

Abstract

Drug resistance in bacterial pathogens is an increasing problem, which stimulates research. In this work, in order to deepen our understanding of drug resistance mechanisms, we investigate the approach of using whole-genome sequences to identifying genetic mutations associated with drug resistance phenotypes in bacterial strains.

In particular, we present GWAMAR, the tool we have developed to support this type of analysis. As a part of this work, we also present **weighted support (WS)** and **tree-generalized hypergeometric (TGH)** score — two statistics we propose for identifying of drug resistance associations, based on phylogenetic information. Additionally, we propose a **rank-based metascore (RBM)** for combining multiple scores into one in order to compromise between different approaches used to define different scores. We present results obtained by applying GWAMAR to two datasets for *M. tuberculosis*, which demonstrate that GWAMAR can be successfully used for identification of drug resistance-associated mutations.

The software, input datasets and results are provided at the website of our project, <http://bioputer.mimuw.edu.pl/gwamar>.

Methodology of GWAMAR

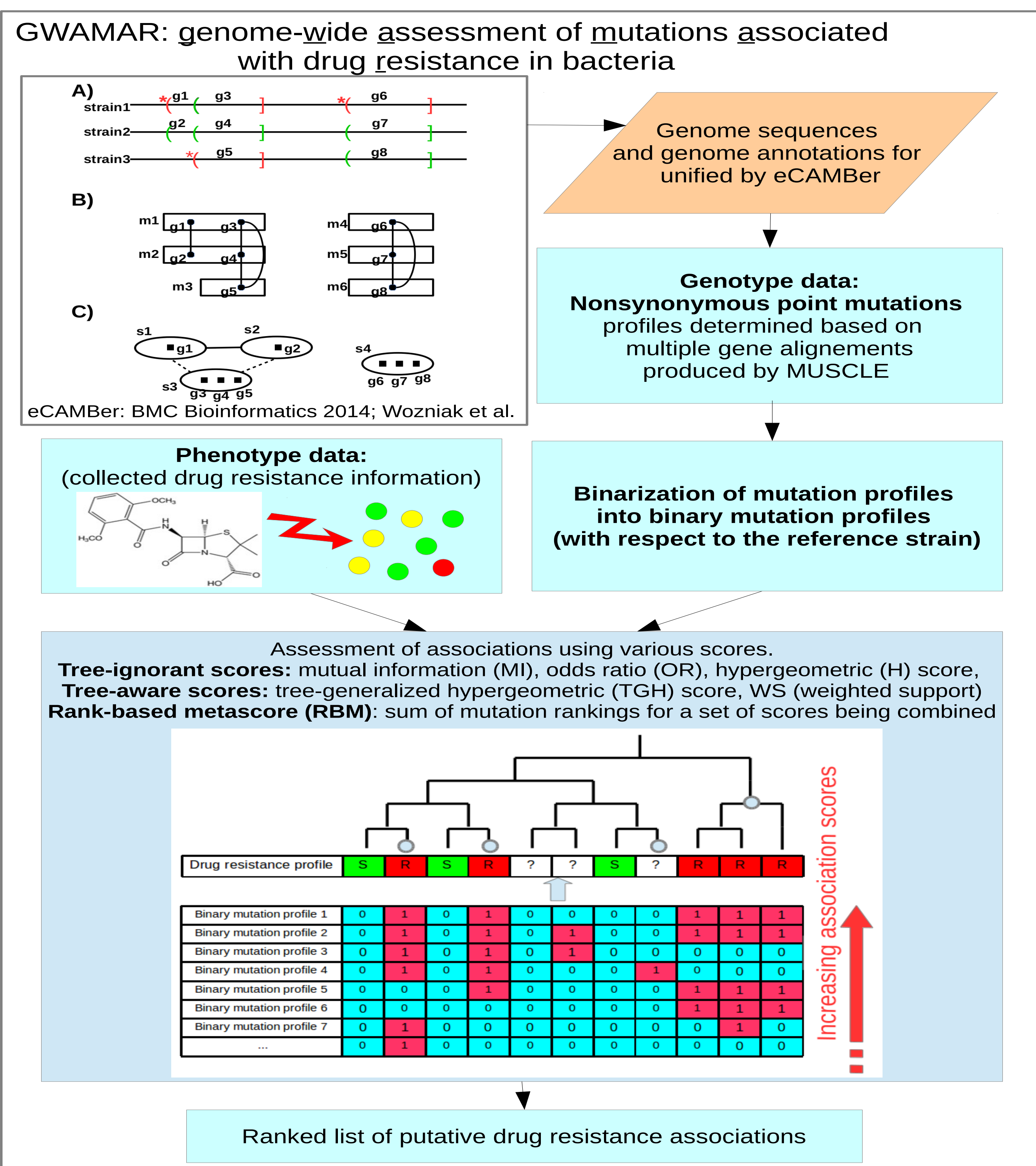
