

hAMRoaster: a tool for comparing performance of AMR gene detection software

Emily F. Wissel^A, Brooke M. Talbot^B, Noriko A. B. Toyosato^C, Robert A Petit III^D, Vicki Hertzberg^A, Anne Dunlop^E, Timothy D. Read^{D,F}

Author affiliations and ORCID

A: Nell Hodgson Woodruff School of Nursing, Emory University, Atlanta, GA, US

B: Population Biology, Ecology, and Evolution Program, Graduate Division of Biological and Biomedical Science, Emory University, Atlanta, GA, US

C: Department of Molecular & Biomedical Biology, University of Nebraska at Omaha, Omaha, NE, US

D: Division of Infectious Diseases, Department of Medicine, School of Medicine, Emory University, Atlanta, GA, USA

E: Department of Gynecology & Obstetrics, Emory University School of Medicine

F: Department of Human Genetics, School of Medicine, Emory University, Atlanta, GA, US

Author ORCIDs

EFW: ewissel@emory.edu <https://orcid.org/0000-0003-2275-8456>

BMT: <https://orcid.org/0000-0001-5246-7209>

NABT: <https://orcid.org/0000-0002-7855-3452>

RAP: <https://orcid.org/0000-0002-1350-9426>

VH: <https://orcid.org/0000-0002-8834-4363>

ALD: <https://orcid.org/0000-0002-5092-8136>

TDR: tread@emory.edu <https://orcid.org/0000-0001-8966-9680>

HAMROASTER

Abstract

Background. The use of shotgun metagenomics for AMR detection is appealing because data can be generated from clinical samples with minimal processing. Detecting antimicrobial resistance (AMR) in clinical genomic data is an important epidemiological task, yet a complex bioinformatic process. Many software tools exist to detect AMR genes, but they have mostly been tested in their detection of genotypic resistance in individual bacterial strains. Further, these tools use different databases, or even different versions of the same databases. Understanding the comparative performance of these bioinformatics tools for AMR gene detection in shotgun metagenomic data is important because this data type is increasingly used in public health and clinical settings.

Methods. We developed a software pipeline, hAMRoaster (Harmonized AMR Output compAriSon Tool ER; <https://github.com/ewissel/hAMRoaster>), for assessing accuracy of prediction of antibiotic resistance phenotypes. For evaluation purposes, we simulated a highly resistant mock community and several low resistance metagenomic short read (Illumina) samples based on sequenced strains with known phenotypes. We benchmarked nine open source bioinformatics tools for detecting AMR genes that 1) were conda or Docker installable, 2) had been actively maintained, 3) had an open source license, and 4) took FASTA or FASTQ files as input. hAMRoaster calculated sensitivity, specificity, precision, and accuracy for each tool, comparing detected AMR genes to susceptibility testing.

Conclusion. Overall, all tools were precise and accurate at all genome coverage levels tested (5x, 50x, 100x sequenced bases / genome length) in the highly resistant mock community with more variability in the low resistance community (1x coverage). This study demonstrated that different bioinformatic tools and pipelines yield differences in AMR gene identification across

HAMROASTER

4

drug classes, and that these differences become important if researchers are interested in resistance to specific drug classes.

Significance. Software selection for metagenomic AMR prediction should be driven by the context of the clinical/research questions and tolerance for true and false negative results. The ability to assess which bioinformatics tool best fits a particular dataset prior to beginning a large-scale project allows for more efficient processing and analysis using optimal tools for a particular research question. As prediction software and databases are in a state of constant refinement, the approach used here—creating synthetic communities containing taxa and phenotypes of interest along with using hAMRoaster to assess performance of candidate software—offers a template to aid researchers in selecting the most appropriate strategy at the time of analysis.

Keywords: antimicrobial resistance, bioinformatics, metagenomics

Tweet: Introducing a new pipeline for comparing results from #AMR tools from @emily_wissel @tdread_emory and others!

hAMRoaster compares detected AMR genes to known resistance, and returns a table with metrics for comparing results across tools.

Introduction

Antibiotic resistant infections pose a serious threat not only to public health but to the agricultural, veterinary, and food safety industries. The misuse of antibiotics in healthcare and livestock production has led to widespread antimicrobial resistance in diverse environments and has emerged as a threat to global health.^{1 2} The burden of multi-drug resistant pathogens is increasing globally, creating complex clinical scenarios in which there are limited (if any) therapeutic options, resulting in increased mortality and healthcare costs for common medical procedures.³ Genes that confer antimicrobial resistance (AMR) are increasingly present in commensal members of the human microbiome and are recognized as an important reservoir for conferring pathogen resistance through horizontal gene transfer.^{4,5}

Two key approaches to mitigating AMR infections are antibiotic stewardship and AMR surveillance. While antibiotic stewardship focuses on using antibiotics appropriately, AMR surveillance focuses on describing AMR genes already present in a community. Currently, AMR surveillance typically relies on phenotypic characterization through culture or genotypic characterization through molecular diagnostics based on PCR and hybridization techniques.⁶ However, there is a move toward genome-based methods⁷ with the Illumina short-read platform being the dominant platform for data generation at the present time.⁸

Sequencing technology has revolutionized research across many disciplines, with more applications found every year as both the technologies and analysis methods advance. This is particularly evident in the use of metagenomic data for the microbial surveillance of antimicrobial resistance (AMR), as microbial communities can be characterized without the need to first isolate and culture the specimen prior to analysis.⁹⁻¹¹ As the cost and time of sequencing has dramatically decreased, petabytes of data are quickly generated, with Illumina short reads

becoming more prevalent.^{8,12,13} Detecting AMR genes potential through non-culture based, high throughput DNA sequencing and bioinformatic approaches is of growing relevance and importance.

There are many bioinformatic tools created to process large amounts of data while following open-science principles.¹⁴ Open science is a term used to describe data that is Findable, Accessible, Interoperable, and Reusable or (FAIR) and that are open-source.¹⁵ With so many options available, it is important that investigators determine which open-source tool is the best suited for their research question. One way to address issues with replicability and variance across studies is to establish standardized bioinformatics pipelines and best practices, as has been done, for example, by the National Microbiome Data Collaborative (NMDC).¹⁶ However, for many researchers, a standardized bioinformatics pipeline may not be the best suited for their data or research question.¹⁴

As shotgun metagenomic sequencing is emerging as a powerful tool for detecting AMR,¹⁷ it is essential to evaluate how well different tools perform. In addition to testing AMR gene prediction tools against widely available metagenome samples, they should be compared in samples with extensive phenotypic resistance (acquired and mutational AMR genes). Here, we describe a software pipeline, hAMRoaster, that provides metrics on tool performance in detecting AMR genes from known resistant phenotypes and can therefore help in decision-making about which tools will be adequate for detecting resistance to the drug classes being studied.

Methods

For a schematic overview of the methods, see **Figure One**.

Development of a software pipeline, hAMRoaster, to assess results of antibiotic resistance prediction

hAMRoaster was written as a conda installable command line tool in a Python script and requires three inputs: a) the text output of AMR tool on a FASTQ or FASTA test file, such as a text file processed through hAMRonization,¹⁸ b) a list of known phenotypes associated with the test file or samples names, and c) (optional) a tab formatted table which matches antibiotic drugs with their drug class. If option c) is not specified a default table is used. The output of the program is a set of performance metrics that include sensitivity and specificity. A conda installable version of the software was deposited in the Bioconda¹⁹ database. The Github site for the software is <https://github.com/ewissel/hAMRoaster>.

hAMRoaster requires, as input, a formatted results table of runs by AMR detection tools. This table is identical to that produced by the hAMRonization¹⁸ software. hAMRonization is conda installable and can compile the output of many AMR tools into a unified format. shortBRED²⁰ and fARGene²¹ are not included in hAMRonization at the time of analysis, so hAMRoaster can take the path to the raw output for these tools and partially match it to the hAMRonization output.

hAMRoaster requires an input to the “known” phenotypic resistance in the mock community (--AMR_key flag of hAMRoaster), such as a result of susceptibility testing tables that are available from NCBI Biosamples. Antibiotics in the table of known phenotypic resistances are matched to their respective drug classes. Results classified as “susceptible” in susceptibility testing are considered “susceptible”, and “intermediate” results are ignored. In cases where susceptibility testing occurred with two or more agents, each agent is considered independently (e.g. resistance to “amoxicillin-tetracycline” is treated as resistance to

“amoxicillin” and “tetracycline” independently). Each identified AMR gene is labeled with its corresponding drug class for comparison. In instances where a gene confers resistance to multiple drug classes, the detected gene is split into multiple rows so that each conferred resistance can be independently compared to the susceptibility testing. Gene to drug class linkage is verified using the CARD database²² when applicable by accession ID. Any genes corresponding to ‘unknown’ or ‘other’ drug classes (including hypothetical resistance genes) are excluded from further analysis. Genes that confer resistance to an antibiotic that is only administered and effective in combination with another drug (e.g. clavulanic acid in amoxicillin-clavulanic acid) are classified as ‘Other’ and excluded from analysis.

A detected AMR gene is labeled as a true positive by hAMRoaster if the drug class matched to an AMR gene corresponds to a drug class that tests “resistant” in the susceptibility testing for the mock community. Similarly, a false positive is coded as a drug class that is called by the software, but tested as susceptible in the mock community (--AMR key parameter). Observed AMR genes are labeled “unknown” if the corresponding drug class is not tested in the mock community and is not included in the AMR key file. Once true/false positives and true/false negatives are determined per tool, hAMRoaster calculates sensitivity, specificity, precision, accuracy, and percent unknown.

Creation of multiple synthetic mock communities of antibiotic resistance bacteria

Simple synthetic community with high resistance

Bacterial members of the base mock community were chosen from NCBI’s BioSample Database²³ and met the following criteria: (1) the strain had extensive antibiotic susceptibility testing data using CLSI or EUCAST testing standards as part of the public NCBI BioSample record; (2) the strain was isolated from human tissue; (3) the strain was the cause of a clinical

infection; (4) the FASTA was available to download from NCBI BioSample Database.²³ Eight bacteria, each representing a different species, with overlapping resistance to 43 antibiotics across 18 drug classes, were selected for the mock community (**Table 1**). The included taxa were *Acinetobacter baumannii* MRSN489669, *Citrobacter freundii* MRSN12115, *Enterobacter cloacae* 174, *Escherichia coli* 222, *Klebsiella pneumoniae* CCUG 70742, *Pseudomonas aeruginosa* CCUG 70744, *Neisseria gonorrhoeae* SW0011, and *Staphylococcus aureus* LAC (Table 1).

Paired-end FASTQs were simulated by NCBI's ART²⁴ using default parameters for HiSeq 2500 at three levels of average sequence coverage (5x, 50x, and 100x sequenced bases / genome length) and are available on FigShare (<https://figshare.com/account/home#/projects/125974>). Simulated FASTQs were subsequently concatenated to resemble shotgun metagenomics reads, and metaSPAdes²⁵ was used to create assembled contigs. The FASTQs were simulated with approximately equal numbers of reads of each genome.

