset used in the Pyseer
359 tutorial [30]. This dataset has previously been used to identify genes associated with
360 penicillin resistance, providing a benchmark for bacterial GWAS. By analysing the
361 identical isolates, we aimed to assess BaGPipe's ability to replicate these findings and
362 evaluate its performance against established methods.

363

364 BaGPipe successfully identified 337,885 unique unitig patterns from a total of 740,945
365 unitigs tested. Recognising that many unitigs are highly correlated due to sequence
366 overlap, we determined a significance threshold of 1.48E-07 using a Bonferroni

367 correction based on the number of unique unitig patterns (0.05/337,885). This
368 approach aligns with the Pyseer tutorial, which suggests adjusting the significance
369 threshold based on the effective number of independent tests. We then applied this
370 threshold to filter significant unitigs — those with p -values below 1.48E-07 were
371 considered significant. By using this adjusted threshold, we effectively controlled for
372 multiple testing without being overly conservative, ensuring that the identified
373 associations are statistically robust. The Q-Q plot generated by BaGPIPE closely
374 resembled that produced in the Pyseer tutorial (**Figure 3**), demonstrating consistency
375 in managing population structure and ensuring the absence of poorly controlled
376 confounders.

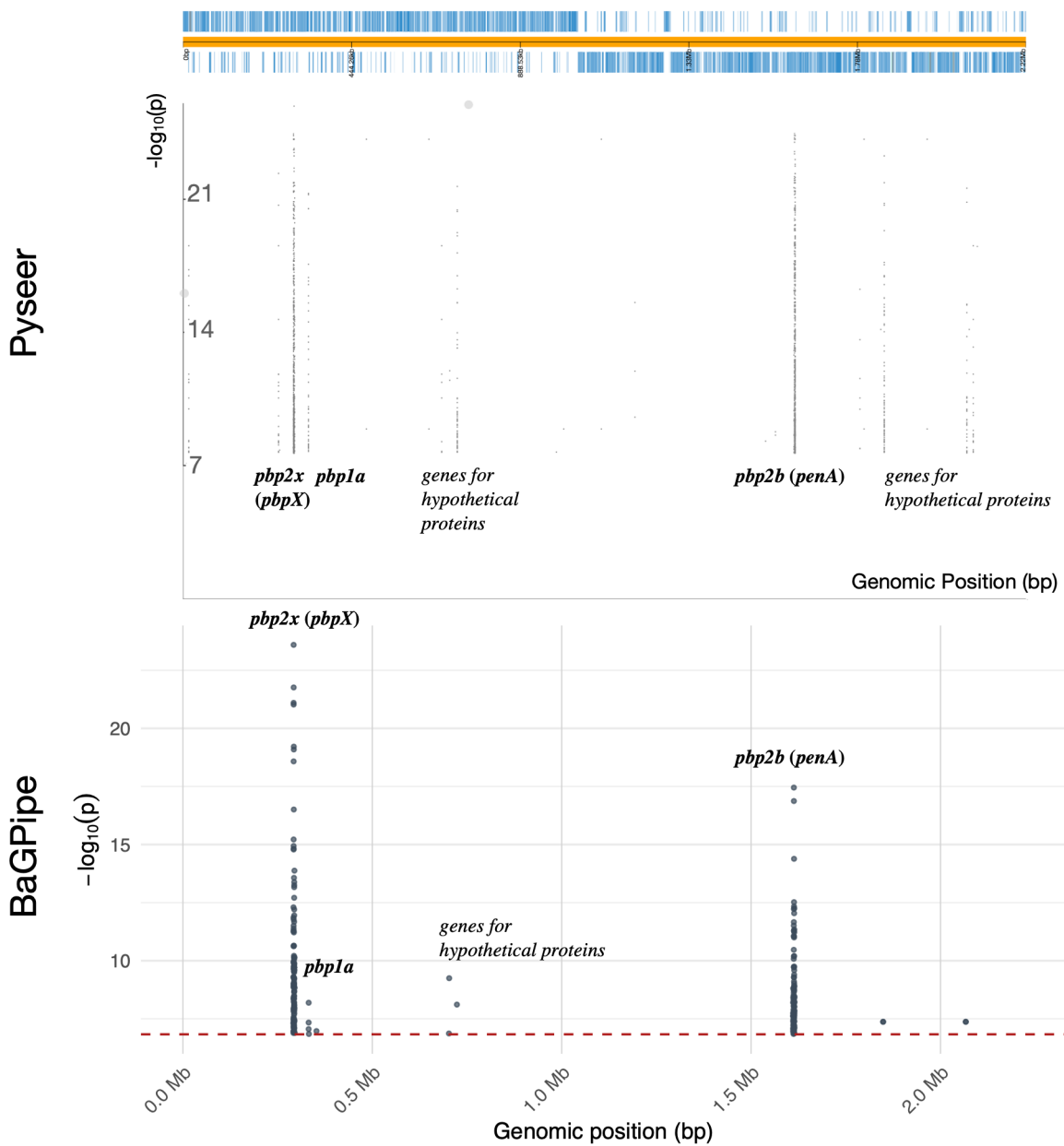


377

378 **Figure 3: Comparison of the Q-Q plots from BaGPIPE and from Pyseer on the**
379 ***Streptococcus pneumoniae* dataset.** The Q-Q plots show that observed $-\log_{10}(p$ -values)
380 are not inflated at low $-\log_{10}(p$ -values) and there is an absence of any poorly controlled
381 confounding population structure (they would appear as big “steps” deviating from the diagonal
382 line). The position of the points being above the null hypothesis (the diagonal line) indicates
383 significant k-mers/unitigs associated with penicillin resistance. The Q-Q plot produced from
384 Pyseer was sourced from the Pyseer tutorial [30].

385

386 The subsequent significant unitig analysis conducted by BaGPipe identified 573
 387 significant unitigs, in contrast to the 5,327 significant k-mers reported in the Pyseer
 388 tutorial. This notable reduction demonstrates the effectiveness of using unitigs to
 389 reduce sequence redundancy while retaining significant biological information. The
 390 significant unitigs were mapped to the reference genome, and a Manhattan plot was
 391 generated, showing two distinct peaks corresponding to the two penicillin-binding
 392 protein (pbp) genes: *pbp2x* and *pbp2b* (**Figure 4**). These are the hits found in the
 393 previous studies and they are known genes involved in penicillin resistance [30–32].



395 **Figure 4: Comparison of the Manhattan plots from BaGPipe and from Pyseer on the**
396 ***Streptococcus pneumoniae* dataset.** The plot is again consistent with that in the Pyseer
397 tutorial, showing two peaks, corresponding to the two *pbp* genes. The strongest peak
398 represents the locus coding for the putative penicillin binding protein 2x (*pbp2x*); the second
399 strongest peak represents *pbp2b* gene. The Manhattan plot produced from Pyseer, visualised
400 via Phandango [45], was sourced from the Pyseer tutorial [30].