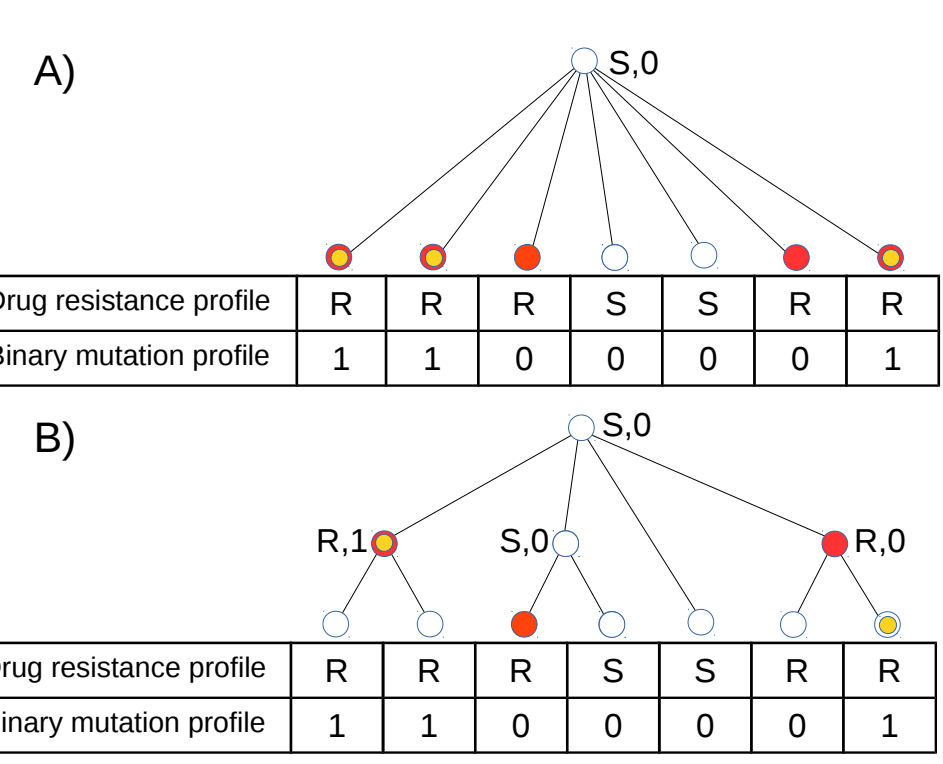


Figure 1: Schema of the pipeline of GWAMAR. For a set of considered bacterial strains, the input data for GWAMAR consists of (i) a set of mutations; (ii) a set of drug resistance profiles; and (iii) optional, phylogenetic tree for the set of bacterial strains. Typically the set of mutation profiles is generated using eCAMBER, which is able to download the genome sequences and annotations for the set of bacterial strains, identify point mutations based on multiple alignments, and reconstruct the phylogenetic tree of the considered bacterial strains. Assuming the genotype data is preprocessed, the first step of GWAMAR is to compute binary mutation profiles for all the mutations. This step significantly reduces the number of profiles considered. Finally, GWAMAR implements several statistical scores to associate drug resistance profiles with mutation profiles. These include: mutual information (MI), odds ratio (OR), hypergeometric (H) score, weighted support (WS), tree-generalized hypergeometric (TGH) score and the rank-based metascore (RBM). As a result, we obtain ordered lists of drug resistance associations, where the top-scored associations are the most likely to be real.

