
Large scale modeling of antimicrobial resistance with interpretable classifiers

Alexandre Drouin^{1,3,†}, Frédéric Raymond^{2,3}, Gaël Letarte St-Pierre^{1,3}, Mario Marchand^{1,3},
Jacques Corbeil^{2,3}, François Laviolette^{1,3}

¹ Department of Computer Science and Software Engineering, ² Infectious Disease Research Center, ³ Big Data Research Center
Université Laval, Québec, Canada

Abstract

Antimicrobial resistance is an important public health concern that has implications in the practice of medicine worldwide. Accurately predicting resistance phenotypes from genome sequences shows great promise in promoting better use of antimicrobial agents, by determining which antibiotics are likely to be effective in specific clinical cases. In healthcare, this would allow for the design of treatment plans tailored for specific individuals, likely resulting in better clinical outcomes for patients with bacterial infections. In this work, we present the recent work of Drouin et al. (2016) on using Set Covering Machines to learn highly interpretable models of antibiotic resistance and complement it by providing a large scale application of their method to the entire PATRIC database. We report prediction results for 36 new datasets and present the Kover AMR platform, a new web-based tool allowing the visualization and interpretation of the generated models.

1 Introduction

Modern medicine relies on antimicrobial drugs to treat infections. However, bacteria have evolved mechanisms to protect themselves from antibiotics [26, 19]. Thus, the use and abuse of antibiotics has led to the selection and to the spread of antibiotic-resistant pathogens. In order to define the right treatment and to reduce clinical failures, the prediction of antimicrobial resistance (AMR) is essential in choosing the right drugs to treat a specific patient [12]. In clinical laboratories, AMR is measured using antibiograms. This method determines the minimal inhibitory concentration (MIC) of an antibiotic by measuring the growth of the infecting microorganism in presence of different concentrations of the drug. The overall objective of this method is to determine if the pathogen will respond to treatment. Clinicians do so by comparing the measured MIC to the guidelines from CLSI or EUCAST, which are constantly being reevaluated by international committees [14]. However, while susceptibility to antibiotics can be predicted using MIC, it does not always hold true, as susceptible isolates can sometimes become phenotypically resistant given the proper conditions [17]. In this case, genomic determinants of resistance can be effective in predicting clinically relevant antimicrobial resistance [18]. In addition, determination of the MIC requires the growth of microorganisms, which generally necessitates one to two days of *in vitro* culture, and even more in the case of slow-growing organisms, such as *Mycobacterium tuberculosis* [28]. Genomic methods, such as polymerase chain reaction or whole genome sequencing, can now be used to predict the resistance phenotypes of pathogens in a more rapid manner [20].

Reanalysis of publicly available genome databases for which AMR phenotypes are available is a useful starting point to improve our understanding of the relationship between genotype and phenotype. Indeed, several groups have used machine learning and statistics to understand and predict antimicrobial resistance ([2, 11, 22, 10, 9]). However, models that predict AMR should not be static and should improve as new genomes are added to databases. For instance, the Pathosystems Resource Integration Center (PATRIC) database is a large-scale aggregation platform for bacterial genomes and their associated metadata [27, 22]. The number of genomes in the PATRIC database nearly doubled between 2014 and 2015, with now more than 52 thousand microbial genome sequences. Recently, Drouin et al. (2016) proposed to use the Set Covering Machine [16] to learn extremely sparse models of antimicrobial resistance that are intelligible for domain experts [10]. They compared their models to more complex predictors, such as linear and kernel-based Support Vector Machines [7, 23, 24], as well as decision trees [3], and showed that Set Covering Machines achieved comparable, often superior, generalization performance, while being significantly sparser. Moreover, they demonstrated that highly accurate models of AMR could be obtained, despite features spaces tens of thousands of times larger than the number of learning examples.

[†]Corresponding author: alexandre.drouin.8@ulaval.ca

Peer-reviewed and accepted for presentation at the Machine Learning for Health Workshop, NIPS 2016, Barcelona, Spain.

The present work summarizes the work of Drouin et al. (2016) and presents extensive new results. Specifically, we present a large scale application of their method to prospectively generate predictive models of AMR from the ever-growing collection of genomes in the PATRIC database. Moreover, we present the Kover AMR Platform (<https://aldro61.github.io/kover-amr-platform/>), a web-based tool that catalogs AMR prediction results for a wide variety of species and antibiotics, providing detailed metrics and allowing the visualization of the generated models. This initiative will allow the interpretation of our results by healthcare researchers, generating new research and treatment opportunities.