Complex synthetic clinical mock community with low resistance

We created a community profile with previously simulated human metagenomes²⁶ and added a single AMR isolate collected from a human infection at 1x coverage to simulate a human metagenome with restrictive phenotypic resistance. We included samples 0 through 5 from CAMISIM,²⁶ a set of previously simulated human metagenomes, and combined these with simulated fastqs from one of two isolates from human infections, SRR17789825²⁷ for even sample numbers and SRR16683675²⁸ for odd sample numbers.

Running antibiotic prediction software on mock communities

All tools for AMR prediction were run on the mock community and restrictive samples at all coverage levels using default settings for either simulated FASTQ or assembled contigs.

Default settings were used as it is what most users use and understand to be the developer recommendations. When both options were available, assembled contigs were run.

Statistical Analysis

Data were analyzed in Python v3.7.7 and plotted in R v4.0.4. hAMRoaster calculated all performance metrics reported in Table 3. Unweighted Cohen's kappa was calculated using R package IRR²⁹ for each pairwise combination of tools to test agreement between tools.

Data Availability

All data and code is available on the hAMRoaster GitHub repository (<https://github.com/ewissel/hAMRoaster>) and figshare (for large files; <https://figshare.com/account/home#/projects/125974>)

Results

Selection of nine open source, conda-installable tools for detection of antibiotic resistance phenotypes

To identify tools for antibiotic resistance prediction, we used a multi-headed search strategy. We searched PubMed using terms “AMR”, “antibiotic resistance genes”, “bioinformatics”, and “antimicrobial resistance”. We also searched GitHub using the same set of terms. Once an initial list of tools was compiled, we performed a second PubMed literature review including the search terms from above plus the names of the tools (“tool 1” OR “tool 2”). We also used Twitter to ask the research community what bioinformatic tools they use to identify

AMR (supplementary text 1). These searches identified 16 potential tools to identify AMR genes (**Table 2**). The search for tools concluded on March 1, 2021.

For an identified tool to be considered eligible for comparison, it had to meet the following criteria: (1) be conda or Docker installable; (2) have source code publicly available in a data repository and be actively maintained (defined as tool updates or GitHub responses within the last year); (3) have an open source license; and (4) take FASTQs or FASTAs as input files. Nine tools met the criteria to be included in this analysis: ABRicate³⁰, fARGene³¹ ResFinder³², shortBRED²⁰, RGI³³, AMRFinderPlus³⁴, starAMR³⁵, sraX³⁶, and deepARG³⁷. PointFinder also qualified³⁸, but was a subtool of ResFinder and only identified mutational resistance for some organisms, so it was excluded from analysis. The code used to install and run all tools is available on the hAMRoaster GitHub.

ABRicate

ABRicate v.1.0.1 took contig FASTA files as inputs and compared reads against NCBI AMRFinder Plus³⁴ by default, though there are options to compare against CARD,³³ ResFinder,³² ARG-ANNOT,³⁹ MEGARES,⁴⁰ EcoOH,⁴¹ PlasmidFinder,⁴² VFDB,⁴³ and Ecoli_VF,⁴⁴ which are also pre-downloaded. ABRicate reported on acquired AMR genes and not mutational resistance.

shortBRED

shortBRED²⁰ v0.9.3 used a set of marker genes to search metagenomic data for protein families of interest. The bioBakery⁴⁵ team published an AMR gene marker database built from 849 AR protein families derived from the ARDB⁴⁶ v1.1 and independent curation alongside shortBRED, which is used in this study.

fARGene

fARGene^{21,31} v.0.1 used Hidden Markov Models to detect AMR genes from short metagenomic data or long read data. This was a different approach from most other tools which compare the reads directly. fARGene has three pre-built models for detecting resistance to quinolone, tetracycline, and beta lactamases, which were tested in this study. fARGene can predict unknown ARGs using its gene models.

RGI

RGI³³ v5.1.1 used protein homology and SNP models to predict ‘resistomes’. It used CARD’s protein homolog models as a database. RGI predicts open reading frames (ORFs) using Prodigal,⁴⁷ detects homologs with BLAST,⁴⁸ and matches to CARD’s database and model cut off values.

ResFinder

ResFinder³² v4.0 was available both as a web-based application or the command line. We used ResFinder 4 in this study, which was specifically designed for detecting genotypic resistance in phenotypically resistant samples. ResFinder aligned reads directly to its own curated database without need for assembly.

deepARG

deepARG³⁷ v.2.0 used a supervised deep learning based approach for antibiotic resistance gene annotation of metagenomic sequences. It combines three databases—CARD, ARDB, and UNIPROT—and categorizes them into resistance categories.

sraX

sraX³⁶ v.1.5 was built as a one step tool; in a single command, sraX downloaded a database and aligned contigs to this database with DIAMOND⁴⁹. By default, sraX used CARD, though other options can be specified. As we use default settings for all tools, only CARD was used in this study for sraX. It should be noted that the one step aspect is convenient, but can become lengthy if there are multiple runs and databases need to be downloaded multiple times.

starAMR

starAMR^{35,50} v.0.7.2 used BLAST+⁵¹ to compare contigs against a combined database with data from ResFinder, PointFinder, and PlasmidFinder.

AMR Finder Plus

AMR Finder Plus³⁴ v.3.9.3 used BLASTX⁴⁸ translated searches and hierarchical tree of gene families to detect AMR genes. The database was derived from the Pathogen Detection Reference Gene Catalog⁵² and was compiled as part of the National Database of Antibiotic Resistant Organisms (NDARO).

Performance of software on synthetic metagenomes with high- and -low-prevalence of AMR phenotypes

Each software tool was run against a synthetic mock community of 8 bacteria at three coverage levels that expressed 43 antibiotic resistance phenotypes. Overall, the number of AMR genes detected across all tools ranged from 13 to over 700 at 100x coverage (**Table 3**). For some tools, genes detected did not correspond to a tested phenotype in the mock community, so the prediction fell into the “unknown” category. Among the tools tested, AMR Finder Plus had the

highest degree of unclassifiable/unknown results (observed AMR gene not tested in the mock community; **Figure 3**). An overview of these results are available in **Figure 2A**.

After filtering out the AMR genes detected in the simulated human metagenomes (for which AMR phenotypes were unknown), detected AMR genes were examined per sample. None of the tools detected true or false positives for one of the AMR isolates in the low resistance samples (**Figure 2b**). Fewer genes were detected overall compared to the highly resistant sample, as expected for samples with a limited resistance phenotype (**Table 3**), though many of these corresponded to unknown AMR phenotypes and not those included in susceptibility testing.

Sensitivity and Specificity

Sensitivity tests what portion of AMR genes are correctly identified by a tool when phenotypic resistance to the drug class that gene confers resistance to is present in the mock community. Specificity tests what portion of known negatives (i.e. susceptible drugs from phenotypic testing) do not have AMR genes detected for that drug class. Sensitivity for phenotype detection ranged from >0.99 (RGI) to 0.23 (sraX) at the lowest coverage levels for the highly resistant, antibiotic resistance gene (ARG)-rich dataset sample (**Fig. 2a**). In general, genome coverage did not greatly affect sensitivity, with the exception of sraX, which increased to 0.53 at the highest level. fARGene and deepARG had a high sensitivity value (>0.90) at all coverage levels. RGI, deepARG, and fARGene are all tools that compare reads to a model of AMR instead of aligning reads directly to a database, indicating that this method may be appropriate when high sensitivity values are preferred. As a note, in this ARG-rich dataset, there were only 2 possible true negatives because only two drug classes were always susceptible to antibiotics in those two drug classes when tested (nitrofurans and polypeptide).

In samples with lower numbers of resistance genes, sensitivity and specificity were variable within- and across-tools for samples, with sensitivity much lower than the high resistance community ($0 - <0.45$; **Fig. 1b**) Specificity was much higher overall, though variable across samples depending on whether any true positives were detected by the tools (**Table 3**). Precision was highly variable across tools with no consistent trend across tools (range 0 - 1), while accuracy was less variable, with most tools having an accuracy between .50 and 0.75

Concordance between tools

An analysis of the agreement between tools of detected resistance to drug classes revealed that overall, agreement was highly variable across tools (0.02 - 0.72 at 5x coverage, **Fig. 5A**) between tools at all coverage levels for the ARG-rich dataset (**Figure 5A, 5B, 5C**). Low agreement was found between most tools in the low AMR samples with the exception of AMR Finder Plus, abricate, and ResFinder4, which had a kappa value > 0.80 (**Figure 5D**).

Discussion

Development of a framework for assessing AMR prediction software performance using synthetic data

There is a considerable research effort to develop new software for predicting AMR using DNA sequence alone. In this dynamic environment, there is a need for researchers and epidemiologists to understand the relative performance of open source software tools . While some tools currently exist for compiling the results of several AMR tools together (hAMRonizer and chARMedDb⁵³), this study was motivated by the lack of an open-source pipeline for comparing the results once compiled.

The central challenge in developing this software was to compare detected AMR genes to resistance phenotypes. Detected AMR genes needed to be classified by their corresponding drug

class(es) so they could be matched to the known phenotypically resistant drug classes. One hurdle in this translation is that tools use different databases, and some databases classify genes differently. For example, shortBRED classifies gene families, while CARD classifies specific genes. While this analysis checked the drug classification via the DNA/Protein Accession value in CARD, only around 300 of the >1,000 genes detected could directly map to genes in CARD by accession value. The hAMRonization tool overcomes this challenge by providing a drug class column and filling in the values from ChEBI ontology⁵⁴ when possible. The hAMRoaster strategy is to assign a CARD drug class value to every detected AMR gene first by accession number, then by gene name. If neither of these methods assign a drug class for an AMR gene, then the drug class provided by hAMRonization is used. Another challenge in converting detected AMR genes to drug classes is that some drugs are only administered in combination, such as clavulanic acid with amoxicillin. For these instances, resistance to the drug only used in combination (e.g. clavulanic acid) is treated as an “other” drug class and excluded from analysis in hAMRoaster. In these cases, we incorporated the experience of practicing clinicians to identify combination antibiotics into the hAMRoaster antibiotic key.

The analysis presented here used synthetic data to compare tool performance. Synthetic data has the benefit of allowing controlled input with known ground truth. Therefore users can focus on the types of organisms and phenotypes they need to detect in their own datasets, perform experiments with real samples, and manipulate a range of factors such as relative abundance and sequencing error. The NCBI BioSample repository (used in this study) is an invaluable resource for creating such datasets as it contains many samples with AMR phenotypes determined by international standards. Researchers could also sequence and phenotype culturable organisms in their own laboratories to provide testing standards to evaluate software.

Here, we exclusively examined synthetic short read Illumina data, but this analysis strategy could be adapted to understand the effect of using data generated on long read technologies such as the Pacific Bioscience and Oxford Nanopore platforms.

Overall trends in performance and reasons for variability between tools

We found the sensitivity of almost all tools to be very good in a highly resistant sample (>0.80), with the exception of sraX, which had a proportionally high number of false negatives compared to true positives. However, sensitivity was lower in low-resistance samples ($0 - <0.45$), indicating that tool selection plays an important role in results for targeted AMR studies. All tools except shortBRED and starAMR detected a large number of genes that were not associated with a lab-determined phenotype in our highly resistant mock community, while this was true for all tools except starAMR in the low-resistance sample. In practice, researchers and epidemiologists may be only interested in a narrow range of AMR phenotypes. Overall, these results indicate when researchers are interested in resistance to a particular drug class as opposed to resistance to a broad range of drug classes, tool selection becomes very important.