TGH score



A subset c of a tree nodes is a coloring, if it satisfies the following two conditions: (i) each path from a leaf to the root contains at most one node from c ; (ii) each internal node in T has at least one immediate child node which does not belong to c .

(A) an example of a pair of a drug resistance profile and a binary mutation profile. Values of the corresponding tree-extended binary mutation profile, and the corresponding tree-extended drug resistance profile are shown next to the nodes. (B) colorings \tilde{c} and \hat{c} induced by the same pair of profiles but for a flat tree.

For a drug resistance profile r and a binary mutation profile b , we denote the colorings induced by the profiles as \tilde{c} and \hat{c} , respectively. Then, we define the TGH score as follows:

$$TGH_T(r, b) = -\log \left(\frac{\sum_{i=k}^n B_{T, \tilde{c}}(i, n)}{W_T(n)} \right). \quad (1)$$

Here, $W_T(n)$ denotes the total number of colorings of T of size n , whereas $B_{T, \tilde{c}}(i, n)$ denotes the total number of colorings of T of size n , such that exactly i of their nodes are visible from coloring \tilde{c} .

Assessment of accuracy

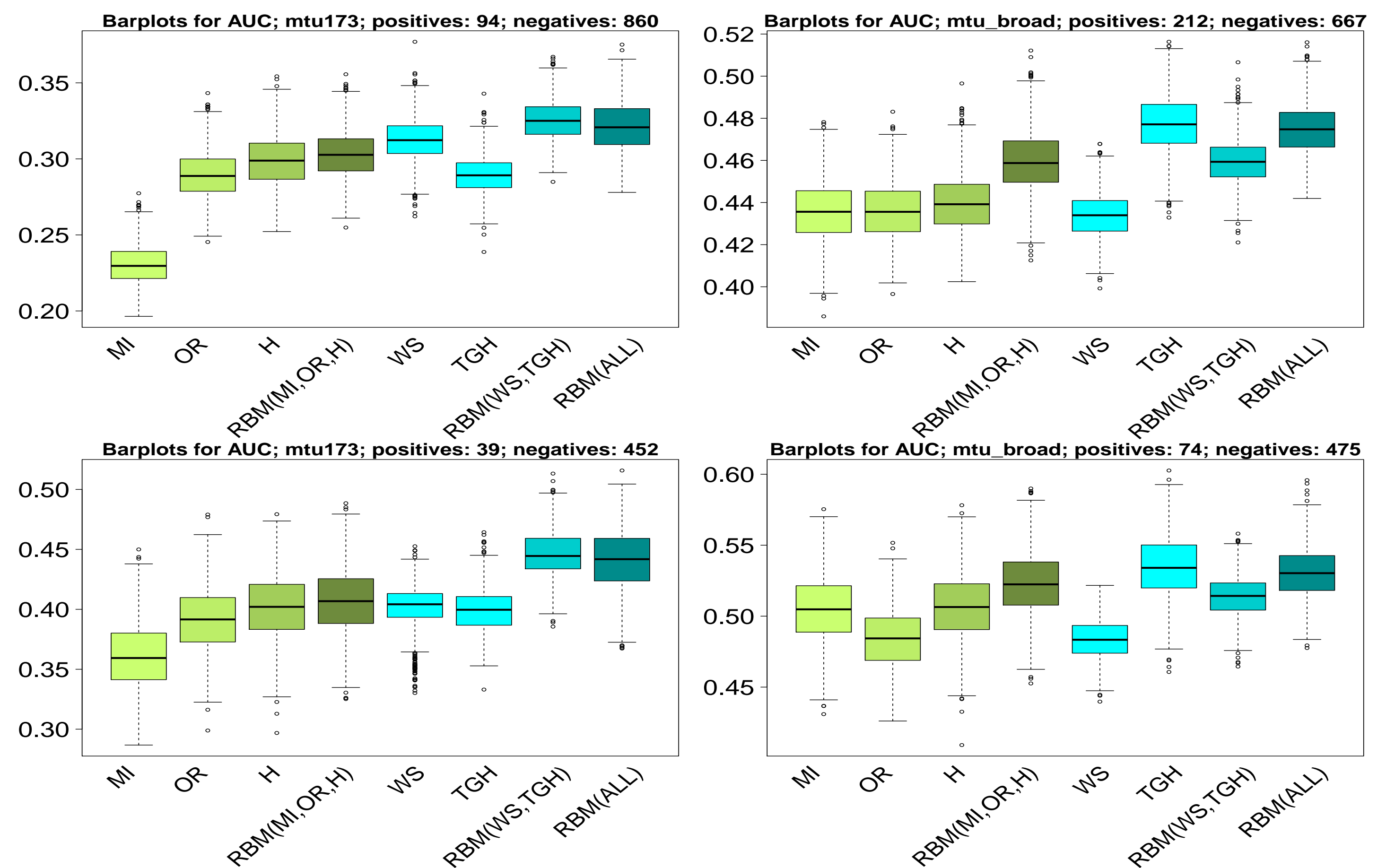


Figure 2: Tree-ignorant vs. tree aware scores: comparison of different association scores implemented in GWAMAR based on the Area Under the Curve (AUC) statistic for the precision-recall curves. Left panels present the results for the *mtu173* dataset; right for the *mtu_broad* dataset. The first row of panels corresponds to the experiments in which all associations present in TBDRaMDB were used as the gold standard, whereas the second row corresponds to the experiments in which only high-confidence associations were used as the gold standard. The process of sampling the set of negatives was repeated 1000 times. The barplots for tree-ignorant and tree-aware scores are shown green and blue, respectively.