401

402 Not all significant unitigs mapped to a single reference, underscoring the need for
403 multiple high-quality reference genomes to ensure comprehensive annotation and
404 reliable downstream analyses. The Manhattan plot generated by BaGPipe was
405 consistent with that from the Pyseer tutorial, although certain peaks observed in the
406 k-mer analysis were attenuated in the unitig analysis. Because unitigs compact
407 overlapping k-mers, the effective number of independent tests and the influence
408 of rare features differ between representations; in this dataset, this primarily narrows
409 signals to *pbp2x/pbp2b*, while neighbouring *recR/ddl* signals appear linkage
410 disequilibrium (LD)-consistent and parameter-sensitive.

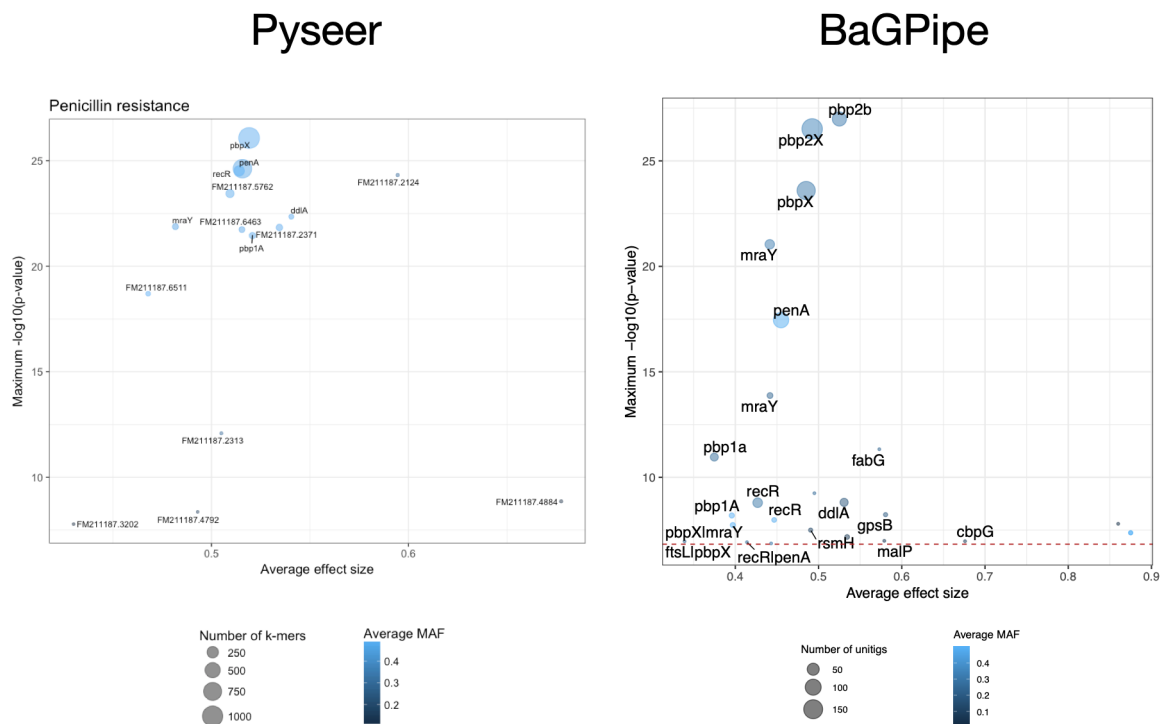
411

412 After annotating the significant unitigs, BaGPipe produced a list of "gene hits," which
413 included unitigs that either resided within or were adjacent to specific coding regions.
414 The gene-hit plot summarising these results closely matched that of the Pyseer tutorial
415 (**Figure 5**). Key hits such as *penA* (*pbp2b*) and *pbpX* (*pbp2x*) were consistently
416 identified across both analyses, demonstrating BaGPipe's reliability. Interestingly, the
417 *recR* and *ddl* gene, identified as significant in the Pyseer analysis, showed a lower
418 significance in the BaGPipe-derived results. *recR* encodes a protein involved in DNA
419 repair, specifically the RecFOR pathway for homologous recombination [46], which
420 may occasionally show spurious associations due to proximity to key resistance genes.
421 This discrepancy between using k-mers and unitigs suggests that differences

422 at *recR* reflect LD-driven, representation-dependent associations: prior work shows
 423 the *pbp2b-recR-ddl* block (*ddl* lies ~0.8 kb from *pbp2b* with *recR* between them)
 424 frequently recombines under β -lactam selection, so signals at *recR* can arise
 425 by hitchhiking with *pbp2b* [47]. Under a unitig representation and pattern-based
 426 thresholds these neighbouring signals are attenuated, while the *pbp2x/pbp2b* peaks
 427 remain.

428

429



430

431 **Figure 5: Comparison of the gene-hit plots from BaGPipe and from Pyseer on the *S.***
 432 ***pneumoniae* dataset.** The x-axis represents the average effect size, while the y-axis shows
 433 the maximum $-\log_{10}(p\text{-value})$, highlighting the statistical significance of gene associations. The
 434 size of the circular dots indicates the number of k-mers/unitigs involved in each association,
 435 providing insight into the genomic support for each hit. The colour scheme, ranging from lighter
 436 to darker shades of blue, illustrates the average minimum allele frequency (MAF), with darker
 437 shades indicating a lower MAF. BaGPipe accurately replicated the findings from Pyseer,
 438 confirming *penA* (*pbp2b*) and *pbpX* (*pbp2x*) as the most significant gene hits. The gene-hit
 439 plot produced from Pyseer was sourced from the Pyseer tutorial [30].

440

441 Reproducing this analysis from raw reads, BaGPipe used approximately 191 hours of
442 total CPU time, with a peak remote-access memory (RAM) usage of 3.5 GB, a mean
443 RAM usage of 574 MB, and an overall run time of about 23 hours and 11 minutes.
444 When excluding steps such as annotation, pangenome analysis, and phylogenetic
445 analysis — assuming these are provided as input — the CPU time would be reduced
446 to approximately 19 hours.

447

448 **Implementation of BaGPipe on a *S. aureus* Dataset**

449 BaGPipe was further evaluated using high-quality sequences from 518
450 *Staphylococcus aureus* samples tested against seven antibiotics [33]. It successfully
451 identified significant genetic associations with resistance phenotypes for multiple
452 antibiotics. To cross-validate these findings, AMRFinderPlus [40] was employed to
453 identify known AMR genes. BaGPipe's results were compared with AMRFinderPlus
454 predictions, confirming the presence of established AMR genes, such as *gyrA*, *parC*,
455 and *parE* for ciprofloxacin, and *ermA* for erythromycin (**Figure 6a**). Set overlap was
456 modest (**Figure 6b**), with a median Jaccard of 0.012 (range 0.005-0.045) reflecting
457 the disparity in set sizes and the broader scope of GWAS compared with a curated
458 catalogue. Precision (fraction of BaGPipe genes present in the catalogue) had
459 a median of 0.017 (range 0.0068-0.0735), whereas recall (fraction of catalogue
460 genes recovered by BaGPipe) had a median of 0.0426 (range 0.021-0.106). The
461 resulting macro-F1 averaged 0.030 across antibiotics (range 0.010-0.087).
462 Intersections comprised 1-5 genes per antibiotic (e.g., ciprofloxacin = 4;
463 erythromycin = 5), consistent with recovery of canonical determinants alongside