We calculated Cohen's kappa to capture the agreement at the drug class level between AMR tools to see if all AMR tools detected resistance to the same drug classes across samples. We found that agreement at the drug class level was surprisingly low across all tools in the high and low resistance data, though some pairs of tools have higher agreement than others (e.g., AMR Finder Plus, abricate, and ResFinder4 in the low resistance samples; **Figure 5**), indicating that some tools may be better suited for detecting different types of resistance. As such, hAMRoaster provided a table with the number of genes detected per drug class for each tool that may help researchers in selecting an AMR gene detection tool that is best suited for their research question.

This research underlines the need for the further development of software tools for the detection of AMR genes in the human microbiome. It is increasingly recognized that the confined location and genetic diversity of this microbial population provides ideal conditions for genetic exchange among residential microbes and between residential and transient microbes, including pathogenic microbes. Notably, rates of horizontal gene transfer among bacteria in the human microbiome (especially the gastrointestinal tract) are estimated to be many times higher than among bacteria in other diverse ecosystems, such as soil.⁵⁵ Refined tools appropriate for use in shotgun metagenomic data will be important for tracking the spread of AMR genes from diverse environmental sources to the human microbiome and across sites in the human body and understanding whether AMR genes are derived from vertical inheritance or via horizontal gene transfer.

In conclusion, this study compared bioinformatics tools for detecting AMR genes in a simulated short read metagenomic sample at three coverage levels at one time point. While tools use slightly different methods and databases, these tools overall had high sensitivity for detection of AMR genes. Moreover, agreement between tools was sometimes low, indicating the importance of careful tool selection. We advocate that researchers should test these software tools using pipelines such as hAMRoaster with a synthetic community that highlights the resistance profiles and sample of interest.

Acknowledgements

We thank Jon Moller for helping to create the hAMRoaster name.

Funding

EFW is supported by the National Science Foundation Graduate Research Fellowship under grant 1937971. NABT is funded through the National Summer Undergraduate Research Program (NSURP) via NSF grant 2149582.

References

- (1) Shao, Y.; Wang, Y.; Yuan, Y.; Xie, Y. A Systematic Review on Antibiotics Misuse in Livestock and Aquaculture and Regulation Implications in China. *Sci. Total Environ.* **2021**, *798*, 149205. <https://doi.org/10.1016/j.scitotenv.2021.149205>.
- (2) Toni Poole; Cynthia Sheffield. Use and Misuse of Antimicrobial Drugs in Poultry and Livestock: Mechanisms of Antimicrobial Resistance. *Pak. Vet. J.* **2013**, *33* (3), 266–271.
- (3) Teillant, A.; Gandra, S.; Barter, D.; Morgan, D. J.; Laxminarayan, R. Potential Burden of Antibiotic Resistance on Surgery and Cancer Chemotherapy Antibiotic Prophylaxis in the USA: A Literature Review and Modelling Study. *Lancet Infect. Dis.* **2015**, *15* (12), 1429–1437. [https://doi.org/10.1016/S1473-3099\(15\)00270-4](https://doi.org/10.1016/S1473-3099(15)00270-4).
- (4) Nji, E.; Kazibwe, J.; Hambridge, T.; Joko, C. A.; Larbi, A. A.; Dampney, L. A. O.; Nkansa-Gyamfi, N. A.; Stålsby Lundborg, C.; Lien, L. T. Q. High Prevalence of Antibiotic Resistance in Commensal Escherichia Coli from Healthy Human Sources in Community Settings. *Sci. Rep.* **2021**, *11* (1), 3372. <https://doi.org/10.1038/s41598-021-82693-4>.
- (5) Brinkac, L.; Voorhies, A.; Gomez, A.; Nelson, K. E. The Threat of Antimicrobial Resistance on the Human Microbiome. *Microb. Ecol.* **2017**, *74* (4), 1001–1008. <https://doi.org/10.1007/s00248-017-0985-z>.
- (6) Anjum, M. F.; Zankari, E.; Hasman, H. Molecular Methods for Detection of Antimicrobial Resistance. *Microbiol. Spectr.* **2017**, *5* (6). <https://doi.org/10.1128/microbiolspec.ARBA-0011-2017>.
- (7) Nutrition, C. for F. S. and A. GenomeTrakr Network. *FDA* **2021**.
- (8) Porter, T. M.; Hajibabaei, M. Over 2.5 Million COI Sequences in GenBank and Growing. *PLOS ONE* **2018**, *13* (9), e0200177. <https://doi.org/10.1371/journal.pone.0200177>.
- (9) Kraemer, S. A.; Ramachandran, A.; Perron, G. G. Antibiotic Pollution in the Environment: From Microbial Ecology to Public Policy. *Microorganisms* **2019**, *7* (6), 180. <https://doi.org/10.3390/microorganisms7060180>.
- (10) Hendriksen, R. S.; Bortolaia, V.; Tate, H.; Tyson, G. H.; Aarestrup, F. M.; McDermott, P. F. Using Genomics to Track Global Antimicrobial Resistance. *Front. Public Health* **2019**, *7*.
- (11) Kumar, D.; Pornsukarom, S.; Thakur, S. Antibiotic Usage in Poultry Production and Antimicrobial-Resistant Salmonella in Poultry. In *Food Safety in Poultry Meat Production*; Venkitanarayanan, K., Thakur, S., Ricke, S. C., Eds.; Food Microbiology and Food Safety; Springer International Publishing: Cham, 2019; pp 47–66. https://doi.org/10.1007/978-3-030-05011-5_3.
- (12) Robinson, T.; Harkin, J.; Shukla, P. Hardware Acceleration of Genomics Data Analysis: Challenges and Opportunities. *Bioinformatics* **2021**, *37* (13), 1785–1795. <https://doi.org/10.1093/bioinformatics/btab017>.
- (13) *GenBank and WGS Statistics*. <https://www.ncbi.nlm.nih.gov/genbank/statistics/> (accessed 2022-07-29).
- (14) de Abreu, V. A. C.; Perdigão, J.; Almeida, S. Metagenomic Approaches to Analyze Antimicrobial Resistance: An Overview. *Front. Genet.* **2021**, *11*.
- (15) Wilkinson, M. D.; Dumontier, M.; Aalbersberg, I. J.; Appleton, G.; Axton, M.; Baak, A.; Blomberg, N.; Boiten, J.-W.; da Silva Santos, L. B.; Bourne, P. E.; Bouwman, J.; Brookes, A. J.; Clark, T.; Crosas, M.; Dillo, I.; Dumon, O.; Edmunds, S.; Evelo, C. T.; Finkers, R.; Gonzalez-Beltran, A.; Gray, A. J. G.; Groth, P.; Goble, C.; Grethe, J. S.; Heringa, J.; 't Hoen, P. A. C.; Hooft, R.; Kuhn, T.; Kok, R.; Kok, J.; Lusher, S. J.; Martone, M. E.; Mons, A.; Packer, A. L.; Persson, B.; Rocca-Serra, P.; Roos, M.; van Schaik, R.; Sansone, S.-A.; Schultes, E.; Sengstag, T.; Slater, T.; Strawn, G.; Swertz, M. A.; Thompson, M.; van der Lei, J.; van Mulligen, E.; Velterop, J.; Waagmeester, A.; Wittenburg, P.; Wolstencroft, K.; Zhao, J.; Mons, B. The FAIR Guiding Principles for Scientific Data Management and Stewardship. *Sci. Data* **2016**, *3* (1), 160018. <https://doi.org/10.1038/sdata.2016.18>.

- (16) Eloë-Fadrosh, E. A.; Ahmed, F.; Anubhav; Babinski, M.; Baumes, J.; Borkum, M.; Bramer, L.; Canon, S.; Christianson, D. S.; Corilo, Y. E.; Davenport, K. W.; Davis, B.; Drake, M.; Duncan, W. D.; Flynn, M. C.; Hays, D.; Hu, B.; Huntemann, M.; Kelliher, J.; Lebedeva, S.; Li, P.-E.; Lipton, M.; Lo, C.-C.; Martin, S.; Millard, D.; Miller, K.; Miller, M. A.; Piehowski, P.; Jackson, E. P.; Purvine, S.; Reddy, T. B. K.; Richardson, R.; Rudolph, M.; Sarrafan, S.; Shakya, M.; Smith, M.; Stratton, K.; Sundaramurthi, J. C.; Vangay, P.; Winston, D.; Wood-Charlson, E. M.; Xu, Y.; Chain, P. S. G.; McCue, L. A.; Mans, D.; Mungall, C. J.; Mouncey, N. J.; Fagnan, K. The National Microbiome Data Collaborative Data Portal: An Integrated Multi-Omics Microbiome Data Resource. *Nucleic Acids Res.* **2022**, *50* (D1), D828–D836. <https://doi.org/10.1093/nar/gkab990>.
- (17) Oniciuc, E. A.; Likotrafiti, E.; Alvarez-Molina, A.; Prieto, M.; Santos, J. A.; Alvarez-Ordóñez, A. The Present and Future of Whole Genome Sequencing (WGS) and Whole Metagenome Sequencing (WMS) for Surveillance of Antimicrobial Resistant Microorganisms and Antimicrobial Resistance Genes across the Food Chain. *Genes* **2018**, *9* (5), 268. <https://doi.org/10.3390/genes9050268>.
- (18) *Issues · pha4ge/hAMRonization*. GitHub. <https://github.com/pha4ge/hAMRonization> (accessed 2021-10-12).
- (19) Grüning, B.; Dale, R.; Sjödin, A.; Chapman, B. A.; Rowe, J.; Tomkins-Tinch, C. H.; Valieris, R.; Köster, J. Bioconda: Sustainable and Comprehensive Software Distribution for the Life Sciences. *Nat. Methods* **2018**, *15* (7), 475–476. <https://doi.org/10.1038/s41592-018-0046-7>.
- (20) Kaminski, J.; Gibson, M. K.; Franzosa, E. A.; Segata, N.; Dantas, G.; Huttenhower, C. High-Specificity Targeted Functional Profiling in Microbial Communities with ShortBRED. *PLOS Comput. Biol.* **2015**, *11* (12), e1004557. <https://doi.org/10.1371/journal.pcbi.1004557>.
- (21) fannyhb. FARGene, 2021. <https://github.com/fannyhb/fargene> (accessed 2021-12-06).
- (22) Alcock, B. P.; Raphenya, A. R.; Lau, T. T. Y.; Tsang, K. K.; Bouchard, M.; Edalatmand, A.; Huynh, W.; Nguyen, A.-L. V.; Cheng, A. A.; Liu, S.; Min, S. Y.; Miroshnichenko, A.; Tran, H.-K.; Werfalli, R. E.; Nasir, J. A.; Oloni, M.; Speicher, D. J.; Florescu, A.; Singh, B.; Faltyn, M.; Hernandez-Koutoucheva, A.; Sharma, A. N.; Bordeleau, E.; Pawlowski, A. C.; Zubyk, H. L.; Dooley, D.; Griffiths, E.; Maguire, F.; Winsor, G. L.; Beiko, R. G.; Brinkman, F. S. L.; Hsiao, W. W. L.; Domselaar, G. V.; McArthur, A. G. CARD 2020: Antibiotic Resistome Surveillance with the Comprehensive Antibiotic Resistance Database. *Nucleic Acids Res.* **2020**, *48* (D1), D517–D525. <https://doi.org/10.1093/nar/gkz935>.
- (23) Barrett, T.; Clark, K.; Gevorgyan, R.; Gorelenkov, V.; Gribov, E.; Karsch-Mizrachi, I.; Kimelman, M.; Pruitt, K. D.; Resenchuk, S.; Tatusova, T.; Yaschenko, E.; Ostell, J. BioProject and BioSample Databases at NCBI: Facilitating Capture and Organization of Metadata. *Nucleic Acids Res.* **2012**, *40* (D1), D57–D63. <https://doi.org/10.1093/nar/gkr1163>.
- (24) Huang, W.; Li, L.; Myers, J. R.; Marth, G. T. ART: A next-Generation Sequencing Read Simulator. *Bioinforma. Oxf. Engl.* **2012**, *28* (4), 593–594. <https://doi.org/10.1093/bioinformatics/btr708>.
- (25) Nurk, S.; Meleshko, D.; Korobeynikov, A.; Pevzner, P. A. MetaSPAdes: A New Versatile Metagenomic Assembler. *Genome Res.* **2017**, *27* (5), 824–834. <https://doi.org/10.1101/gr.213959.116>.
- (26) Fritz, A.; Hofmann, P.; Majda, S.; Dahms, E.; Dröge, J.; Fiedler, J.; Lesker, T. R.; Belmann, P.; DeMaere, M. Z.; Darling, A. E.; Sczyrba, A.; Bremges, A.; McHardy, A. C. CAMISIM: Simulating Metagenomes and Microbial Communities. *Microbiome* **2019**, *7* (1), 17. <https://doi.org/10.1186/s40168-019-0633-6>.
- (27) Biggs, S. L.; Jennison, A. V.; Bergh, H.; Graham, R.; Nimmo, G.; Whiley, D. Limited Evidence of Patient-to-Patient Transmission of Staphylococcus Aureus Strains between