2 Methods

2.1 Problem statement

We address the problem of predicting antimicrobial resistance as a supervised learning problem. The goal is to learn a model that accurately discriminates genomes that are resistant or susceptible to an antibiotic based on genomic characteristics. Formally, we assume that we are given a dataset $S \stackrel{\text{def}}{=} \{(x_i, y_i)\}_{i=1}^m \sim D^m$, where $x_i \in \{A, C, G, T\}^*$ is a bacterial genome, $y_i \in \{0, 1\}$ is its associated phenotype (0 for susceptible and 1 for resistant) and D is an unknown data generating distribution from which the dataset is sampled. We start by defining an alternative representation for the genomic sequences, where each genome is characterized by the presence or absence of every possible k -mer, i.e. every possible sequence of k DNA nucleotides. This representation is obtained through a mapping function $\phi : \{A, C, G, T\}^* \rightarrow \{0, 1\}^{4^k}$, such that $\phi_j(x) \stackrel{\text{def}}{=} 1$ if the k -mer k_j is in the genome x and 0 otherwise. This yields the transformed dataset $S' \stackrel{\text{def}}{=} \{(\phi(x_i), y_i)\}_{i=1}^m$, which is then used to train the learning algorithm.

The goal is then to find a model h that has a good generalization performance, i.e. that minimizes the probability $R(h)$ of making a prediction error for any example drawn according to distribution D , i.e.,

$$R(h) \stackrel{\text{def}}{=} \Pr_{(x,y) \sim D} [h(\phi(x)) \neq y]. \quad (1)$$

Furthermore, we seek highly interpretable models from which biologically relevant knowledge can be extracted.

2.2 The Set Covering Machine

Such interpretable models are obtained through the Set Covering Machine algorithm (SCM) [16, 10], which produces models that are logical combinations (conjunctions or disjunctions) of boolean-valued rules that are generated from the data. We now briefly present the algorithm and direct the reader to [10, 16] for further explanations.

The input of the SCM algorithm is a set of learning examples S and a set of boolean-valued rules \mathcal{R} . In our context, S is composed of genomes, in the k -mer form induced by ϕ , and their labels. Let \mathcal{K} be the set of all, possibly overlapping, k -mers that are present in at least one genome of S . For each k -mer $k_j \in \mathcal{K}$, we consider a presence rule, defined as $p_{k_j}(\phi(x)) \stackrel{\text{def}}{=} I[\phi_j(x) = 1]$ and an absence rule, defined as $a_{k_j}(\phi(x)) \stackrel{\text{def}}{=} I[\phi_j(x) = 0]$, where $I[\text{true}] = 1$ and 0 otherwise. These boolean-valued rules constitute the set \mathcal{R} . Given S and \mathcal{R} , the SCM attempts to find the model that relies on the smallest possible set of rules $\mathcal{R}^* = \{r_1^*, \dots, r_n^*\} \subseteq \mathcal{R}$, while minimizing Equation (1). The models generated can be conjunctions $h(\phi(x)) \stackrel{\text{def}}{=} r_1^*(\phi(x)) \wedge \dots \wedge r_n^*(\phi(x))$ or disjunctions $h(\phi(x)) \stackrel{\text{def}}{=} r_1^*(\phi(x)) \vee \dots \vee r_n^*(\phi(x))$. Hence, they directly highlight the importance of a small set of genomic sequences for predicting AMR phenotypes.

However, it is important to note that the distribution D is unknown; therefore, it is not possible to directly minimize Equation (1). Instead, the algorithm constructs a model that achieves an appropriate trade-off between the empirical risk (i.e., the fraction of training errors) and the number of rules it uses. A model containing many rules is likely to overfit the data, whereas a model containing too few rules is likely to underfit. To find the appropriate trade-off between the classifier’s size and its accuracy on the training set, the SCM relies on a modified version of the set covering greedy algorithm of Chvátal [6], which has a worst-case guarantee. The running time and space complexities of this algorithm are linear in the number of examples and rules, thus linear in the number of genomes and k -mers. Consequently, this algorithm is particularly well-suited for learning from large datasets of extreme dimensionalities, such as the ones that often occur in healthcare applications.