Top-scoring mutations

drug name	gene id	gene name	mutation	all	h.c.	TGH	drug name	gene id	gene name	mutation	all	h.c.	TGH
Fluoroquinolones	Rv0006	gyrA	D94H ₁ A ₃ N ₂ Y ₂ G ₁₂	Y	Y	14.184	Fluoroquinolones	Rv0006	gyrA	D94Y ₂ H ₂ A ₂₀ G ₂₀ N ₁₄	Y	Y	128.323
Isoniazid	Rv1908c	katG	S315N ₁ G ₂ T ₇₅	Y	Y	9.045	Rifampicin	Rv0667	rpoB	S450T ₁₄₃ W ₂₂	Y	Y	72.284
Rifampicin	Rv0667	rpoB	S450L ₂₁	Y	Y	8.602	Ethambutol	Rv3795	embB	M300T ₁₁₆ V ₂₀₀ L ₃₁₃	Y	Y	70.217
Streptomycin	Rv0682	rpsL	K43R ₁₅	Y	Y	8.323	Fluoroquinolones	Rv0006	gyrA	A80G ₂ V ₁₆	Y	Y	41.699
Ethambutol	Rv3795	embB	M300H ₁₂ V ₁₈	Y	Y	8.259	Streptomycin	Rv0682	rpsL	K43R ₂₂₈	Y	Y	30.012
Isoniazid	Rv1483	fabG1	C-15T ₃₀	Y	Y	5.845	Isoniazid	Rv1908c	katG	S315T ₃₀₅ G ₂ L ₃ R ₄ N ₂₇	Y	Y	27.966
Rifampicin	Rv0667	rpoB	D435Y ₂ F ₁₁ V ₁₁ G ₃ A ₁	Y	Y	5.040	Ethambutol	Rv3795	embB	Q497H ₂ K ₂₈ F ₁₀ R ₄₃	Y	Y	17.081
Streptomycin	Rv0682	rpsL	K88R ₃ M ₁	Y	Y	4.164	Streptomycin	Rv0682	rpsL	K88Q ₁ R ₂₈ T ₁₀ M ₇	Y	Y	16.327
Ethambutol	Rv3795	embB	E504G ₁ D ₁	N	N	3.351	Fluoroquinolones	Rv0005	gyrB	N338K ₁ S ₁ T ₃ D ₂	Y	Y	12.605
Pyrazinamide	Rv2043c	pncA	H51P ₁	Y	Y	2.708	Rifampicin	Rv0667	rpoB	H445F ₂ Q ₂₁ L ₂₇ Y ₃₃ R ₁₂ D ₂₅ N ₇	Y	Y	12.252
Streptomycin	Rv0667	rpoB	W68L ₁	Y	Y	2.708	Streptomycin	Rvnr01	rrs	A140I ₂ C ₂₄	Y	N	9.509
Rifampicin	Rv0667	rpoB	H445D ₂ Y ₂ R ₁	Y	Y	2.530	Streptomycin	Rvnr01	rrs	A514C ₅₀	Y	Y	8.940
Streptomycin	Rvnr01	rrs	G1108C ₂	N	N	1.717	Pyrazinamide	Rv2043c	pncA	T135A ₁ P ₂₂	Y	Y	8.814
Ethambutol	Rv3795	embB	D869G ₁	N	N	1.688	Fluoroquinolones	Rv0006	gyrA	S91P ₂	Y	Y	7.557
Ethambutol	Rv3795	embB	A305T ₁	N	N	1.688	Rifampicin	Rv0667	rpoB	D455H ₁ N ₂ A ₂ Y ₃₇ G ₁ V ₁₀	Y	Y	7.480
Ethambutol	Rv3795	embB	D1024N ₁	Y	N	1.688	Ethambutol	Rv3795	embB	G406C ₂ A ₂₀ D ₂ S ₁₃	Y	Y	7.057