- Children with Cystic Fibrosis, Queensland, Australia. *PLOS ONE* **2022**, *17* (10), e0275256. <https://doi.org/10.1371/journal.pone.0275256>.
- (28) SRR16683675 NCBI. https://trace.ncbi.nlm.nih.gov/Traces/?view=run_browser&acc=SRR16683675&display=metadata.
- (29) Gamer, M.; Lemon, J.; Fellows, I.; Singh, P. IRR.
- (30) Seemann, T. ABRicate, 2021. <https://github.com/tseemann/abricate> (accessed 2021-10-12).
- (31) Berglund, F.; Österlund, T.; Boulund, F.; Marathe, N. P.; Larsson, D. G. J.; Kristiansson, E. Identification and Reconstruction of Novel Antibiotic Resistance Genes from Metagenomes. *Microbiome* **2019**, *7* (1), 52. <https://doi.org/10.1186/s40168-019-0670-1>.
- (32) Bortolaia, V.; Kaas, R. S.; Ruppe, E.; Roberts, M. C.; Schwarz, S.; Cattoir, V.; Philippon, A.; Allesoe, R. L.; Rebelo, A. R.; Florensa, A. F.; Fagelhauer, L.; Chakraborty, T.; Neumann, B.; Werner, G.; Bender, J. K.; Stingl, K.; Nguyen, M.; Coppens, J.; Xavier, B. B.; Malhotra-Kumar, S.; Westh, H.; Pinholt, M.; Anjum, M. F.; Duggett, N. A.; Kempf, I.; Nykäsenoja, S.; Olkkola, S.; Wiczorek, K.; Amaro, A.; Clemente, L.; Mossong, J.; Losch, S.; Ragimbeau, C.; Lund, O.; Aarestrup, F. M. ResFinder 4.0 for Predictions of Phenotypes from Genotypes. *J. Antimicrob. Chemother.* **2020**, *75* (12), 3491–3500. <https://doi.org/10.1093/jac/dkaa345>.
- (33) Alcock, B. P.; Raphenya, A. R.; Lau, T. T. Y.; Tsang, K. K.; Bouchard, M.; Edalatmand, A.; Huynh, W.; Nguyen, A.-L. V.; Cheng, A. A.; Liu, S.; Min, S. Y.; Miroshnichenko, A.; Tran, H.-K.; Werfalli, R. E.; Nasir, J. A.; Oloni, M.; Speicher, D. J.; Florescu, A.; Singh, B.; Faltyn, M.; Hernandez, A.; Koutoucheva; Sharma, A. N.; Bordeleau, E.; Pawlowski, A. C.; Zubyk, H. L.; Dooley, D.; Griffiths, E.; Maguire, F.; Winsor, G. L.; Beiko, R. G.; Brinkman, F. S. L.; Hsiao, W. W. L.; Domselaar, G. V.; McArthur, A. G. CARD 2020: Antibiotic Resistance Surveillance with the Comprehensive Antibiotic Resistance Database, 2020. <https://doi.org/10.1093/nar/gkz935>.
- (34) NCBI Antimicrobial Resistance Gene Finder (AMRFinderPlus), 2021. <https://github.com/ncbi/amr> (accessed 2021-10-12).
- (35) Staramr, 2021. <https://github.com/phac-nml/staramr> (accessed 2021-10-12).
- (36) Panunzi, L. G. SraX: A Novel Comprehensive Resistome Analysis Tool. *Front. Microbiol.* **2020**, *11*, 52. <https://doi.org/10.3389/fmicb.2020.00052>.
- (37) Arango-Argoty, G.; Garner, E.; Pruden, A.; Heath, L. S.; Vikesland, P.; Zhang, L. DeepARG: A Deep Learning Approach for Predicting Antibiotic Resistance Genes from Metagenomic Data. *Microbiome* **2018**, *6* (1), 23. <https://doi.org/10.1186/s40168-018-0401-z>.
- (38) *PointFinder: a novel web tool for WGS-based detection of antimicrobial resistance associated with chromosomal point mutations in bacterial pathogens | Journal of Antimicrobial Chemotherapy | Oxford Academic.* <https://academic.oup.com/jac/article/72/10/2764/3979530?login=true> (accessed 2021-10-12).
- (39) Gupta, S. K.; Padmanabhan, B. R.; Diene, S. M.; Lopez-Rojas, R.; Kempf, M.; Landraud, L.; Rolain, J.-M. ARG-ANNOT, a New Bioinformatic Tool To Discover Antibiotic Resistance Genes in Bacterial Genomes. *Antimicrob. Agents Chemother.* **2014**, *58* (1), 212–220. <https://doi.org/10.1128/AAC.01310-13>.
- (40) Doster, E.; Lakin, S. M.; Dean, C. J.; Wolfe, C.; Young, J. G.; Boucher, C.; Belk, K. E.; Noyes, N. R.; Morley, P. S. MEGARes 2.0: A Database for Classification of Antimicrobial Drug, Biocide and Metal Resistance Determinants in Metagenomic Sequence Data. *Nucleic Acids Res.* **2020**, *48* (D1), D561–D569. <https://doi.org/10.1093/nar/gkz1010>.
- (41) Ingle, D. J.; Valcanis, M.; Kuzevski, A.; Tauschek, M.; Inouye, M.; Stinear, T.; Levine, M. M.; Robins-Browne, R. M.; Holt, K. E. In Silico Serotyping of E. Coli from Short Read Data

- Identifies Limited Novel O-Loci but Extensive Diversity of O:H Serotype Combinations within and between Pathogenic Lineages. *Microb. Genomics* **2016**, 2 (7), e000064. <https://doi.org/10.1099/mgen.0.000064>.
- (42) Carattoli, A.; Hasman, H. PlasmidFinder and In Silico PMLST: Identification and Typing of Plasmid Replicons in Whole-Genome Sequencing (WGS). *Methods Mol. Biol. Clifton NJ* **2020**, 2075, 285–294. https://doi.org/10.1007/978-1-4939-9877-7_20.
- (43) Chen, L.; Zheng, D.; Liu, B.; Yang, J.; Jin, Q. VFDB 2016: Hierarchical and Refined Dataset for Big Data Analysis—10 Years On. *Nucleic Acids Res.* **2016**, 44 (D1), D694–D697. <https://doi.org/10.1093/nar/gkv1239>.
- (44) Escherichia Coli Virulence Factors, 2021. https://github.com/phac-nml/ecoli_vf (accessed 2021-12-09).
- (45) McIver, L. J.; Abu-Ali, G.; Franzosa, E. A.; Schwager, R.; Morgan, X. C.; Waldron, L.; Segata, N.; Huttenhower, C. BioBakery: A Meta'omic Analysis Environment. *Bioinformatics* **2018**, 34 (7), 1235–1237. <https://doi.org/10.1093/bioinformatics/btx754>.
- (46) Liu, B.; Pop, M. ARDB--Antibiotic Resistance Genes Database. *Nucleic Acids Res.* **2009**, 37 (Database), D443–D447. <https://doi.org/10.1093/nar/gkn656>.
- (47) Hyatt, D.; Chen, G.-L.; LoCascio, P. F.; Land, M. L.; Larimer, F. W.; Hauser, L. J. Prodigal: Prokaryotic Gene Recognition and Translation Initiation Site Identification. *BMC Bioinformatics* **2010**, 11 (1), 119. <https://doi.org/10.1186/1471-2105-11-119>.
- (48) McGinnis, S.; Madden, T. L. BLAST: At the Core of a Powerful and Diverse Set of Sequence Analysis Tools. *Nucleic Acids Res.* **2004**, 32 (suppl_2), W20–W25. <https://doi.org/10.1093/nar/gkh435>.
- (49) Buchfink, B.; Reuter, K.; Drost, H.-G. Sensitive Protein Alignments at Tree-of-Life Scale Using DIAMOND. *Nat. Methods* **2021**, 18 (4), 366–368. <https://doi.org/10.1038/s41592-021-01101-x>.
- (50) Zankari, E.; Hasman, H.; Cosentino, S.; Vestergaard, M.; Rasmussen, S.; Lund, O.; Aarestrup, F. M.; Larsen, M. V. Identification of Acquired Antimicrobial Resistance Genes. *J. Antimicrob. Chemother.* **2012**, 67 (11), 2640–2644. <https://doi.org/10.1093/jac/dks261>.
- (51) Camacho, C.; Coulouris, G.; Avagyan, V.; Ma, N.; Papadopoulos, J.; Bealer, K.; Madden, T. L. BLAST+: Architecture and Applications. *BMC Bioinformatics* **2009**, 10 (1), 421. <https://doi.org/10.1186/1471-2105-10-421>.
- (52) *Reference Gene Catalog - Pathogen Detection - NCBI.* <https://www.ncbi.nlm.nih.gov/pathogens/refgene/#> (accessed 2021-09-13).
- (53) *Anthony Underwood / chAMReDb.* GitLab. <https://gitlab.com/antunderwood/chamredb> (accessed 2021-10-12).
- (54) Hastings, J.; Owen, G.; Dekker, A.; Ennis, M.; Kale, N.; Muthukrishnan, V.; Turner, S.; Swainston, N.; Mendes, P.; Steinbeck, C. ChEBI in 2016: Improved Services and an Expanding Collection of Metabolites. *Nucleic Acids Res.* **2016**, 44 (D1), D1214–D1219. <https://doi.org/10.1093/nar/gkv1031>.
- (55) Huttenhower, C.; Gevers, D.; Knight, R.; Abubucker, S.; Badger, J. H.; Chinwalla, A. T.; Creasy, H. H.; Earl, A. M.; FitzGerald, M. G.; Fulton, R. S.; Giglio, M. G.; Hallsworth-Pepin, K.; Lobos, E. A.; Madupu, R.; Magrini, V.; Martin, J. C.; Mitreva, M.; Muzny, D. M.; Sodergren, E. J.; Versalovic, J.; Wollam, A. M.; Worley, K. C.; Wortman, J. R.; Young, S. K.; Zeng, Q.; Aagaard, K. M.; Abolude, O. O.; Allen-Vercoe, E.; Alm, E. J.; Alvarado, L.; Andersen, G. L.; Anderson, S.; Appelbaum, E.; Arachchi, H. M.; Armitage, G.; Arze, C. A.; Ayvaz, T.; Baker, C. C.; Begg, L.; Belachew, T.; Bhonagiri, V.; Bihan, M.; Blaser, M. J.; Bloom, T.; Bonazzi, V.; Paul Brooks, J.; Buck, G. A.; Buhay, C. J.; Busam, D. A.; Campbell, J. L.; Canon, S. R.; Cantarel, B. L.; Chain, P. S. G.; Chen, I.-M. A.; Chen, L.; Chhibba, S.; Chu, K.; Ciulla, D. M.; Clemente, J. C.; Clifton, S. W.; Conlan, S.; Crabtree, J.; Cutting, M. A.; Davidovics, N. J.; Davis, C. C.; DeSantis, T. Z.; Deal, C.; Delehaunty, K. D.; Dewhirst, F. E.; Deych, E.; Ding, Y.; Dooling, D. J.; Dugan, S. P.;