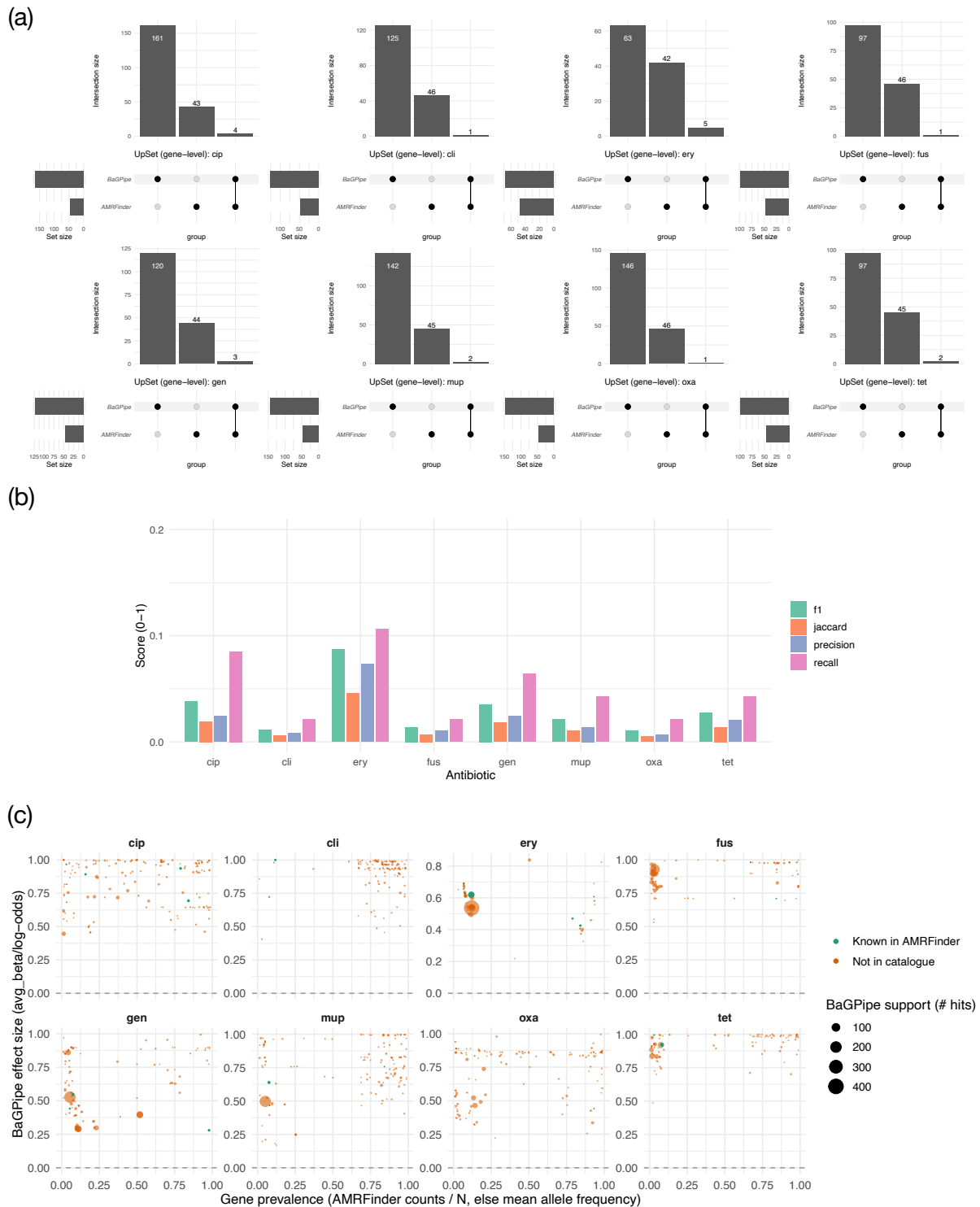
464 additional BaGPipe-only candidates, suggesting novel resistance mechanisms that
465 warrant further investigation.

466

467 Some genes have lower prevalence but exhibit moderate to strong effects (**Figure 6c**).
468 Upon inspection, some genes identified by BaGPipe are known to confer resistance
469 to other antibiotic classes, likely due to LD within a shared mobile genetic element
470 (MGE). This reflects the potential for cross-resistance or false-positive associations.
471 However, these observations could also indicate mechanisms of resistance that are
472 not yet fully understood, e.g. a regulatory gene. A detailed comparison of gene hits
473 predicted by AMRFinderPlus and those identified by BaGPipe is summarised in
474 **Table 1** (see also **Supplementary Table 4**).

475

476 The counts of resistant and susceptible phenotypes for all eight antibiotics, along with
477 Manhattan plots and gene hit plots produced from BaGPipe, are provided in the
478 **Supplementary Materials** and Zenodo (<https://doi.org/10.5281/zenodo.14947249>).



479

480 **Figure 6. Agreement between BaGPipe GWAS hits and AMRFinderPlus across eight**
 481 **antibiotics in *S. aureus*.** (a) UpSet plots display, for each antibiotic, the gene-level
 482 intersection and exclusive sets: BaGPipe-only (significant GWAS hits, $FDR \leq 0.05$),
 483 AMRFinder-only (catalogue genes), and the intersection. Gene symbols were harmonised as
 484 described in Methods. (b) Overlap metrics per antibiotic (Jaccard index, precision, recall, F1)
 485 computed on the gene sets in (a). (c) BaGPipe effect size (avg_beta/log-odds) versus gene

486 prevalence (AMRFinder counts / 518 isolates), coloured by AMRFinder status; filled symbols
 487 denote $FDR \leq 0.05$. AMRFinder counts provide prevalence for catalogue genes; prevalence
 488 is derived from mean allele frequency across tests contributing to the gene for BaGPipe-only
 489 genes. Point size encodes BaGPipe support (gene hits). cip, ciprofloxacin; cli, clindamycin;
 490 ery, erythromycin; fus, fusidic acid; gen, gentamicin; mup, mupirocin; oxa, oxacillin; tet,
 491 tetracycline.

492

493 **Table 1: Gene hits predicted from AMRFinderPlus and identified from bacterial GWAS**
 494 **by BaGPipe for each of the eight tested antibiotics.** For each antibiotic, the

495 corresponding class was identified, and AMRFinderPlus was used to screen for known AMR
 496 genes within that class in its database. For example, for ciprofloxacin (cip), a quinolone,
 497 genes like *gyrA*, *parC*, and *parE* were predicted from the sequences. There are some gene
 498 hits identified by BaGPipe that are not predicted from the known database; these hits are
 499 manifested as identifiers extracted from GFF files during annotation. Only the most common
 500 genes are shown here, and the full comparison table is **Supplementary Table 4**.