The experiments in this work were performed using Kover, the SCM implementation of Drouin et al. (2016), which has been tailored for learning from k -mer represented genomes [10]. In Kover, the SCM algorithm is trained *out-of-core*, which means that the data resides on the disk and is accessed in blocks by the algorithm. The implementation exploits a compressed data representation and atomic CPU instructions to speed up computations. It is open-source and available at <https://github.com/aldro61/kover>.

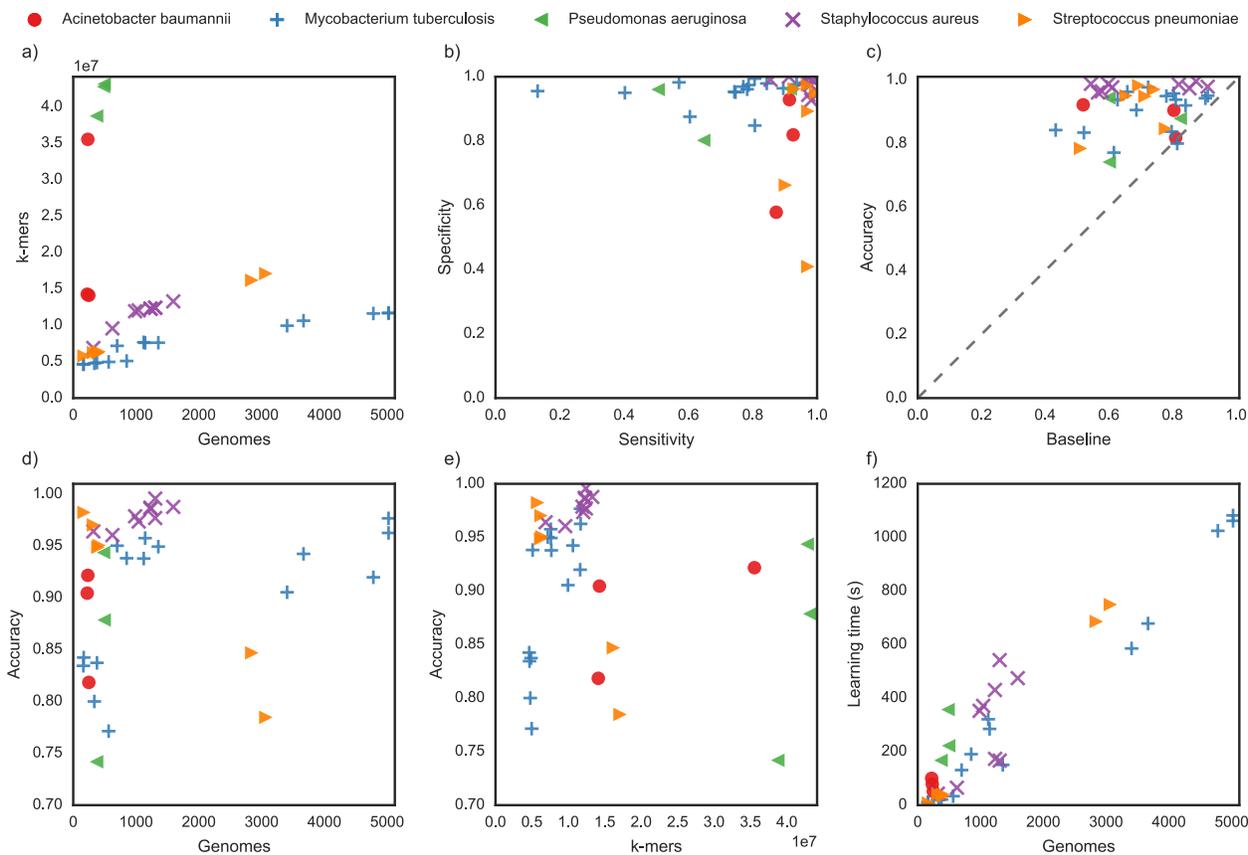


Figure 1: Each point corresponds to a dataset. Colored symbols identify bacterial species. a) Size of the datasets in terms of genomes and k -mers; b) Specificity of the models with respect to their sensitivity; c) Accuracy of the models with respect to the baseline predictor (see text for details); d) Accuracy of the models with respect to the number of genomes; e) Accuracy of the models with respect to the number of k -mers; f) Time required to train the algorithm with respect to the number of genomes.