Fluoroquinolones	Rv0005	gyrB	N538T ₁	Y	Y	1.685	Pyrazinamide	Rv2043c	pncA	T-11G ₁ C ₂₄	Y	Y	6.766
Fluoroquinolones	Rv0006	gyrA	S91P ₁	Y	Y	1.685	Fluoroquinolones	Rv0006	gyrA	D89G ₂ N ₄	Y	N	6.253
Fluoroquinolones	Rv0005	gyrB	T530I ₁	N	N	1.685	Pyrazinamide	Rv2043c	pncA	L120T ₂₀ R ₅	Y	N	6.146
Streptomycin	Rvnr01	rrs	A140I ₁₇	Y	N	1.288	Streptomycin	Rvnr01	rrs	C517I ₂₆	Y	Y	5.169
Ethambutol	Rv3795	embB	Y334H ₂	Y	Y	1.054	Pyrazinamide	Rv2043c	pncA	Q10H ₁ R ₁₀ P ₁₂	Y	Y	5.053
Ethambutol	Rv3795	embB	Q497R ₂	Y	Y	1.054	Pyrazinamide	Rv2043c	pncA	V139M ₂ G ₂ A ₂ L ₁	Y	Y	5.053
Rifampicin	Rv0667	rpoB	E250G ₁	N	N	1.047	Ethambutol	Rv3795	embB	D328G ₂ H ₁ Y ₉	Y	N	5.032
Fluoroquinolones	Rv0006	gyrA	A90V ₂ G ₃	Y	Y	1.035	Streptomycin	Rvnr01	rrs	A908C ₂ G ₁	Y	N	4.779
Streptomycin	Rvnr01	rrs	C517I ₃₃	Y	Y	0.915	Pyrazinamide	Rv2043c	pncA	D12E ₁ G ₂ N ₁ A ₁₂	Y	Y	4.725

Figure 3: 25 top-scoring associations between drug resistance profiles and point mutations in the case study on 173 fully sequenced *M. tuberculosis* strains (left table) and 1398 *M. tuberculosis* strains for the Broad Institute dataset (right table). The associations are restricted to only these genes which are associated with drug resistance to the corresponding drugs. Each row corresponds to one association, whereas the consecutive columns describe: drug name, gene identifier, gene name, mutation, association presence in the TBDRaMDB database, status indicating whether the association is categorized as high-confidence in TBDRaMDB, and the TGH score. Lower indexes in the mutation descriptions indicate the numbers of strains possessing the corresponding amino acid or nucleotide variant.