Michael Dunne, W.; Scott Durkin, A.; Edgar, R. C.; Erlich, R. L.; Farmer, C. N.; Farrell, R. M.; Faust, K.; Feldgarden, M.; Felix, V. M.; Fisher, S.; Fodor, A. A.; Forney, L. J.; Foster, L.; Di Francesco, V.; Friedman, J.; Friedrich, D. C.; Fronick, C. C.; Fulton, L. L.; Gao, H.; Garcia, N.; Giannoukos, G.; Giblin, C.; Giovanni, M. Y.; Goldberg, J. M.; Goll, J.; Gonzalez, A.; Griggs, A.; Gujja, S.; Kinder Haake, S.; Haas, B. J.; Hamilton, H. A.; Harris, E. L.; Hepburn, T. A.; Herter, B.; Hoffmann, D. E.; Holder, M. E.; Howarth, C.; Huang, K. H.; Huse, S. M.; Izard, J.; Jansson, J. K.; Jiang, H.; Jordan, C.; Joshi, V.; Katancik, J. A.; Keitel, W. A.; Kelley, S. T.; Kells, C.; King, N. B.; Knights, D.; Kong, H. H.; Koren, O.; Koren, S.; Kota, K. C.; Kovar, C. L.; Kyrpides, N. C.; La Rosa, P. S.; Lee, S. L.; Lemon, K. P.; Lennon, N.; Lewis, C. M.; Lewis, L.; Ley, R. E.; Li, K.; Liolios, K.; Liu, B.; Liu, Y.; Lo, C.-C.; Lozupone, C. A.; Dwayne Lunsford, R.; Madden, T.; Mahurkar, A. A.; Mannon, P. J.; Mardis, E. R.; Markowitz, V. M.; Mavromatis, K.; McCorrison, J. M.; McDonald, D.; McEwen, J.; McGuire, A. L.; McInnes, P.; Mehta, T.; Mihindukulasuriya, K. A.; Miller, J. R.; Minx, P. J.; Newsham, I.; Nusbaum, C.; O’Laughlin, M.; Orvis, J.; Pagani, I.; Palaniappan, K.; Patel, S. M.; Pearson, M.; Peterson, J.; Podar, M.; Pohl, C.; Pollard, K. S.; Pop, M.; Priest, M. E.; Proctor, L. M.; Qin, X.; Raes, J.; Ravel, J.; Reid, J. G.; Rho, M.; Rhodes, R.; Riehle, K. P.; Rivera, M. C.; Rodriguez-Mueller, B.; Rogers, Y.-H.; Ross, M. C.; Russ, C.; Sanka, R. K.; Sankar, P.; Fah Sathirapongsasuti, J.; Schloss, J. A.; Schloss, P. D.; Schmidt, T. M.; Scholz, M.; Schriml, L.; Schubert, A. M.; Segata, N.; Segre, J. A.; Shannon, W. D.; Sharp, R. R.; Sharpton, T. J.; Shenoy, N.; Sheth, N. U.; Simone, G. A.; Singh, I.; Smillie, C. S.; Sobel, J. D.; Sommer, D. D.; Spicer, P.; Sutton, G. G.; Sykes, S. M.; Tabbaa, D. G.; Thiagarajan, M.; Tomlinson, C. M.; Torralba, M.; Treangen, T. J.; Truty, R. M.; Vishnivetskaya, T. A.; Walker, J.; Wang, L.; Wang, Z.; Ward, D. V.; Warren, W.; Watson, M. A.; Wellington, C.; Wetterstrand, K. A.; White, J. R.; Wilczek-Boney, K.; Wu, Y.; Wylie, K. M.; Wylie, T.; Yandava, C.; Ye, L.; Ye, Y.; Yooseph, S.; Youmans, B. P.; Zhang, L.; Zhou, Y.; Zhu, Y.; Zoloth, L.; Zucker, J. D.; Birren, B. W.; Gibbs, R. A.; Highlander, S. K.; Methé, B. A.; Nelson, K. E.; Petrosino, J. F.; Weinstock, G. M.; Wilson, R. K.; White, O.; The Human Microbiome Project Consortium. Structure, Function and Diversity of the Healthy Human Microbiome. *Nature* **2012**, *486* (7402), 207–214. <https://doi.org/10.1038/nature11234>.

Figure 1: Schematic I Methods

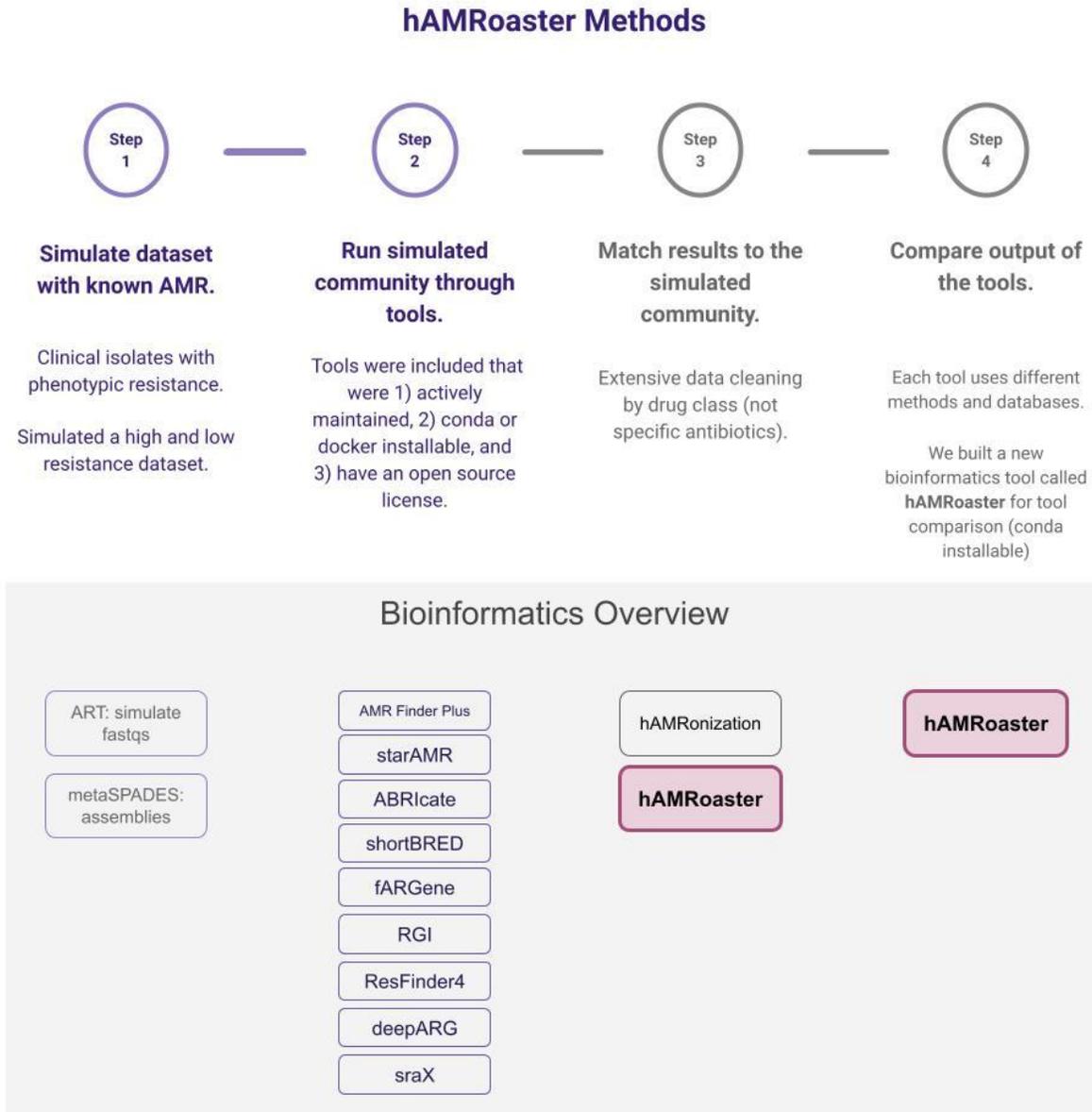


Figure 2: Antimicrobial Resistance (AMR) Genes Detected By Software Tools by Drug Class

AMR Genes detected by each tool across coverage levels, grouped into drug class to which the genes confer resistance with the color coding indicating whether the detection was true positive (green), false positive (purple) or unknown (yellow). Clear spaces in the plot indicate that AMR genes were not detected for the drug class on the x-axis by the tool on the y-axis. Plot A contains the high AMR Data, while plot B contains the low AMR data.