2.3 Data acquisition

The data used in this study were extracted from the PATRIC database via its FTP server (<ftp.patricbrc.org>). First, the latest AMR metadata were acquired. These data consisted of genome identifiers and measured resistance phenotypes (resistant or susceptible) for various antibiotics. The data were segmented by species and antibiotic to form datasets, and those with at least 50 genomes in each class were retained. Finally, for each dataset, the assembled genomes were downloaded and their k -mer representation was obtained using the DSK k -mer counter [21]. The value of k was set to 31, since extensive experimentation has shown that this value is appropriate for the task at hand [10].

3 Results and discussion

A total of 36 datasets, comprising five bacterial species, were extracted from PATRIC. A detailed list is available in Appendix - Table 1. Figure 1a shows the size of these datasets, in which the number of examples ranged from 161 to 5022 and the number of k -mers from 4.6 to 43.1 millions. For most species, an increase in the number of genomes does not lead to a large increase in the number of k -mers. This reflects the limited genomic diversity that exists within species (e.g.: *M. tuberculosis* [13]). Contrasting results were observed for *P. aeruginosa* and *A. baumannii*, which have higher genomic diversities [15, 25].

The SCM's ability to predict AMR phenotypes was evaluated on held-out subsets of the data. First, each dataset was randomly partitioned into a training set \mathcal{S} (80%) and a testing set \mathcal{T} (20%). The SCM was then trained on \mathcal{S} and metrics were computed based on the model's predictions on \mathcal{T} . This procedure was repeated 10 times, on different partitions,

and the resulting metrics were averaged. The values of the algorithm’s hyperparameters¹ (HPs) were selected by bound selection (see [16, 10]). Bound selection uses a probabilistic upper bound on Equation (1), computed from the training data, to score each of the HP combinations. For each of the latter, a single training of the algorithm is required; hence, bound selection is much faster than standard cross-validation. Drouin et al (2016). proposed such a bound for conjunctions/disjunctions of presence/absence rules of k -mers and found that bound selection yielded results comparable to those of 5-fold cross-validation.

The results are summarized in Figure 1b-e and detailed in Appendix - Table 2. Figure 1b shows the specificity of the models with respect to their sensitivity. A perfect model would score 1 for each of these metrics. Specificities superior to 80% were achieved for 33/36 datasets and comparable sensitivities were achieved for 25/36 datasets. In general, the obtained models are more specific than sensitive, meaning that they sometimes fail to identify resistant isolates, but very rarely mark a susceptible isolate as resistant. Figure 1c compares the SCM models to a baseline predictor that predicts the most abundant class in the training data. The SCM models, which achieve accuracies greater than 80% on 33/36 datasets, generally surpass the baseline predictors, indicating that the algorithm extracts relevant patterns of antibiotic resistance. Of note, the models learned by the SCM are extremely sparse, using an average of 2.5 rules (std: 2.2), which makes them well-suited for further review and experimental validation. Moreover, Figure 1c highlights the strong class imbalance that exists in some of the datasets considered, as the baseline predictors often achieve high accuracies. Furthermore, based on Figure 1d, we observe that the accuracy of the models is generally higher for datasets that contain more examples. There are notable exceptions, such as *S. pneumoniae*, where the accuracies are lower for larger datasets. However, the two largest datasets for this species correspond to combinations of drugs (Beta-lactams and Trimethoprim/Sulfamethoxazole), which could complexify the learning task. Also, notice that there are cases where the algorithm achieves near perfect accuracies with very few examples, despite disproportionately large feature spaces. In fact, Figure 1e illustrates that the accuracy of the models is not related to the number of k -mers. This supports the theoretical and empirical results of Drouin et al. (2016), which found that the SCM could avoid overfitting, even in such high dimensional settings. Finally, Figure 1f shows the time required to train the algorithm with respect to the number of genomes in each dataset. The time, which grows linearly with the number of genomes, varied between 8 seconds and 18 minutes, using a single CPU core and less than 1 GB of memory.

