

Aims & Background

- This pilot study aims to utilise a novel machine learning (ML) platform, VariantSpark (VS)¹, to investigate genetic variations contributing to *Neisseria gonorrhoeae* (NG) antimicrobial resistance (AMR).
- NG is the second-most prevalent sexually-transmitted bacterial pathogen. It is listed as a WHO 'high-priority' resistance organism, due to its current high AMR distribution and historical ability to develop AMR².
- Applying ML to genomic data can improve identification of variants important for surveillance and diagnostics.

Methods

- Downloaded NG sequence and resistance data for 314 patients from a New Zealand publication³.
- Processed Illumina reads to generate variant call files (VCFs) with NZ_AP023069.1 as reference.
- Input VCFs into VS for random forest analysis and Hail2.0 for Firth's logistic regression analysis as a comparison.
- Utilised tetracycline AMR data because of its equal distribution of resistant and susceptible samples.

Results & Discussion

- VS (Figure 2) provided high importance scores for known tetracycline allele variants⁴: *rpsJ* (tetracycline target) and *porB* (membrane pore). However, *mtrCDE* complex genes (5 genes for membrane efflux-pump) had low scores (<0.05).
- Interestingly, allele variants in two chaperone proteins (protein stabilisers), *hscA* and *hemW*, also had high importance scores. While these are not currently known to cause AMR, other chaperones have been linked to NG resistance⁵.
- Compared to logistic regression (Figure 3), VS had less 'noise', however both assigned high values to 'hypothetical' proteins.