B

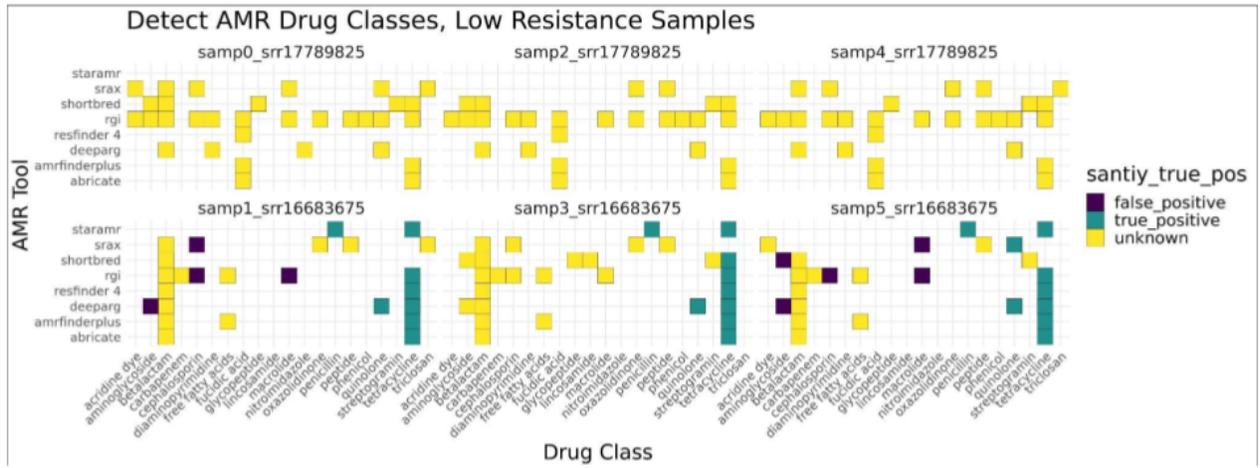


Figure 3: Sensitivity of Software Tools for Detection of Antimicrobial Resistance (AMR)

Genes Across Coverage Levels

Sensitivity was calculated as (true positives) / (true positives + false negatives). Most tools were highly sensitive (greater than 0.80). All genes corresponding to “Other” or “Unknown” drug classes were not included in these calculations. Similarly, AMR genes corresponding to phenotypic resistance that was not tested in the mock community was considered “Unknown” and not included in the sensitivity analysis.

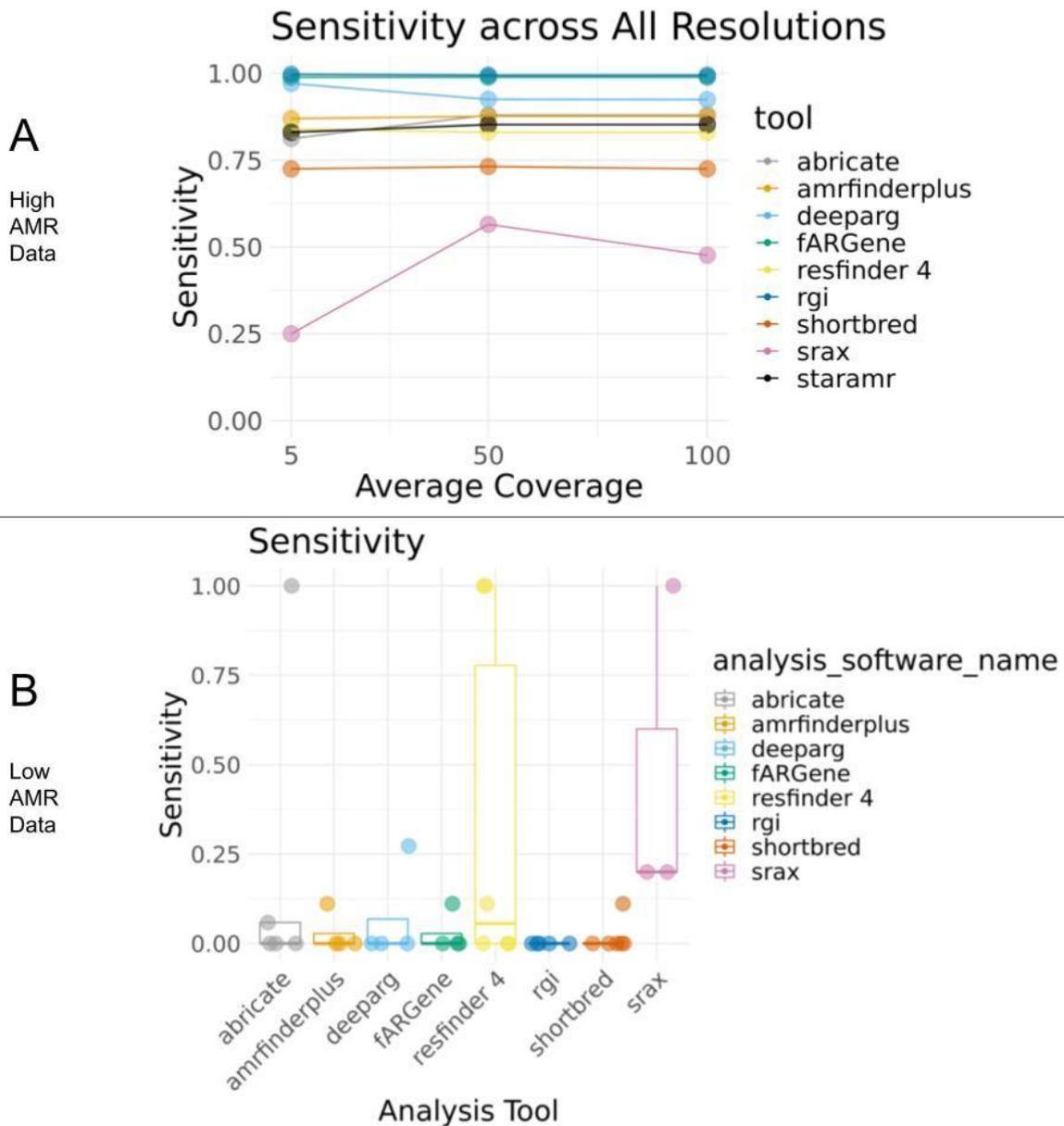


Figure 4: Percent Detection of Unknown Antimicrobial (AMR) Resistance Genes Across Coverage

The percent detection of AMR genes that could not be classified because the drug class the gene confers resistance to was not tested for the high AMR (A) and low AMR (b) data. A black dashed line is placed at 20%, indicating where at least 20% of the detected AMR genes could not be classified.

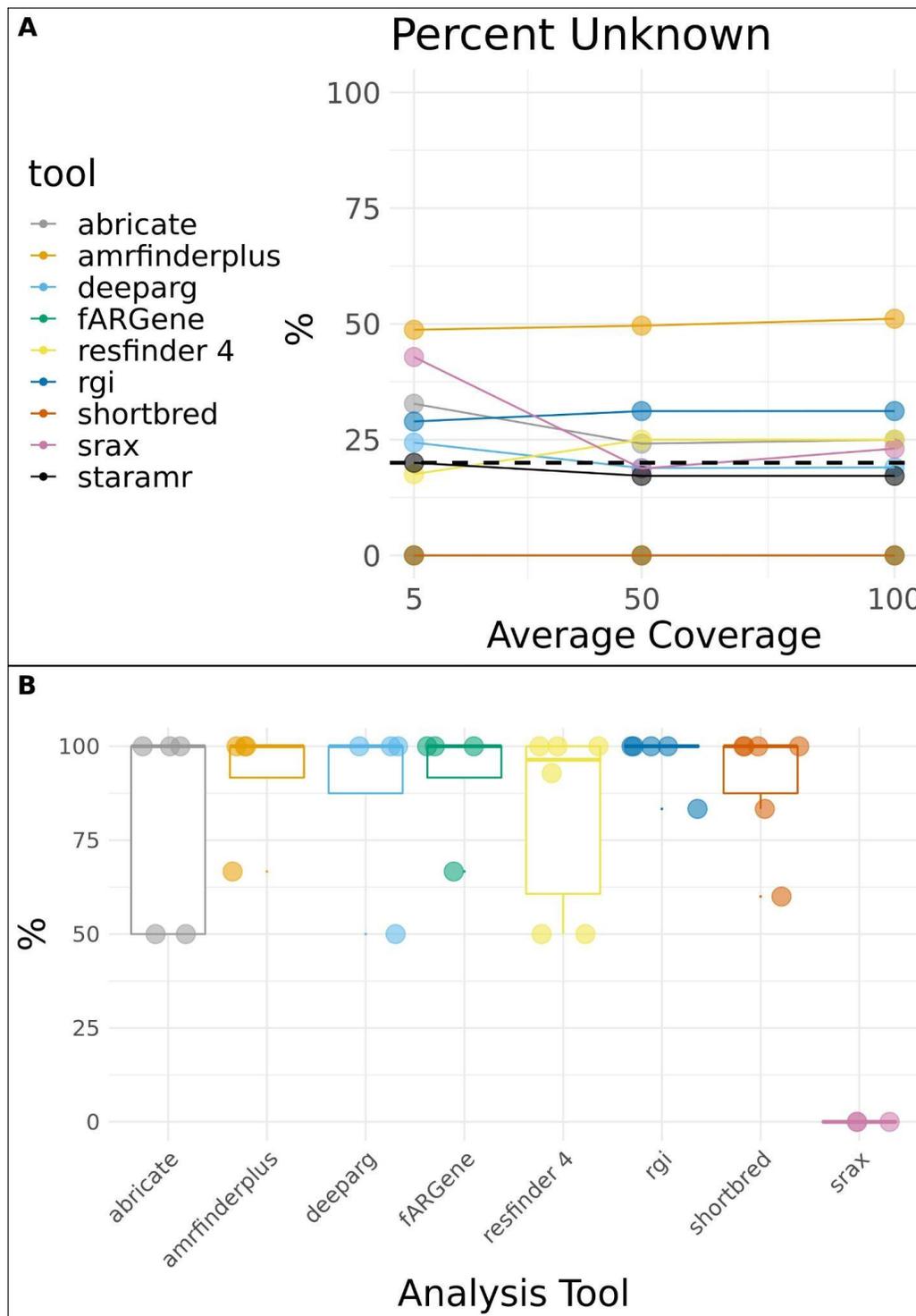


Figure 5: Agreement (Cohen’s Kappa) values between tools across coverage levels calculated in R using the kappa2 function

Agreement between tools in detecting resistance to drug classes is shaded across all plots while kappa values are bolded when the p-value is less than 0.05. A, B, and C display the agreement between tools for the 5x, 50x, and 100x coverage high AMR datasets, respectively. D displays the agreement between tools for the low AMR samples.