501 BaGPipe-only gene functions (where annotated) and concise literature notes on any
 502 plausible links to resistance mechanisms are provided in **Supplementary Table 7**. *: Not
 503 significant gene hit. cip, ciprofloxacin; cli, clindamycin; ery, erythromycin; fus, fusidic acid;
 504 gen, gentamicin; mup, mupirocin; oxa, oxacillin; tet, tetracycline.

505

Antibiotic	Class	Known AMR gene(s) predicted from AMRFinderPlus	Gene hit(s) identified from BaGPipe
cip	Quinolone	<i>gyrA</i>	<i>gyrA</i>
cip	Quinolone	<i>parC</i>	<i>parC</i>
cip	Quinolone	<i>parE</i>	Not found
cli	Lincosamide	<i>ermA</i>	<i>ermA</i> *
cli	Lincosamide	<i>ermC</i> , <i>Inu(A)</i>	Not found
ery	Macrolide	<i>ermA</i>	<i>ermA</i>
ery	Macrolide	<i>parC</i>	<i>parC</i> *
ery	Macrolide	<i>gyrA</i>	<i>gyrA</i> *
ery	Macrolide	<i>ermC</i> , <i>mph(C)</i> , <i>msr(A)</i>	Not found
ery	Macrolide	Not found	<i>GKPKGGHK_02596</i>

fus	Fusidic acid	<i>fusA, fusB, fusC</i>	<i>fus</i>
fus	Fusidic acid	<i>parC</i>	<i>parC*</i>
fus	Fusidic acid	Not found	<i>LFPMDBMC_02247, LFPMDBMC_02248, SA0059</i>
gen	Aminoglycoside	<i>aadD1</i>	<i>aadD*</i>
gen	Aminoglycoside	<i>qacA</i>	<i>qacA*</i>
gen	Aminoglycoside	<i>aac(6')-Ie/aph(2'')-Ia, ant(6)-Ia, ant(9)-Ia, aph(3')-IIIa</i>	Not found
mup	Mupirocin	<i>ileS</i>	<i>ileS_2, ileS*</i>
mup	Mupirocin	<i>parC</i>	<i>parC*</i>
mup	Mupirocin	<i>fusA, fusB, fusC</i>	<i>fus*</i>
mup	Mupirocin	<i>aadD1</i>	<i>aadD*</i>
mup	Mupirocin	<i>mupA</i>	Not found
oxa	Beta-lactam	<i>mecR1</i>	<i>mecR1*</i>
oxa	Beta-lactam	<i>blaZ, blaI, blaPC1, blaR1, mecA, mecl</i>	Not found
oxa	Beta-lactam	Not found	<i>fmtC, fruA</i>
tet	Tetracycline	<i>tet(K)</i>	<i>tet(K)</i>
tet	Tetracycline	<i>fusA, fusB, fusC</i>	<i>fus*</i>
tet	Tetracycline	<i>rpoB</i>	<i>rpoB*</i>
tet	Tetracycline	<i>tet(M), tet(38), tet(L)</i>	Not found
tet	Tetracycline	Not found	<i>cds-M013TW_04140, cds-M013TW_04125</i>

506

507 **Benchmarking BaGPipe against other GWAS tools**

508 To evaluate BaGPipe's computational performance and analytical capabilities, we
509 benchmarked it over the *S. pneumoniae* dataset against two state-of-the-art bacterial
510 GWAS tools: AURORA, a recent machine learning-based method [34], and
511 microGWAS, the latest Snakemake pipeline for bacterial genome-wide association

512 studies [28]. **Table 2** summarises the key differences between these approaches
 513 across multiple dimensions.

514 **Table 2: Comparison of bacterial GWAS tools**

515

Feature	BaGPipe	AURORA	microGWAS
Implementation	Nextflow pipeline	R package	Snakemake pipeline
Statistical Method	Linear Mixed Models (pyseer)	Machine learning: Random Forest, AdaBoost, Logistic Regression	Linear Mixed Models (pyseer)
Genotype Inputs	K-mers/Unitigs, Gene presence/absence, or Variants	Gene presence/absence, or k-mers	Unitigs, Gene presence/absence, Variants, Gene cluster k-mers, or Whole-genome combined (machine learning) unitigs
Required Input Files	Assemblies, Phenotypes	Annotated Genomes, Panaroo/Roary gene presence/absence matrix, phylogeny	Assemblies, Phenotypes, Annotations (GFF)
Annotation	Automatic (Bakta/Prokka), but user can skip this by providing GFF files.	External requirement (Panaroo/Roary input)	External requirement (GFF input)
Population Structure	Linear mixed models with kinship matrix (based on phylogenetic correction)	Phylogenetic correction and auto/allopatric strain identification	Linear mixed models with unitig-based kinship
Visualisation	Automatic (Q-Q, Manhattan, Gene-hit plots)	Externally manual	Automatic (Q-Q, Manhattan, Functional enrichment plots)
AMR Integration	Built-in optional module (AMRFinderPlus)	None	Built-in optional module (abritamr)
Screening Significant	Built-in optional module	None	None

Genes on External Dataset			
Containerisation	Docker/Singularity	None	Conda environments
Level of Automation	End-to-end	GWAS only	Near end-to-end. (requires pre-annotation)

