

GWAMAR: Genome-wide assessment of mutations associated with drug resistance in bacteria

Michal Wozniak^{1,2}, Limsoon Wong² and Jerzy Tiuryn¹

¹University of Warsaw

²National University of Singapore

16 December, 2014



Introduction

- Mechanisms of drug action against bacteria
- Mechanisms of drug resistance in bacteria

Methods

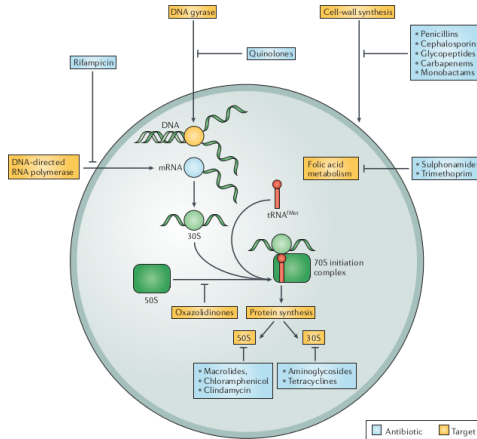
- Schema of the approach
- Input data
- Association scores

Results

- Input datasets
- Comparison of different association scores
- Top-scoring mutations
- Compensatory mutations

Summary

Drug action mechanisms



Adopted from: Platforms for antibiotic discovery; Kim Lewis; Nature Reviews; 2013

Timeline of antibiotics

Antibiotic class; example	Discovery	Introduction	Resistance	Mechanism of action	Activity or target species
Sulfadruugs; protonosil	1932	1936	1942	Inhibition of dihydro- pteroate synthetase	Gram-positive bacteria
β -lactams; penicillin	1928	1938	1945	Inhibition of cell wall biosynthesis	Broad-spectrum activity
Aminoglycosides; streptomycin	1943	1946	1946	Binding of 30S ribosomal subunit	Broad-spectrum activity
Chloramphenicols; chloramphenicol	1946	1948	1950	Binding of 50S ribosomal subunit	Broad-spectrum activity
Macrolides; erythromycin	1948	1951	1955	Binding of 50S ribosomal subunit	Broad-spectrum activity
Tetracyclines; chlortetracycline	1944	1952	1950	Binding of 30S ribosomal subunit	Broad-spectrum activity
Rifamycins; rifampicin	1957	1958	1962	Binding of RNA polymerase β -subunit	Gram-positive bacteria
Glycopeptides; vancomycin	1953	1958	1960	Inhibition of cell wall biosynthesis	Gram-positive bacteria
Quinolones; ciprofloxacin	1961	1968	1968	Inhibition of DNA synthesis	Broad-spectrum activity
Streptogramins; streptogramin B	1963	1998	1964	Binding of 50S ribosomal subunit	Gram-positive bacteria
Oxazolidinones; linezolid	1955	2000	2001	Binding of 50S ribosomal subunit	Gram-positive bacteria
Lipopetides; daptomycin	1986	2003	1987	Depolarization of cell membrane	Gram-positive bacteria
Fidaxomicin	1948	2011	1977	Inhibition of RNA polymerase	Gram-positive bacteria
Diarylquinolines; bedaquiline	1997	2012	2006	Inhibition of F_1F_0 -ATPase	Narrow-spectrum activity

Timeline of the discovery and introduction of antibiotics (based on Platforms for antibiotic discovery; Kim Lewis; Nature Reviews; 2013).

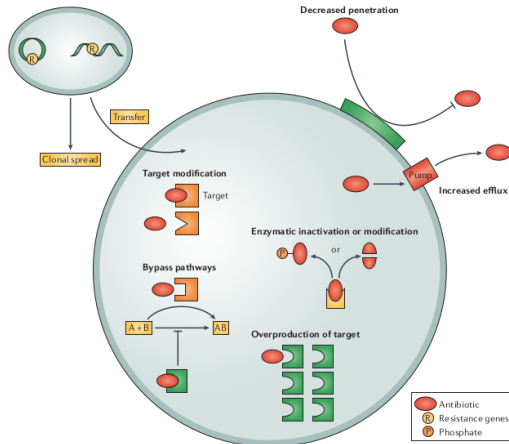
Drug resistance mechanisms I

There are several known drug resistance mechanisms which can be categorized as follows (adopted from: *Wright GD, Chem. Comm., 2011*):

- ▶ drug target modification;
- ▶ drug molecule modification by specialized enzymes
- ▶ reduced accumulation of the drug inside a bacteria cell by decreased cell wall permeability or by pumping out the drug
- ▶ alternative metabolic pathways

These drug resistance mechanisms can be acquired either by **chromosomal mutations** or **horizontal gene transfer**.

Drug resistance mechanisms II



Adopted from: Platforms for antibiotic discovery; Kim Lewis; Nature Reviews; 2013

GWAMAR: drug resistance-associated mutations

Goal: identify drug resistance-associated mutations (primary and secondary)

General approach implemented in GWAMAR:

- ▶ we use whole-genome comparative approach to identify genetic variations among multiple bacterial strains,
- ▶ we retrieve from literature and databases information of the drug resistance phenotypes of the strains,
- ▶ we associate the identified mutations with drug resistance-phenotypes based on association scores,
- ▶ we propose a new association score, called TGH, which implements scores phylogenetic information.