The models generated using the SCM are available through the Kover AMR Platform. The k -mer sequences of the rules in these models were annotated using BLAST [1]. This revealed that, for most antibiotics, the SCM was accurate in identifying known resistance mechanisms. For example, the absence of one k -mer located in the DNA gyrase subunit A (*gyrA*) predicts resistance to moxifloxacin in *M. tuberculosis* with an error rate as low as 4%. This k -mer refers to the amino-acids 88 to 94 of GyrA, the mutation of which confers resistance to fluoroquinolones [5]. Thus, the model relies on the absence of the susceptible genotype to account for the presence of resistant variants, which are more diverse. This behavior was also observed for predicting resistance to isoniazid in *M. tuberculosis*, where the model rightly targets a region of the *katG* gene where multiple mutations are known to induce resistance [4, 8]. The model also relies on a k -mer in the *rpoB* gene, a known rifampicin resistance determinant [8]. This could result from the frequent combined use of antituberculosis drugs.

We have briefly demonstrated the accuracy and interpretability of models produced by the SCM. Detailed results are available in the Appendix and the Kover AMR platform, which allows the visualization and further investigation of AMR models generated at an unprecedented scale.

4 Conclusion

In summary, we have outlined the recently published work of Drouin et al. (2016), while complementing their analysis with a large-scale application of their method to the ever-growing PATRIC database. Kover, an extremely efficient out-of-core implementation of the SCM, allowed the rapid generation of these results with limited computational resources. Moreover, our results show that the method yields accurate results for predicting AMR in most datasets and that, due to their strong interpretability, the obtained models can generate biologically relevant insight into these phenotypes.

On another hand, contrasting results were obtained, which pave the way to extensions of their method. In fact, the obtained models were generally highly specific, but some lacked sensitivity. This could result from seeking the sparsest model that detects resistant genomes in the entire population of isolates. Hence, the deconvolution of resistance mechanisms based on population structure could provide a deeper understanding of antibiotic resistance and increase the sensitivity of the models. Future work will therefore involve the development of algorithmic extensions to the Set Covering Machine, which allow the inclusion of prior knowledge of the population structure and the biological structures present in the data (e.g., gene functions, pathways).

¹ $p \in \{0.1, 0.178, 0.316, 0.562, 1, 1.778, 3.162, 5.623, 10, +\infty\}$, $s \in \{1, \dots, 10\}$, $\text{model.type} \in \{\text{conjunction}, \text{disjunction}\}$ (notation of [10])

We have only scratched the surface of the biological knowledge that can be generated from these results. The fact that our approach generates interpretable predictors, together with our proposed Kover AMR Platform (<https://aldro61.github.io/kover-amr-platform/>) will allow further analysis of these results by researchers with diverse backgrounds, bridging the gap between machine learning and healthcare research.

Acknowledgements

The authors acknowledge Sébastien Giguère and Pier-Luc Plante for helpful comments and suggestions. Computations were performed on the Colosse supercomputer at Université Laval (resource allocation project: nne-790-ae), under the auspices of Calcul Québec and Compute Canada. AD is recipient of an Alexander Graham Bell Canada Graduate Scholarship Doctoral Award of the National Sciences and Engineering Research Council of Canada (NSERC). This work was supported in part by the NSERC Discovery Grants (FL; 262067, MM; 122405). JC acknowledges the Canada Research Chair in Medical Genomics.