Compensatory mutations

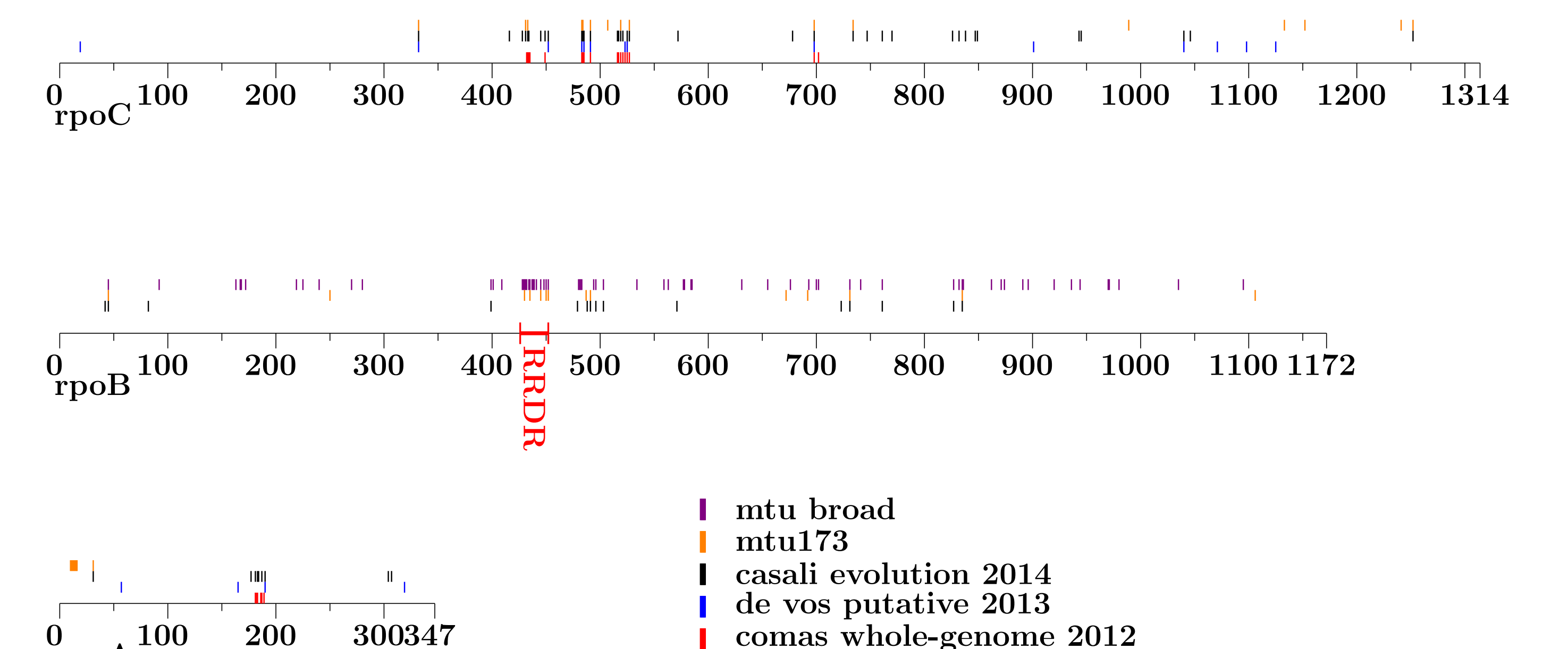


Figure 4: Comparison of the sets of putative compensatory mutations within the *rpoA*, *rpoB* and *rpoC* genes, reported in various sources and detected in our two datasets. Each mutation's position is indicated by a vertical line of the color corresponding to the source it was reported in. In particular orange and violet lines indicate positions of mutations identified by our approach applied to the *mtu173* and *mtu_broad* datasets, respectively. The other lines indicate mutations reported in the recent articles.

References

- Woźniak M., Tiuryn J., Wong L. An approach to identifying drug resistance associated mutations in bacterial strains BMC Genomics; 2012
- Woźniak M., Wong L., Tiuryn J. eCAMBER: efficient support for large-scale comparative analysis of multiple bacterial strains BMC Bioinformatics; 2014
- Woźniak M., Tiuryn J., Wong L. GWAMAR: Genome-wide assessment of mutations associated with drug resistance in bacteria BMC Genomics; 2014
- Woźniak M. Computational aspects of the presence of drug resistance mechanisms PhD thesis; 2015