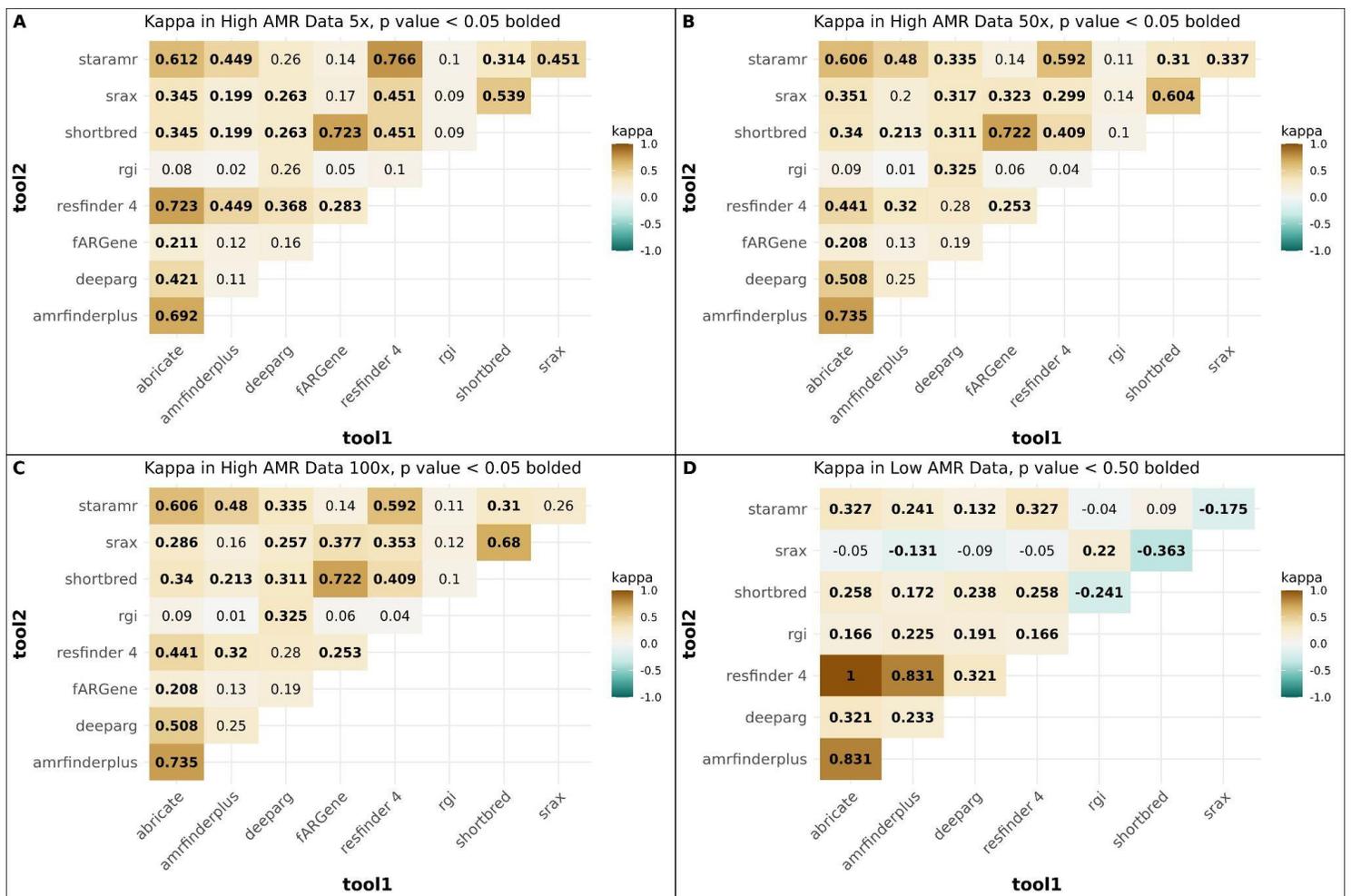


Table 1A: Clinical isolates included in the high resistance simulated community. (susceptibility test is in the spreadsheet, will have to be supplemental bc so big)

| Strain | Testing Standard (CLSI or EUCAST) | BioSample ID | Link |
|--|--------------------------------------|--------------|--|
| <i>Neisseria gonorrhoeae</i> SW0011 | CLSI | SAMN15960549 | https://www.ncbi.nlm.nih.gov/biosample/SAMN15960549 |
| <i>Klebsiella pneumoniae</i> CCUG 70742 | EUCAST | SAMN07602587 | https://www.ncbi.nlm.nih.gov/biosample/SAMN07602587 |
| <i>Pseudomonas aeruginosa</i> CCUG 70744 | EUCAST | SAMN07602569 | https://www.ncbi.nlm.nih.gov/biosample/SAMN07602569 / |
| <i>Acinetobacter baumannii</i> MRSN489669 | CLSI | SAMN12087686 | https://www.ncbi.nlm.nih.gov/biosample/SAMN12087686 |
| <i>Enterobacter cloacae</i> 174 | CLSI | SAMN04456586 | https://www.ncbi.nlm.nih.gov/biosample/SAMN04456586 |
| <i>Citrobacter freundii</i> MRSN12115 | CLSI | SAMN13412315 | https://www.ncbi.nlm.nih.gov/biosample/SAMN13412315 |

| | | | |
|-------------------------------------|------|--------------|---|
| <i>Staphylococcus aureus</i> LAC | CLSI | SAMN08391108 | https://www.ncbi.nlm.nih.gov/biosample/SAMN08391108 |
| <i>Escherichia coli</i> 222 | CLSI | SAMN05194390 | https://www.ncbi.nlm.nih.gov/biosample/SAMN05194390 |

Table 1B: Clinical isolates included in the low resistance simulated community. (susceptibility test is in the spreadsheet, will have to be supplemental bc so big)

| Strain | Testing Standard (CLSI or EUCAST) | BioSample ID | Link |
|------------------------------|--------------------------------------|--------------|---|
| <i>Staphylococcus aureus</i> | EUCAST | SAMN25295985 | https://www.ncbi.nlm.nih.gov/biosample/25295985 |
| <i>Neisseria gonorrhoeae</i> | CLSI | SAMN22824038 | https://www.ncbi.nlm.nih.gov/biosample/22824038 |

Table 2: Tools identified from search methods with the selection criteria and whether they subsequently worked or not.

| Tool | Conda / Docker Installable? | Actively Maintained? | Input format? | Included in Analysis? | Implementation Method | Database |
|------------------|------------------------------------|-----------------------------|----------------------|------------------------------|-----------------------------------|---|
| ABRicate | Yes - conda | Yes | FASTA | Yes | Align reads to specified database | NCBI (default), AMRFinder Plus, CARD, ResFinder, ARG-ANNOT, MEGARES, EcoOH, PlasmidFinder, VFDB, and Ecoli_VF |
| shortBRED | Yes - docker & conda | Yes | FASTA | Yes | Align reads to database | AMR gene marker database from 849 AR protein families from the ARDB19 and independent curation |

HAMROASTER

| | | | | | | |
|--------------------|------------------------------|---------|-------|-----|-------------------------|--|
| fARGene | Yes - conda | Yes | FASTQ | Yes | Compare to AMR model | Hidden markov models for quinolone, tetracycline, and beta lactamases |
| RGI | Yes -docker (conda outdated) | Yes | FASTQ | Yes | Compare to AMR model | Prodigal predicts ORF and compared to CARD and WildCARD |
| ResFinder 4 | Yes - docker (conda broken) | Yes | FASTA | Yes | Align reads to database | ResFinder 4 database |
| DeepARG | Yes - docker | Unclear | FASTA | Yes | Compare to AMR model | Supervised deep learning compares reads to antibiotic resistance categories created from CARD, ARDB, and UNIPROT |
| sraX | Yes - docker & | Yes | FASTA | Yes | Align reads to database | CARD by default |

HAMROASTER

| | | | | | | |
|--|--------------|-----|---------------------------------|-----|-------------------------|--|
| | conda | | | | | |
| starAMR | Yes - conda | Yes | FASTA | Yes | Align reads to database | ResFinder, PointFinder, and PlasmidFinder |
| AMR Finder Plus | Yes - conda | Yes | FASTA | Yes | Align reads to database | Pathogen Detection Reference Gene Database |
| ResPipe | No | Yes | FASTQ or BAM | No | | |
| PointFinder | Yes - docker | Yes | FASTA | No | | |
| PCM: Pairwise Comparative Modelling | No | Yes | FASTA - protein | No | | |
| SRST2 | No | No | FASTQ | No | | |
| Arg_Ranker | Yes - conda | Yes | Requires special metadata input | No | | |
| MetaCherchant | Yes - conda | No | FASTA - | No | | |

| | | | | | | |
|----------------------------|--------------|----|------------------------|---|--|--|
| | | | genomic | | | |
| ARIBA | Yes - docker | No | Paired end FASTQ | No | | |
| ARG-ANN OT | No | No | Unclear | No | | |
| kmerresista nce | No | No | - | No | | |
| c-sstar | No | No | Unkno wn | No - could not track down github | | |

Table 3A: Summary Statistics for the high resistance data from hAMRoaster: These are the counts and metrics as calculated by the hAMRoaster pipeline. Formulas for all metrics are as follows:

$$\text{Specificity} = \text{TN} / (\text{TN} + \text{FP})$$

$$\text{Sensitivity} = \text{TP} / (\text{TP} + \text{FN})$$

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{FP} + \text{TN} + \text{FN})$$

$$\text{Proportion Unknown} = \text{unknown} / (\text{TP} + \text{FP} + \text{unknowns})$$

| High Resistance Data, 100x Coverage | | | | | | | | | | |
|-------------------------------------|----------------|---------------|---------|----------------|---------------|-------------|-------------|-----------|----------|--------------------|
| tool | False positive | True positive | unknown | False negative | True negative | sensitivity | specificity | precision | accuracy | Proportion Unknown |
| abricate | 0 | 66 | 22 | 9 | 2 | 0.8800 | 1.0000 | 1.0000 | 0.8831 | 0.2500 |
| amrfinder plus | 2 | 62 | 71 | 9 | 1 | 0.8732 | 0.3333 | 0.9688 | 0.8514 | 0.5259 |
| deeparg | 0 | 98 | 23 | 8 | 2 | 0.9245 | 1.0000 | 1.0000 | 0.9259 | 0.1901 |
| fARGene | 0 | 713 | 0 | 13 | 2 | 0.9821 | 1.0000 | 1.0000 | 0.9821 | 0.0000 |
| resfinder 4 | 1 | 43 | 15 | 9 | 1 | 0.8269 | 0.5000 | 0.9773 | 0.8148 | 0.2542 |
| rgi | 4 | 559 | 255 | 6 | 1 | 0.9894 | 0.2000 | 0.9929 | 0.9825 | 0.3117 |
| shortbred | 0 | 29 | 0 | 11 | 2 | 0.7250 | 1.0000 | 1.0000 | 0.7381 | 0.0000 |
| srax | 0 | 10 | 3 | 11 | 2 | 0.4762 | 1.0000 | 1.0000 | 0.5217 | 0.2308 |
| staramr | 1 | 52 | 11 | 9 | 1 | 0.8525 | 0.5000 | 0.9811 | 0.8413 | 0.1719 |
| High Resistance Data, 50x Coverage | | | | | | | | | | |
| tool | False positive | True positive | unknown | False negative | True negative | sensitivity | specificity | precision | accuracy | Proportion Unknown |

| | | | | | | | | | | |
|--|---------------------------|--------------------------|----------------|---------------------------|--------------------------|-------------------------|-------------------------|------------------|-----------------|-------------------------------|
| abricate | 0 | 66 | 21 | 9 | 2 | 0.8800 | 1.0000 | 1.0000 | 0.8831 | 0.2414 |
| amrfinder plus | 2 | 62 | 67 | 9 | 1 | 0.8732 | 0.3333 | 0.9688 | 0.8514 | 0.5115 |
| deeparg | 0 | 99 | 23 | 8 | 2 | 0.9252 | 1.0000 | 1.0000 | 0.9266 | 0.1885 |
| fARGene | 0 | 702 | 0 | 13 | 2 | 0.9818 | 1.0000 | 1.0000 | 0.9819 | 0.0000 |
| resfinder 4 | 1 | 43 | 15 | 9 | 1 | 0.8269 | 0.5000 | 0.9773 | 0.8148 | 0.2542 |
| rgi | 4 | 557 | 254 | 6 | 1 | 0.9893 | 0.2000 | 0.9929 | 0.9824 | 0.3117 |
| shortbred | 0 | 30 | 0 | 11 | 2 | 0.7317 | 1.0000 | 1.0000 | 0.7442 | 0.0000 |
| srax | 0 | 13 | 3 | 10 | 2 | 0.5652 | 1.0000 | 1.0000 | 0.6000 | 0.1875 |
| staramr | 1 | 52 | 11 | 9 | 1 | 0.8525 | 0.5000 | 0.9811 | 0.8413 | 0.1719 |
| High Resistance Data, 5x Coverage | | | | | | | | | | |
| tool | False positive | True positive | unknown | False negative | True negative | sensitivit y | specificit y | precision | accuracy | Proportion Unknown |