References

- [1] Stephen F Altschul et al. “Basic local alignment search tool”. In: *Journal of Molecular Biology* 215.3 (Oct. 1990), pp. 403–410.
- [2] Phelim Bradley et al. “Rapid antibiotic-resistance predictions from genome sequence data for *Staphylococcus aureus* and *Mycobacterium tuberculosis*”. In: *Nature communications* 6 (Dec. 2015), p. 10063.
- [3] Leo Breiman et al. *Classification and regression trees*. New York: CRC press, 1984.
- [4] Christine E Cade et al. “Isoniazid-resistance conferring mutations in *Mycobacterium tuberculosis* KatG: Catalase, peroxidase, and INH-NADH adduct formation activities”. In: *Protein Science* 19.3 (2010), pp. 458–474.
- [5] Jung-Yien Chien et al. “Mutations in *gyrA* and *gyrB* among Fluoroquinolone- and Multidrug-resistant *Mycobacterium tuberculosis* Isolates”. In: *Antimicrobial agents and chemotherapy* 60.4 (2016), pp. 2090–2096.
- [6] V Chvátal. “A Greedy Heuristic for the Set-Covering Problem”. In: *Mathematics of Operations Research* 4.3 (Aug. 1979), pp. 233–235.
- [7] Corinna Cortes and Vladimir Vapnik. “Support-vector networks”. In: *Machine learning* 20.3 (1995), pp. 273–297.
- [8] Pedro Eduardo Almeida Da Silva and Juan Carlos Palomino. “Molecular basis and mechanisms of drug resistance in *Mycobacterium tuberculosis*: classical and new drugs”. In: *Journal of antimicrobial chemotherapy* 66.7 (2011), pp. 1417–1430.
- [9] James J Davis et al. “Antimicrobial resistance prediction in PATRIC and RAST”. In: *Scientific reports* 6 (2016).
- [10] Alexandre Drouin et al. “Predictive computational phenotyping and biomarker discovery using reference-free genome comparisons”. In: *BMC Genomics* 17.1 (2016), pp. 1–15.
- [11] Sarah G Earle et al. “Identifying lineage effects when controlling for population structure improves power in bacterial association studies”. In: *Nature Microbiology* 1 (2016), p. 16041.
- [12] N Deborah Friedman, Elizabeth Temkin, and Yehuda Carmeli. “The negative impact of antibiotic resistance”. In: *Clinical Microbiology and Infection* 22.5 (2016), pp. 416–422.
- [13] Ruth Hershberg et al. “High functional diversity in *Mycobacterium tuberculosis* driven by genetic drift and human demography”. In: *PLoS Biol* 6.12 (2008), e311.
- [14] Gunnar Kahlmeter. “The 2014 Garrod Lecture: EUCAST—are we heading towards international agreement?”. In: *Journal of Antimicrobial Chemotherapy* 70.9 (2015), pp. 2427–2439.
- [15] Veronica N Kos et al. “The resistome of *Pseudomonas aeruginosa* in relationship to phenotypic susceptibility.” In: *Antimicrobial Agents and Chemotherapy* 59.1 (Jan. 2015), pp. 427–436.
- [16] Mario Marchand and John Shawe-Taylor. “The set covering machine”. In: *The Journal of Machine Learning Research* 3 (2002), pp. 723–746.
- [17] Francis Martineau et al. “Correlation between the Resistance Genotype Determined by Multiplex PCR Assays and the Antibiotic Susceptibility Patterns of *Staphylococcus aureus* and *Staphylococcus epidermidis*”. In: *Antimicrobial Agents and Chemotherapy* 44.2 (2000), pp. 231–238.
- [18] Andrew G McArthur and Gerard D Wright. “Bioinformatics of antimicrobial resistance in the age of molecular epidemiology”. In: *Current opinion in microbiology* 27 (2015), pp. 45–50.
- [19] Jose M Munita, Arnold S Bayer, and Cesar A Arias. “Evolving resistance among Gram-positive pathogens”. In: *Clinical Infectious Diseases* 61.suppl 2 (2015), S48–S57.

- [20] NV Punina et al. “Whole-genome sequencing targets drug-resistant bacterial infections”. In: *Human genomics* 9.1 (2015), p. 1.
- [21] Guillaume Rizk, Dominique Lavenier, and Rayan Chikhi. “DSK: k-mer counting with very low memory usage”. In: *Bioinformatics* (2013), btt020.
- [22] John W Santerre et al. “Machine Learning for Antimicrobial Resistance”. In: *arXiv.org* (July 2016). arXiv:1607.01224v1.
- [23] Bernhard Schölkopf, Koji Tsuda, and Jean-Philippe Vert. *Kernel methods in computational biology*. Cambridge, Massachusetts: MIT press, 2004.
- [24] John Shawe-Taylor and Nello Cristianini. *Kernel methods for pattern analysis*. Cambridge, United Kingdom: Cambridge university press, 2004.
- [25] Lalena Wallace et al. “Use of Comparative Genomics To Characterize the Diversity of *Acinetobacter baumannii* Surveillance Isolates in a Health Care Institution”. In: *Antimicrobial Agents and Chemotherapy* 60.10 (2016), pp. 5933–5941.
- [26] Christopher Walsh. “Molecular mechanisms that confer antibacterial drug resistance”. In: *Nature* 406.6797 (2000), pp. 775–781.
- [27] Alice R Wattam et al. “PATRIC, the bacterial bioinformatics database and analysis resource”. In: *Nucleic acids research* (2013), gkt1099.
- [28] Adam A Witney et al. “Clinical use of whole genome sequencing for *Mycobacterium tuberculosis*”. In: *BMC medicine* 14.1 (2016), p. 1.