516

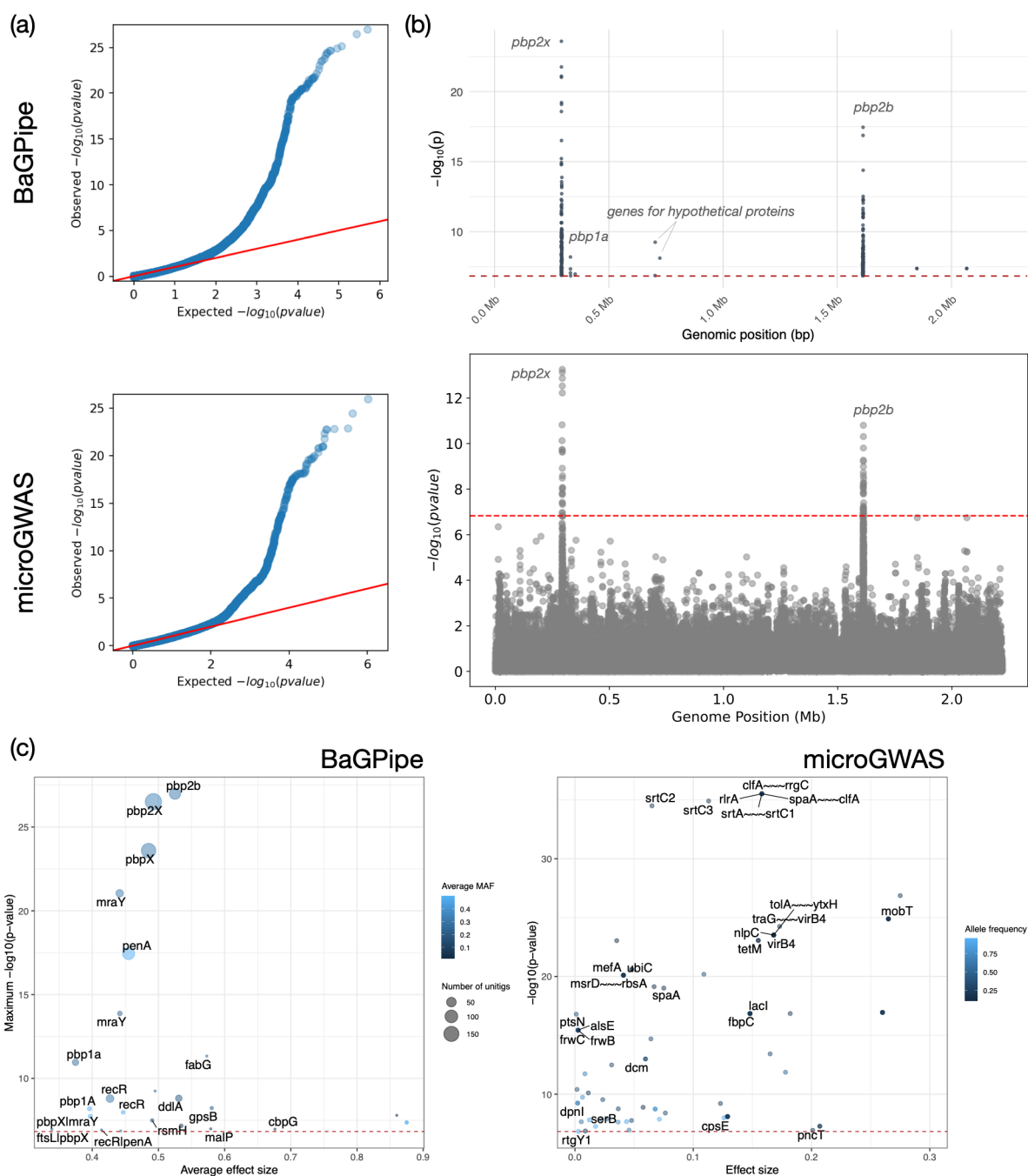
517 After mapping AURORA features to gene symbols via the Panaroo Roary file and
518 restricting to named genes, AURORA yielded 2,801 distinct genes (**Supplementary**
519 **Table 8**) and BaGPipe produced 20 genes. The intersection comprised 10 genes,
520 corresponding to a Jaccard index of 0.0036. Top-K overlap was 0 at all K values
521 from 10 to 100, reflecting the fact that the very top of each list is dominated by
522 method-specific features. Spearman rank correlation across the overlapping set
523 was $\rho = 0.1405$ ($p = 0.6986$), indicating weak pairwise association (**Supplementary**
524 **Table 9; Supplementary Fig. 10**). Despite the modest global overlap, the three
525 canonical β -lactam targets were detected by both methods: *pbp2x*, *pbp2b*,
526 and *pbp1a*. In the BaGPipe results these genes had hits of 178, 74,
527 and 16 respectively; in AURORA they achieved R-class F1 ≈ 0.667 for all three
528 (**Supplementary Table 9**). Additional overlaps
529 included *mraY*, *recR*, *gpsB*, *rsmH*, *fabG*, *ftsL*, and *malP*, consistent with broader
530 cell-wall and replication/repair associations in the cohort.

531

532 In comparison, microGWAS yielded 1,407 distinct genes, whereas the intersection
533 with BaGPipe's results comprised 2 genes, corresponding to a Jaccard index of
534 0.0014 (**Supplementary Table 11**). The two overlapping genes were *pbp2x* and
535 *fabG*, however, it was clear that microGWAS also picked up *pbp2b* from the
536 Manhattan plot (**Figure 7**), but the association might be lost in the downstream gene

537 mapping. In BaGPipe, *pbp2x* had the highest rank (hits = 178; rank = 1) whereas
 538 *fabG* was mid-list (rank = 14), but both appeared deep in the microGWAS ranking
 539 (microGWAS ranks = 678 for *pbp2x* and 947 for *fabG*). Interestingly, the third
 540 highest signal identified from BaGPipe is the *pbp1a* gene, which was not identified
 541 by microGWAS (**Figure 7**).

542



543

544 **Figure 7: Comparison of the results from BaGPipe and from microGWAS on the *S.***
545 ***pneumoniae* dataset.** (a) Q-Q plots, showing that observed $-\log_{10}(p\text{-values})$ are not inflated
546 at low $-\log_{10}(p\text{-values})$ and there is an absence of any poorly controlled confounding
547 population structure (they would appear as big “steps” deviating from the diagonal line). The
548 position of the points being above the null hypothesis (the diagonal line) indicates significant
549 k-mers/unitigs associated with penicillin resistance. (b) Manhattan plots. In both plots, the
550 strongest peak represents the locus coding for the putative penicillin binding protein 2x (*pbp2x*
551 or *pbpX*); the second strongest peak represents *pbp2b* (or *penA*) gene. The red dashed lines
552 are the p -value thresholds used by the pipelines to determine an association significant. (c)
553 Gene-hit plots. The x-axis represents the average effect size (or beta), while the y-axis shows
554 the $-\log_{10}(p\text{-value})$, highlighting the statistical significance of gene associations. The size of
555 the circular dots on the BaGPipe’s plot indicates the number of unitigs involved in each
556 association, providing insight into the genomic support for each hit. This information was not
557 obtained from a default microGWAS run, hence lacking on the right. The colour scheme,
558 ranging from lighter to darker shades of blue, illustrates the average MAF for BaGPipe (or
559 allele frequency for microGWAS), with darker shades indicating a lower MAF. The red dashed
560 lines are the p -value thresholds used by the pipelines to determine an association significant.
561 Genes which names are only strain-specific locus tags are not labelled on the plots.

562

563 Taken together, our comparative analysis on the *S. pneumoniae* penicillin resistance
564 dataset revealed that while the three methods show relatively low global feature
565 overlap (reflecting their distinct statistical frameworks and filtering strategies, and
566 differences in significant gene annotation mapping), they consistently identify the
567 most significant well-established biological associations, the canonical β -lactam
568 resistance determinants. AURORA's machine learning framework excels at
569 predictive accuracy and capturing non-linear associations; microGWAS offers
570 flexibility through six complementary statistical approaches; and BaGPipe prioritises
571 interpretability and automation for gene-level biological discovery.

572

573 From a computational perspective, the successful main microGWAS run completed
574 with approximately 44 CPU-hours with peak memory usage of 7.2 GB on this

575 moderately-sized dataset (616 strains). Resource requirements are slightly higher
576 than BaGPipe (without genome annotation). The analysis generated 3.0 GB of
577 output data with 8.0 GB for intermediate Snakemake files. Both tools demonstrated
578 resource profiles typical of multi-step bacterial genomics pipelines on HPC
579 environments. To get the workflow running for new users, both pipelines required
580 some debugging attempts, though BaGPipe's modular Nextflow design with
581 extensive caching drastically reduces re-computation costs for parameter
582 exploration, as does microGWAS's Snakemake rule-based dependency tracking.