| | | | | | | | | | | |
|---------------------------|----|----|------|-----|---|--------|--------|--------|--------|--------|
| abricate | 0 | 9 | 39 | 19 | 2 | 0.8125 | 1.0000 | 1.0000 | 0.8200 | 0.3276 |
| amrfinder plus | 1 | 9 | 60 | 58 | 1 | 0.8696 | 0.5000 | 0.9836 | 0.8592 | 0.4874 |
| deeparg | 0 | 8 | 267 | 86 | 2 | 0.9709 | 1.0000 | 1.0000 | 0.9711 | 0.2436 |
| fARGene | 0 | 13 | 470 | 0 | 2 | 0.9731 | 1.0000 | 1.0000 | 0.9732 | 0.0000 |
| resfinder 4 | 0 | 9 | 43 | 10 | 2 | 0.8269 | 1.0000 | 1.0000 | 0.8333 | 0.1887 |
| rgi | 12 | 6 | 1015 | 418 | 1 | 0.9941 | 0.0769 | 0.9883 | 0.9826 | 0.2893 |
| shortbred | 0 | 11 | 29 | 0 | 2 | 0.7250 | 1.0000 | 1.0000 | 0.7381 | 0.0000 |
| srax | 0 | 12 | 4 | 3 | 2 | 0.2500 | 1.0000 | 1.0000 | 0.3333 | 0.4286 |
| staramr | 0 | 9 | 44 | 11 | 2 | 0.8302 | 1.0000 | 1.0000 | 0.8364 | 0.2000 |

Table 3B: Summary Statistics for the low resistance data from hAMRoaster: These are the counts and metrics as calculated by the hAMRoaster pipeline.

| input_file_name | AMR Isolate | tool | True Positive | False Positive | Unknown | True Negative | False Negative | Sensitivity | Specificity | Precision | Accuracy | Proportion Unknown |
|-------------------|-------------|---------------|---------------|----------------|---------|---------------|----------------|-------------|-------------|-----------|----------|--------------------|
| samp0_srr17789825 | SRR177898 | amrfinderplus | 0 | 0 | 4 | 5 | 4 | 1.0000 | 0.0000 | 0.0000 | 0.5556 | 1.0000 |
| samp0_srr17789825 | SRR177898 | deeparg | 0 | 0 | 7 | 5 | 4 | 1.0000 | 0.0000 | 0.0000 | 0.5556 | 1.0000 |
| samp0_srr17789825 | SRR177898 | resfinder 4 | 0 | 0 | 1 | 5 | 4 | 1.0000 | 0.0000 | 0.0000 | 0.5556 | 1.0000 |
| samp0_srr17789825 | SRR177898 | rgi | 0 | 0 | 22 | 5 | 4 | 1.0000 | 0.0000 | 0.0000 | 0.5556 | 1.0000 |

| | | | | | | | | | | | | |
|---------------|-----------|-------------|---|---|---|---|---|--------|--------|-------|-------|--------|
| samp0_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.555 | |
| 9825 | 29 | shortbred | 0 | 0 | 8 | 5 | 4 | 1.0000 | 0.0000 | 0 | 6 | 1.0000 |
| samp0_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.555 | |
| 9825 | 30 | srax | 0 | 0 | 6 | 5 | 4 | 1.0000 | 0.0000 | 0 | 6 | 1.0000 |
| samp1_srr1668 | SRR166836 | amrfinderpl | | | | | | | | 1.000 | 0.555 | |
| 3675 | 75 | us | 1 | 0 | 2 | 4 | 4 | 1.0000 | 0.2000 | 0 | 6 | 0.6667 |
| samp1_srr1668 | SRR166836 | | | | | | | | | 0.750 | 0.583 | |
| 3675 | 76 | deeparg | 3 | 1 | 2 | 4 | 4 | 0.8000 | 0.4286 | 0 | 3 | 0.3333 |
| samp1_srr1668 | SRR166836 | | | | | | | | | 1.000 | 0.555 | |
| 3675 | 77 | resfinder 4 | 1 | 0 | 2 | 4 | 4 | 1.0000 | 0.2000 | 0 | 6 | 0.6667 |
| samp1_srr1668 | SRR166836 | | | | | | | | | 0.142 | 0.333 | |
| 3675 | 78 | rgi | 1 | 6 | 7 | 4 | 4 | 0.4000 | 0.2000 | 9 | 3 | 0.5000 |
| samp1_srr1668 | SRR166836 | | | | | | | | | 0.000 | 0.500 | |
| 3675 | 79 | shortbred | 0 | 0 | 4 | 4 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |

| | | | | | | | | | | | | |
|---------------|-----------|-------------|---|---|----|----|---|--------|--------|-------|-------|--------|
| samp1_srr1668 | SRR166836 | | | | | | | | | 0.000 | 0.444 | |
| 3675 | 80 | srax | 0 | 1 | 5 | 4 | 4 | 0.8000 | 0.0000 | 0 | 4 | 0.8333 |
| samp1_srr1668 | SRR166836 | | | | | | | | | 1.000 | 0.600 | |
| 3675 | 81 | staramr | 2 | 0 | 0 | 4 | 4 | 1.0000 | 0.3333 | 0 | 0 | 0.0000 |
| samp2_srr1778 | SRR177898 | amrfinderpl | | | | | | | | 0.000 | 0.750 | |
| 9825 | 30 | us | 0 | 0 | 4 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp2_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 31 | deeparg | 0 | 0 | 6 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp2_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 32 | resfinder 4 | 0 | 0 | 1 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp2_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 33 | rgi | 0 | 0 | 22 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp2_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 34 | shortbred | 0 | 0 | 4 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |

| | | | | | | | | | | | | |
|---------------|-----------|-------------|---|---|----|----|---|--------|--------|-------|-------|--------|
| samp2_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 35 | srax | 0 | 0 | 3 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp3_srr1668 | SRR166836 | amrfinderpl | | | | | | | | 1.000 | 0.555 | |
| 3675 | 75 | us | 1 | 0 | 2 | 4 | 4 | 1.0000 | 0.2000 | 0 | 6 | 0.6667 |
| samp3_srr1668 | SRR166836 | | | | | | | | | 1.000 | 0.636 | |
| 3675 | 76 | deeparg | 3 | 0 | 3 | 4 | 4 | 1.0000 | 0.4286 | 0 | 4 | 0.5000 |
| samp3_srr1668 | SRR166836 | | | | | | | | | 1.000 | 0.555 | |
| 3675 | 77 | resfinder 4 | 1 | 0 | 2 | 4 | 4 | 1.0000 | 0.2000 | 0 | 6 | 0.6667 |
| samp3_srr1668 | SRR166836 | | | | | | | | | 1.000 | 0.555 | |
| 3675 | 78 | rgi | 1 | 0 | 13 | 4 | 4 | 1.0000 | 0.2000 | 0 | 6 | 0.9286 |
| samp3_srr1668 | SRR166836 | | | | | | | | | 1.000 | 0.600 | |
| 3675 | 79 | shortbred | 2 | 0 | 9 | 4 | 4 | 1.0000 | 0.3333 | 0 | 0 | 0.8182 |
| samp3_srr1668 | SRR166836 | | | | | | | | | 0.000 | 0.500 | |
| 3675 | 80 | srax | 0 | 0 | 5 | 4 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |

| | | | | | | | | | | | | |
|---------------|-----------|-------------|---|---|----|----|---|--------|--------|-------|-------|--------|
| samp3_srr1668 | SRR166836 | | | | | | | | | 1.000 | 0.600 | |
| 3675 | 81 | staramr | 2 | 0 | 0 | 4 | 4 | 1.0000 | 0.3333 | 0 | 0 | 0.0000 |
| samp4_srr1778 | SRR177898 | amrfinderpl | | | | | | | | 0.000 | 0.750 | |
| 9825 | 30 | us | 0 | 0 | 4 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp4_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 31 | deeparg | 0 | 0 | 6 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp4_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 32 | resfinder 4 | 0 | 0 | 1 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp4_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 33 | rgi | 0 | 0 | 22 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp4_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 34 | shortbred | 0 | 0 | 7 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |
| samp4_srr1778 | SRR177898 | | | | | | | | | 0.000 | 0.750 | |
| 9825 | 35 | srax | 0 | 0 | 6 | 12 | 4 | 1.0000 | 0.0000 | 0 | 0 | 1.0000 |

| | | | | | | | | | | | | |
|---------------|-----------|-------------|---|---|---|---|---|--------|--------|-------|-------|--------|
| samp5_srr1668 | SRR166836 | amrfinderpl | | | | | | | | 1.000 | 0.555 | |
| 3675 | 75 | us | 1 | 0 | 2 | 4 | 4 | 1.0000 | 0.2000 | 0 | 6 | 0.6667 |
| samp5_srr1668 | SRR166836 | | | | | | | | | 0.750 | 0.583 | |
| 3675 | 76 | deeparg | 3 | 1 | 2 | 4 | 4 | 0.8000 | 0.4286 | 0 | 3 | 0.3333 |
| samp5_srr1668 | SRR166836 | | | | | | | | | 1.000 | 0.555 | |
| 3675 | 77 | resfinder 4 | 1 | 0 | 2 | 4 | 4 | 1.0000 | 0.2000 | 0 | 6 | 0.6667 |
| samp5_srr1668 | SRR166836 | | | | | | | | | 0.142 | 0.333 | |
| 3675 | 78 | rgi | 1 | 6 | 7 | 4 | 4 | 0.4000 | 0.2000 | 9 | 3 | 0.5000 |
| samp5_srr1668 | SRR166836 | | | | | | | | | 0.000 | 0.444 | |
| 3675 | 79 | shortbred | 0 | 1 | 5 | 4 | 4 | 0.8000 | 0.0000 | 0 | 4 | 0.8333 |
| samp5_srr1668 | SRR166836 | | | | | | | | | 0.500 | 0.500 | |
| 3675 | 80 | srax | 1 | 1 | 3 | 4 | 4 | 0.8000 | 0.2000 | 0 | 0 | 0.6000 |
| samp5_srr1668 | SRR166836 | | | | | | | | | 1.000 | 0.600 | |
| 3675 | 81 | staramr | 2 | 0 | 0 | 4 | 4 | 1.0000 | 0.3333 | 0 | 0 | 0.0000 |

HAMROASTER

47

Supplementary text 1: URL link to tweet

https://twitter.com/emily_wissel/status/1336013892116488195

HAMROASTER

48

Supplementary table 1: tidy table of data

<https://docs.google.com/spreadsheets/d/1bfACqEh0nkS65vCUj5DfMg4PvW0fHxbtrv0PgKt1gT4/edit#gid=53644837>