

Title: Large-Scale Evaluation of AST Prediction using Resistance Marker Presence/Absence vs. Machine Learning on WGS Data

Authors: Arolyn Conwill, Mohamad Sater, Nick Worley Jason Wittenbach, Miriam H. Huntley

Background

Whole genome sequencing (WGS) of bacterial pathogens directly from clinical samples or cultured isolates is a promising technique for diagnosing infections, yet challenges remain in using WGS to guide antimicrobial therapy. The standard approach predicts antimicrobial susceptibility test (AST) profiles using the presence/absence of known resistance markers. Here we measure the performance of resistance markers compared to a machine learning method for predicting AST.

Methods

We assessed ResFinder, a publicly available bioinformatics tool that detects a curated set of resistance genes and point mutations, to test the utility of resistance marker presence or absence for predicting resistance/susceptibility on 107 species/drug combinations relevant for bloodstream infections. We also tested the performance of Keynome *g*AST, a machine learning method we developed that predicts AST from WGS data using the entirety of the bacterial genome without limiting to known resistance markers. Performance was assessed on >36,000 bacterial strains from MicrohmDB®, a large-scale database containing WGS data and phenotypic AST results for tens of thousands of clinical bacterial isolates.

Results

ResFinder performance varied widely across species/drug combinations, with a median balanced accuracy of 80% (52%-92%, 1st and 3rd quartiles). ResFinder detected a resistance marker in only 72.3% of resistant and 50.3% of intermediate sample-drug combinations, indicating that many phenotypically nonsusceptible organisms may be mischaracterized by relying on resistance marker presence. Conversely resistance markers were present in 14.2% of susceptible sample-drug combinations. In contrast, across the same species/drug combinations Keynome *g*AST had a median balanced accuracy of 92% (87%-96%, 1st and 3rd quartiles), and was superior to ResFinder in 69% and equivalent to ResFinder in 28% of species/drug combinations.

Conclusions

The lack of robust correlation between resistance markers and true phenotype highlights the limitations of predicting resistance/susceptibility based on the use of curated resistance markers. A more flexible machine learning approach that accesses the entire bacterial genome allows for complex combinations of genomic regions to inform AST prediction.