583

584 To further validate BaGPipe's performance, we tested it on a published *Escherichia*
585 *coli* virulence dataset [44] originally reported by the microGWAS authors [28].

586 BaGPipe successfully identified *irp1*, *irp2*, and *ybtU* as significantly associated with
587 the phenotype. These genes comprise the yersiniabactin biosynthesis cluster
588 (located at approximately 2 Mb), a well-characterized pathogenicity island associated
589 with enhanced virulence in extraintestinal pathogenic *E. coli*. BaGPipe also detected
590 significant associations at location 1.15 Mb, corresponding to the high-pathogenicity
591 island (HPI), in agreement with microGWAS findings. These results demonstrate
592 BaGPipe's ability to recapitulate known biology and detect established genetic
593 determinants across different bacterial species and phenotypes (Results included in
594 **Supplementary Materials 12-15**).

595

596 **Discussion**

597 BaGPipe addresses the challenges associated with analysing large-scale bacterial
598 datasets by integrating a cutting-edge workflow management tool, Nextflow [29], along
599 with a modular design that allows researchers to efficiently conduct end-to-end

600 analyses. This pipeline offers an accessible solution to a major bottleneck in bacterial
601 genomics: performing reproducible and efficient GWAS on large datasets with diverse
602 computational requirements.

603

604 Validation of BaGPipe through replication of a published *S. pneumoniae* GWAS study
605 [30–32] demonstrated optimised computational performance. Further application to a
606 *S. aureus* dataset illustrated its versatility across different microbial contexts, with
607 cross-validation using AMRFinderPlus [40] identifying both established AMR genes
608 and previously unreported genetic elements. While some of these loci are not
609 confirmed resistance genes, they represent promising leads for further experimental
610 investigation, as these associations do not necessarily indicate false positives. For
611 instance, mutations in *rpoB*, initially linked to rifaximin resistance, were recently shown
612 to confer daptomycin resistance in *Enterococcus faecium* [19]. These findings highlight
613 both BaGPipe’s sensitivity and robustness in hypothesis generation for discovering
614 unexplored genetic mechanisms.

615

616 Our comprehensive benchmarking of BaGPipe against two state-of-the-art bacterial
617 GWAS tools, AURORA and microGWAS, shows that BaGPipe offers a distinctive,
618 complementary approach to bacterial association analysis. Rather than a single
619 optimal method, the field benefits from multiple specialised tools suited to different
620 research aims, computational constraints, and levels of user experience. BaGPipe’s
621 strengths lie in its modularity, flexibility, and automated resource management, which
622 streamline analyses without manual adjustments. When isolating the association
623 step, BaGPipe finished in under eight minutes of CPU time while producing all
624 standard outputs, whereas the machine learning based AURORA required ~58

625 minutes of wall-time and additional downstream mapping from pangenome clusters
626 to gene symbols. Biologically, BaGPipe prioritised the canonical penicillin-resistance
627 determinants in *S. pneumoniae* (*pbp2x*, *pbp2b*, *pbp1a*), consistent with established
628 literature and the pyseer tutorial, enabling rapid interpretation. Although AURORA
629 detected these loci, they were not ranked among the top features, and interpretation
630 was delayed by the need for manual feature annotation. microGWAS agreed with
631 BaGPipe in identifying the top penicillin-resistance genes (*pbp2x* and *pbp2b*), but did
632 not recover the weaker *pbp1a* signal detected by BaGPipe. microGWAS provides
633 greater flexibility through six complementary statistical approaches and may deliver
634 additional insights if allocated more resources. On this moderately-sized dataset
635 (616 strains), the main microGWAS run completed in ~44 CPU-hours with a peak
636 memory usage of 7.2 GB; resource requirements were slightly higher than BaGPipe
637 (without genome annotation). Excluding genome annotation to match the
638 microGWAS setup, BaGPipe required ~19 CPU-hours on the same dataset.
639 Uniquely, BaGPipe integrates genome annotation, AMR prediction, GWAS analysis,
640 and an external validation workflow within a single, reproducible framework that
641 requires only raw assemblies and phenotype data. This end-to-end automation
642 substantially reduces the bioinformatics burden, making BaGPipe particularly
643 valuable for research groups without dedicated computational expertise.

644

645 From a computational perspective, BaGPipe's implementation in Nextflow provides
646 significant practical advantages. Nextflow is one of the most widely adopted workflow
647 managers in genomics, ensuring immediate familiarity to the bioinformatics
648 community and access to extensive documentation, training resources, and
649 community support. The highly modular architecture of BaGPipe enables

650 researchers to customize analyses: users can provide their own phylogenetic trees
651 (critical for population structure correction), swap tools to leverage newer methods,
652 extend pipelines, or integrate BaGPipe with other established workflows in a
653 standardised manner. This modularity makes BaGPipe a robust long-term choice
654 that can evolve alongside methodological advances, new input formats, and
655 pathogen-specific requirements, essential for a field advancing as rapidly as
656 bacterial genomics.

657

658 BaGPipe is actively maintained by the Pathogen Informatics team at the Wellcome
659 Sanger Institute, ensuring regular updates, responsive issue tracking, and integration
660 of user feedback from the research community. This institutional commitment to
661 sustainability differentiates BaGPipe from many academic pipelines that lack long-
662 term support structures. The combination of comprehensive functionality,
663 computational efficiency, active maintenance, and integration with the broader
664 Nextflow ecosystem positions BaGPipe as an accessible and sustainable tool for
665 advancing bacterial genetic research, with the potential to become indispensable for
666 researchers seeking an integrated solution for microbial association studies. With
667 continued improvements and community-driven development, BaGPipe has the
668 potential to foster new discoveries and deepen our understanding of bacterial
669 genetics and pathogenesis.

670

671 **Conclusions**

672 BaGPipe is a powerful, flexible, and efficient pipeline for conducting bacterial GWAS
673 with the well-established Pyseer framework. We show it is capable of reliably handling
674 complex datasets, reducing redundancy, and replicating or uncovering both known

675 and novel genetic associations related to antibiotic resistance. BaGPipe lowers the
676 barriers to performing bacterial GWAS and offers a comprehensive, reproducible
677 framework, vital to better understand the genetic basis of AMR and bacterial traits.

678

679 **List of abbreviations**

680 **AMR:** Antimicrobial Resistance

681 **GWAS:** Genome-wide Association Study

682 **HPC:** High-performance Computing

683 **HPI:** High-pathogenicity Island

684 **Indel:** Insertions and Deletions

685 **LD:** Linkage Disequilibrium

686 **LSF:** Load Sharing Facility

687 **MAF:** Minimum Allele Frequency

688 **MDS:** Multidimensional Scaling

689 **MGE:** Mobile Genetic Element

690 **MRSA:** Methicillin-Resistant *Staphylococcus aureus*

691 **Q-Q plot:** Quantile-Quantile plot

692 **RAM:** Remote-Access Memory

693 **SNP:** Single Nucleotide Polymorphism

694 **VCF:** Variant Calling Format

695

696 **Declarations**

697 ***Ethics approval and consent to Participate:***

698 Not applicable

699 ***Consent to Publish declarations:***

700 Not applicable

701 **Clinical trial number:**

702 Not applicable

703 **Availability of data and materials:**

704 BaGPipe is freely available at <https://github.com/sanger-pathogens/BaGPipe>. The
705 *Streptococcus pneumoniae* input dataset is available from the Pyseer tutorial
706 (<https://pyseer.readthedocs.io/en/master/tutorial.html#>). The *Staphylococcus aureus*
707 sequencing assemblies can be sourced from their ERS accession numbers provided
708 in supplementary data. The reference assemblies, listed in the supplementary, can be
709 sourced from NCBI. Supplementary data are available at
710 <https://doi.org/10.5281/zenodo.14947249>.