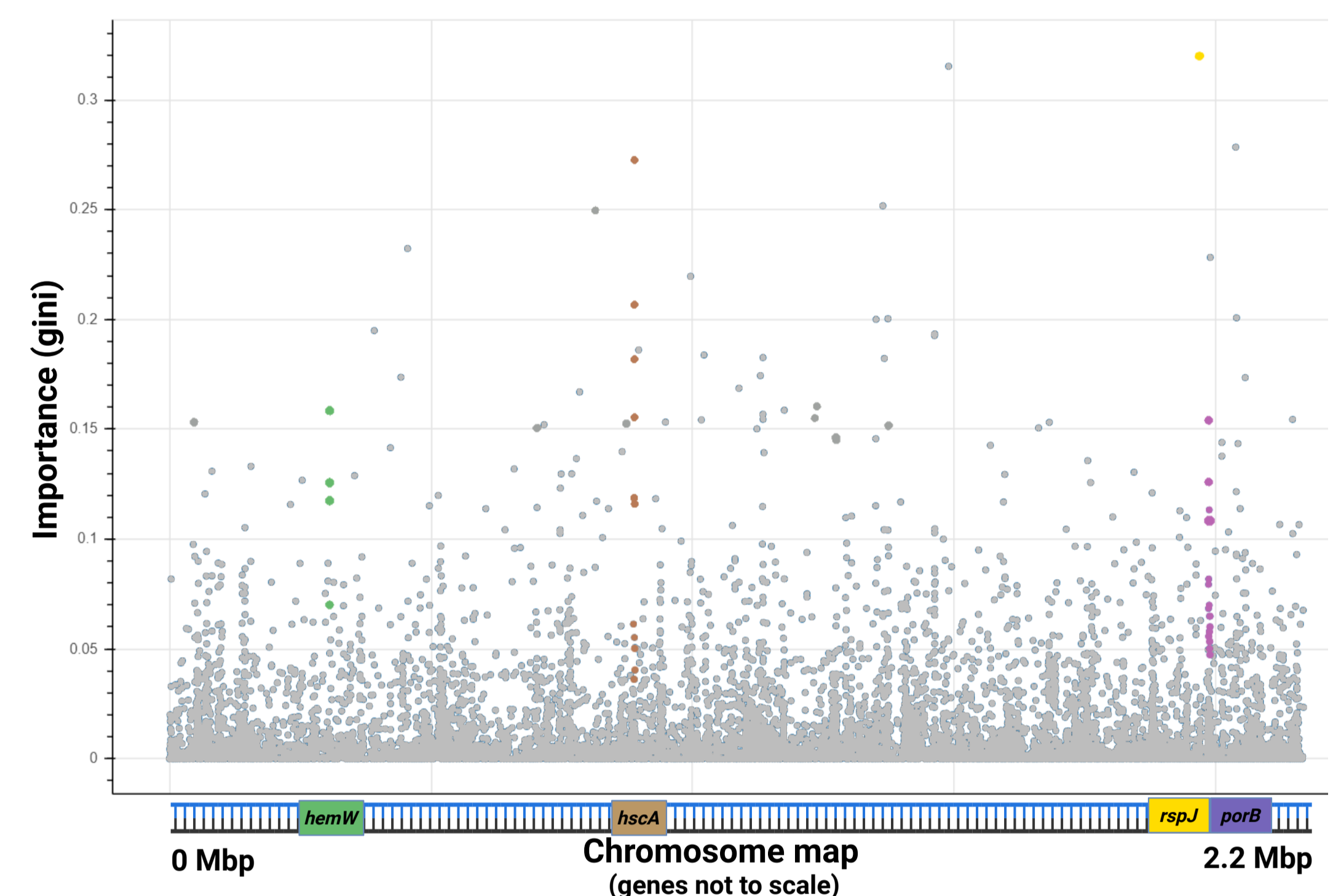


Figure 2. VariantSpark random forest importance score assignments to allele variants. Each dot signifies a unique allele variant.

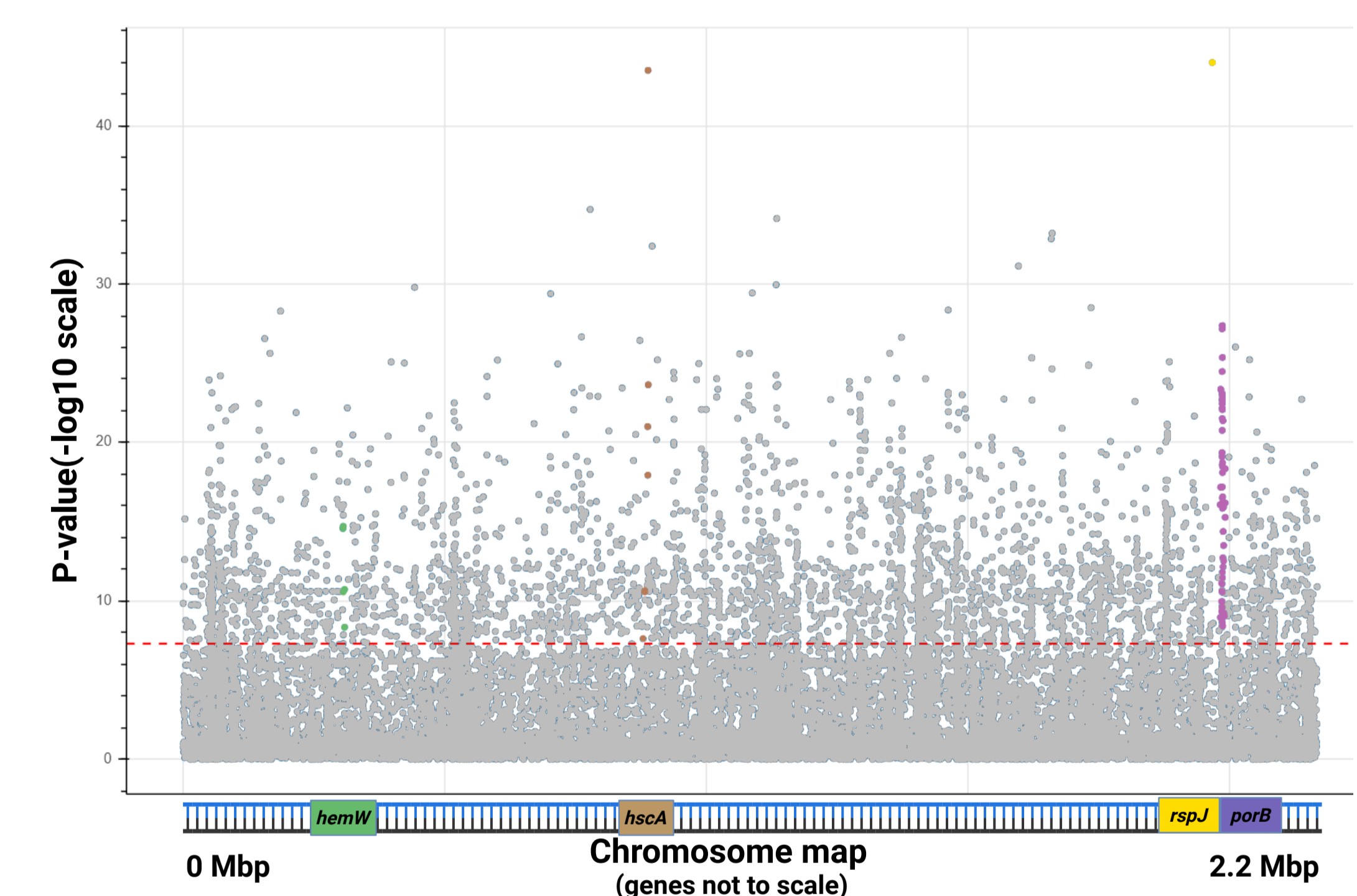


Figure 3. Hail2.0 logistic regression P-value assignments to allele variants. Each dot signifies a unique allele variant.

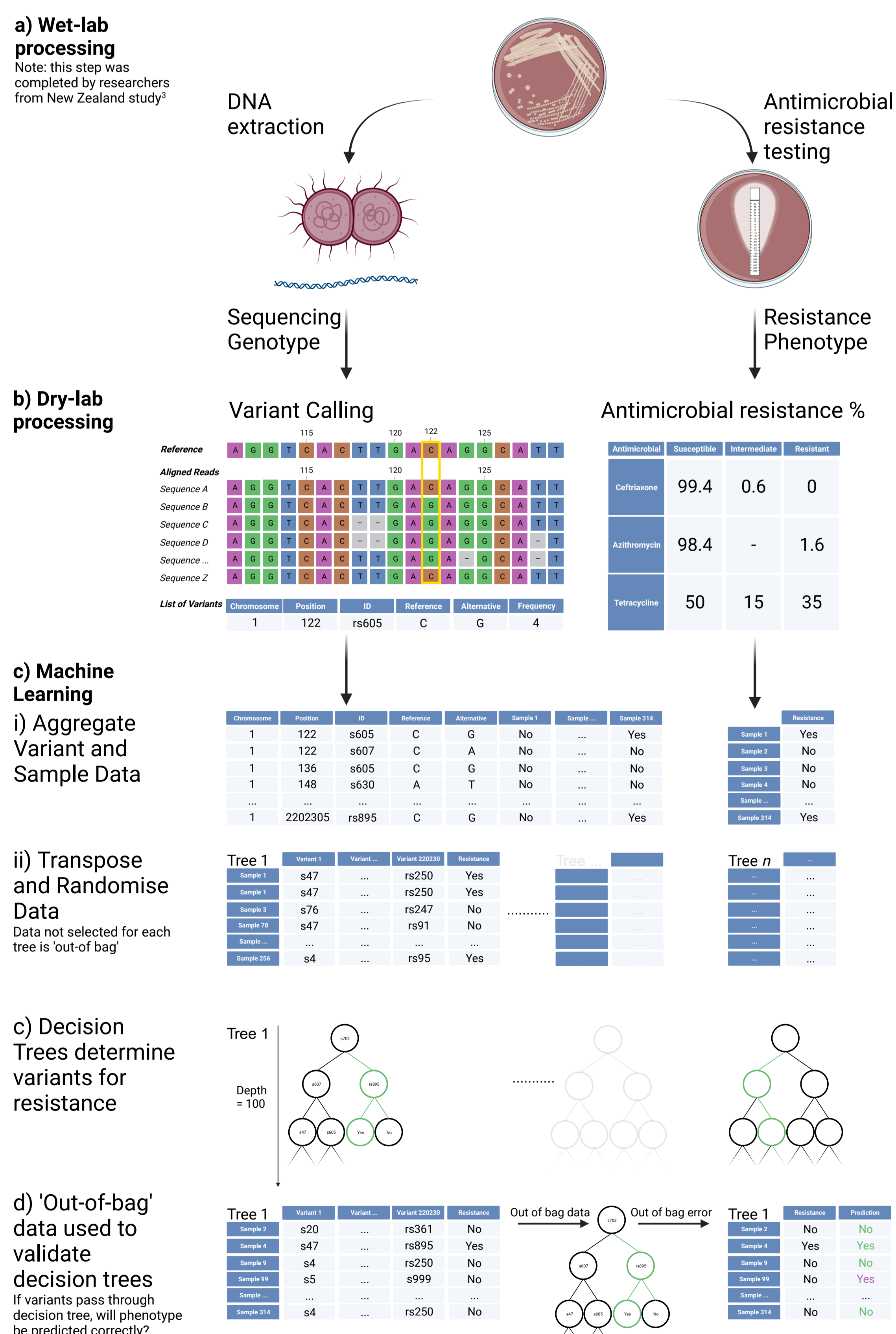


Figure 1. Pipeline for preparing data and executing VariantSpark random forest prediction. Example numbers shown, splits count per node and variable weight not covered.

In conclusion, machine learning successfully determined known and potential novel resistance genes.

Future Directions

- Integrate multiple, public NG datasets for VS analysis. Importantly, combining large datasets offers varied resistance profiles, thus allowing prediction for other antimicrobials.
- Incorporate patient data (e.g., gender) into VS analysis to examine these as potential confounding factors.