711 *Software availability and requirements*

712 **Project name:** BaGPipe

713 **Project home page:** <https://github.com/sanger-pathogens/BaGPipe>

714 **Operating system:** Platform independent

715 **Programming language:** Nextflow, R, Python, shell

716 **Other requirements:** Java 11 or later (required by Nextflow), Nextflow

717 **License:** MIT licence

718 **Any restrictions to use by non-academics:** None

719 **Competing interests**

720 The authors declare that they have no competing interests.

721 **Funding**

722 This work was supported by Wellcome Grant reference: 220540/Z/20/A, 'Wellcome
723 Sanger Institute Quinquennial Review 2021-2026' – core funding of Wellcome Sanger
724 Institute. J.T. was supported by the European Molecular Biology Laboratory. S.R. is

725 funded by the German Ministry of Education and Research (BMBF) through grant
726 01KI2018. E.E. was supported by a Clinician Scientist Fellowship funded by the
727 Academy of Medical Sciences and the Health Foundation. D.A. has been supported
728 by the Wellcome Trust (Grant number: 222903/Z/21/Z). DA, Clinical Lecturer, is
729 funded by Health Education England (HEE) / NIHR for this research project. The views
730 expressed in this publication are those of the author(s) and not necessarily those of
731 the NIHR, NHS or the UK Department of Health and Social Care.

732 ***Authors' contributions***

733 D.A., W.R.S., E.M.H. designed the study; K.C.W. built the pipeline with input from
734 W.R.S. and undertook the bioinformatic analyses with contribution from D.A. and
735 W.R.S.; W.R.S., D.A., and E.M.H. supervised the study; J.T. provided intellectual input
736 during the conceptual phase of the study; B.B. curated and provided the *S. aureus*
737 dataset; S.R., M.S.T., M.E.T., S.J.P. contributed significantly to data generation and
738 analysis of the original study which produced the *S. aureus* dataset; K.C.W. wrote the
739 first draft of the manuscript; D.A., W.R.S., E.M.H., K.C.W. edited the manuscript. All
740 authors had access to the data and read, contributed and approved the final
741 manuscript.

742 ***Acknowledgements***

743 We would like to appreciate John Lees and his team for making Pyseer and its tutorial
744 available. We would also like to appreciate Francesc Coll for their encouragement and
745 feedback.

746

747 **References**

- 748 1. **Xue A, Wu Y, Zhu Z, Zhang F, Kemper KE, et al.** Genome-wide association
749 analyses identify 143 risk variants and putative regulatory mechanisms for type 2
750 diabetes. *Nat Commun* 2018;9:2941.

- 751 2. **Benafif S, Kote-Jarai Z, Eeles RA.** A Review of Prostate Cancer Genome Wide
752 Association Studies (GWAS). *Cancer Epidemiol Biomarkers Prev* 2018;27:845–857.
- 753 3. **Cotsapas C, Mitrovic M.** Genome-wide association studies of multiple sclerosis.
754 *Clin Transl Immunology* 2018;7:e1018.
- 755 4. **Klein RJ, Zeiss C, Chew EY, Tsai J-Y, Sackler RS, et al.** Complement Factor H
756 Polymorphism in Age-Related Macular Degeneration. *Science* 2005;308:385–389.
- 757 5. **Siontis KCM, Patsopoulos NA, Ioannidis JPA.** Replication of past candidate loci
758 for common diseases and phenotypes in 100 genome-wide association studies. *Eur*
759 *J Hum Genet* 2010;18:832–837.
- 760 6. **Power RA, Parkhill J, de Oliveira T.** Microbial genome-wide association studies:
761 lessons from human GWAS. *Nat Rev Genet* 2017;18:41–50.
- 762 7. **Lees J.** The background of bacterial GWAS.
763 https://figshare.com/articles/thesis/The_background_of_bacterial_GWAS/5550037
764 (2017).
- 765 8. **Falush D.** Bacterial genomics: Microbial GWAS coming of age. *Nat Microbiol*
766 2016;1:1–2.
- 767 9. **Vermeulen S.** *Bacterial GWAS: A Comprehensive Assessment of Challenges,*
768 *Methods and Alternatives.* Master Thesis.
769 <https://studenttheses.uu.nl/handle/20.500.12932/43914> (2023, accessed 22 April
770 2024).
- 771 10. **Coll F, Gouliouris T, Bruchmann S, Phelan J, Raven KE, et al.** PowerBacGWAS: a
772 computational pipeline to perform power calculations for bacterial genome-wide
773 association studies. *Commun Biol* 2022;5:1–12.
- 774 11. **San JE, Baichoo S, Kanzi A, Moosa Y, Lessells R, et al.** Current Affairs of Microbial
775 Genome-Wide Association Studies: Approaches, Bottlenecks and Analytical
776 Pitfalls. *Front Microbiol*;10. Epub ahead of print 30 January 2020. DOI:
777 10.3389/fmicb.2019.03119.
- 778 12. **Saber MM, Shapiro BJ.** Benchmarking bacterial genome-wide association study
779 methods using simulated genomes and phenotypes. *Microb Genom*
780 2020;6:e000337.
- 781 13. **Brynildsrud O, Bohlin J, Scheffer L, Eldholm V.** Rapid scoring of genes in microbial
782 pan-genome-wide association studies with Scoary. *Genome Biol* 2016;17:238.
- 783 14. **Collins C, Didelot X.** A phylogenetic method to perform genome-wide association
784 studies in microbes that accounts for population structure and recombination.
785 *PLOS Computational Biology* 2018;14:e1005958.

- 786 15. **Earle SG, Wu C-H, Charlesworth J, Stoesser N, Gordon NC, et al.** Identifying
787 lineage effects when controlling for population structure improves power in
788 bacterial association studies. *Nat Microbiol* 2016;1:1–8.
- 789 16. **Lees JA, Vehkala M, Välimäki N, Harris SR, Chewapreecha C, et al.** Sequence
790 element enrichment analysis to determine the genetic basis of bacterial
791 phenotypes. *Nat Commun* 2016;7:12797.
- 792 17. **Lees JA, Galardini M, Bentley SD, Weiser JN, Corander J.** pyseer: a
793 comprehensive tool for microbial pangenome-wide association studies.
794 *Bioinformatics* 2018;34:4310–4312.
- 795 18. **Zhou X, Stephens M.** Genome-wide efficient mixed-model analysis for association
796 studies. *Nat Genet* 2012;44:821–824.
- 797 19. **Turner AM, Li L, Monk IR, Lee JYH, Ingle DJ, et al.** Rifaximin prophylaxis causes
798 resistance to the last-resort antibiotic daptomycin. *Nature* 2024;1–9.
- 799 20. **Negrete-Paz AM, Vázquez-Marrufo G, Gutiérrez-Moraga A, Vázquez-**
800 **Garcidueñas MS.** Pangenome Reconstruction of Mycobacterium tuberculosis as a
801 Guide to Reveal Genomic Features Associated with Strain Clinical Phenotype.
802 *Microorganisms* 2023;11:1495.
- 803 21. **Hyun JC, Monk JM, Szubin R, Hefner Y, Palsson BO.** Global pathogenomic
804 analysis identifies known and candidate genetic antimicrobial resistance
805 determinants in twelve species. *Nat Commun* 2023;14:7690.
- 806 22. **Cuénod A, Agnetti J, Seth-Smith HMB, Roloff T, Wälchli D, et al.** Bacterial
807 genome-wide association study substantiates papGII of Escherichia coli as a major
808 risk factor for urosepsis. *Genome Med* 2023;15:89.
- 809 23. **Jaillard M, Lima L, Tournoud M, Mahé P, van Belkum A, et al.** A fast and agnostic
810 method for bacterial genome-wide association studies: Bridging the gap between k-
811 mers and genetic events. *PLoS Genet* 2018;14:e1007758.
- 812 24. **Olawoye IB, Frost SDW, Happi CT.** The Bacteria Genome Pipeline (BAGEP): an
813 automated, scalable workflow for bacteria genomes with Snakemake. *PeerJ*
814 2020;8:e10121.
- 815 25. **Chukamnerd A, Jeenkeawpiam K, Chusri S, Pomwised R, Singkhamanan K, et**
816 **al.** BacSeq: A User-Friendly Automated Pipeline for Whole-Genome Sequence
817 Analysis of Bacterial Genomes. *Microorganisms*;11. Epub ahead of print July 2023.
818 DOI: 10.3390/microorganisms11071769.
- 819 26. **Wu J, Pipistrelle J.** jessiewu/bacterialGWAS.
820 <https://github.com/jessiewu/bacterialGWAS> (2024, accessed 9 May 2024).

- 821 27. **Consortium TCr**. Genome-wide association studies of global Mycobacterium
822 tuberculosis resistance to 13 antimicrobials in 10,228 genomes identify new
823 resistance mechanisms. *PLoS Biology* 2022;20:e3001755.
- 824 28. **Burgaya J, Damaris BF, Fiebig J, Galardini M**. microGWAS: a computational
825 pipeline to perform large scale bacterial genome-wide association studies.
826 2024;2024.07.08.602456.
- 827 29. **Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, et al**. Nextflow
828 enables reproducible computational workflows. *Nat Biotechnol* 2017;35:316–319.
- 829 30. **Lees J, Galardini M**. pyseer documentation. *pyseer documentation*.
830 <https://pyseer.readthedocs.io/en/master> (2018, accessed 22 April 2024).
- 831 31. **Chewapreecha C, Marttinen P, Croucher NJ, Salter SJ, Harris SR, et al**.
832 Comprehensive Identification of Single Nucleotide Polymorphisms Associated with
833 Beta-lactam Resistance within Pneumococcal Mosaic Genes. *PLoS Genetics*
834 2014;10:e1004547.
- 835 32. **Croucher NJ, Finkelstein JA, Pelton SI, Parkhill J, Bentley SD, et al**. Population
836 genomic datasets describing the post-vaccine evolutionary epidemiology of
837 *Streptococcus pneumoniae*. *Sci Data* 2015;2:150058.
- 838 33. **Toleman MS, Reuter S, Jamrozny D, Wilson HJ, Blane B, et al**. Prospective genomic
839 surveillance of methicillin-resistant *Staphylococcus aureus* (MRSA) associated with
840 bloodstream infection, England, 1 October 2012 to 30 September 2013.
841 *Eurosurveillance* 2019;24:1800215.
- 842 34. **Bujdoš D, Walter J, O’Toole PW**. aurora: a machine learning gwas tool for analyzing
843 microbial habitat adaptation. *Genome Biol* 2025;26:66.
- 844 35. **Seemann T**. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*
845 2014;30:2068–2069.
- 846 36. **Schwengers O, Jelonek L, Dieckmann MA, Beyvers S, Blom J, et al**. Bakta: rapid
847 and standardized annotation of bacterial genomes via alignment-free sequence
848 identification. *Microbial Genomics* 2021;7:000685.
- 849 37. **Tonkin-Hill G, MacAlasdair N, Ruis C, Weimann A, Horesh G, et al**. Producing
850 polished prokaryotic pangenomes with the Panaroo pipeline. *Genome Biology*
851 2020;21:180.
- 852 38. **Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ**. IQ-TREE: A Fast and Effective
853 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular*
854 *Biology and Evolution* 2015;32:268–274.
- 855 39. **Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, et al**. Mash: fast
856 genome and metagenome distance estimation using MinHash. *Genome Biology*
857 2016;17:132.

- 858 40. **Feldgarden M, Brover V, Gonzalez-Escalona N, Frye JG, Haendiges J, et al.**
859 AMRFinderPlus and the Reference Gene Catalog facilitate examination of the
860 genomic links among antimicrobial resistance, stress response, and virulence. *Sci*
861 *Rep* 2021;11:12728.
- 862 41. **Chen S, Zhou Y, Chen Y, Gu J.** fastp: an ultra-fast all-in-one FASTQ preprocessor.
863 *Bioinformatics* 2018;34:i884–i890.
- 864 42. **Wood DE, Lu J, Langmead B.** Improved metagenomic analysis with Kraken 2.
865 *Genome Biology* 2019;20:257.
- 866 43. **Lu J, Breitwieser FP, Thielen P, Salzberg SL.** Bracken: estimating species
867 abundance in metagenomics data. *PeerJ Comput Sci* 2017;3:e104.
- 868 44. **Galardini M, Clermont O, Baron A, Busby B, Dion S, et al.** Major role of iron uptake
869 systems in the intrinsic extra-intestinal virulence of the genus *Escherichia* revealed
870 by a genome-wide association study. *PLOS Genetics* 2020;16:e1009065.
- 871 45. **Hadfield J, Croucher NJ, Goater RJ, Abudahab K, Aanensen DM, et al.**
872 Phandango: an interactive viewer for bacterial population genomics. *Bioinformatics*
873 2018;34:292–293.
- 874 46. **Clark AJ.** rec genes and homologous recombination proteins in *Escherichia coli*.
875 *Biochimie* 1991;73:523–532.
- 876 47. **Enright MC, Spratt BG.** Extensive variation in the *ddl* gene of penicillin-resistant
877 *Streptococcus pneumoniae* results from a hitchhiking effect driven by the penicillin-
878 binding protein 2b gene. *Mol Biol Evol* 1999;16:1687–1695.
- 879