Genotype and phenotype data

Genotype data

We consider two kinds of genetic variations (determined by eCAMBer based on gene families and their multiple alignments):

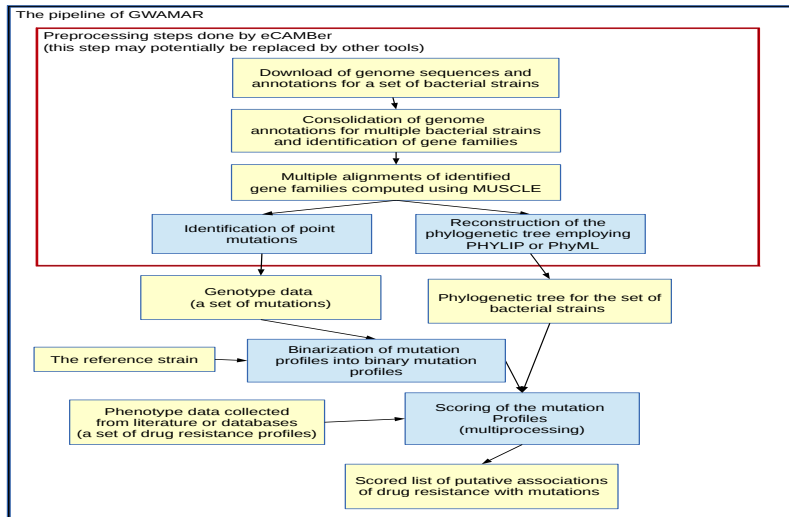
- ▶ gene gain/loss,
- ▶ amino acid point mutation.

These genetic variations are represented as '0'-'1' vectors (called **mutation profiles**), where '0' denotes the reference state and '1' denotes some change.

Phenotype data (drug susceptibility)

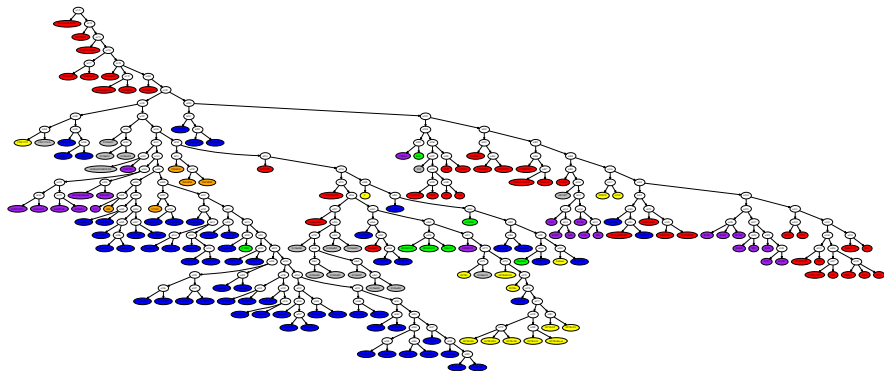
Phenotype data are represented as vectors, called **drug resistance profiles**, with possible states: 'S', 'R', 'I', '?'.

Schema of the framework



Tree-aware scores

We observe that subtrees of the phylogenetic tree very often correspond to geographic locations. Since drug resistance mutations are subject to evolutionary pressure caused by the drug treatment they should be independent of geographic location and therefore be more widely distributed over the tree, as opposed to mutations driven by other environmental factors which tend to rather concentrate in small subtrees.



Classical scores (tree-ignorant) association scores

The classical scores used in genotype-phenotype association studies and co-evolution studies are tree-ignorant.

- ▶ odds ratio:

$$\text{OR}(b, r) = \frac{n_1^R \cdot n_0^S}{\max(1, n_0^R) \cdot \max(1, n_1^S)}$$

- ▶ mutual information:

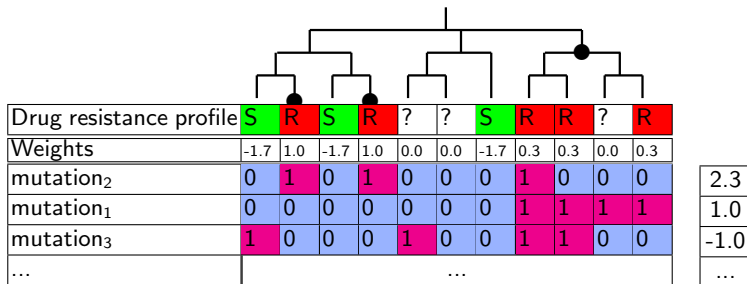
$$\text{MI}(b, r) = \sum_{x \in \{0, 1\}} \sum_{y \in \{S, I, R\}} \frac{n_x^y}{n} \cdot \log\left(\frac{n_x^y \cdot n}{n_x \cdot n^y}\right)$$

- ▶ hypergeometric score

$$H(b, r) = -\log\left(\sum_{i=n^R}^n H(n, n^R, n_1, i)\right)$$

Weighted support

Weighted support rewards for drug-resistant strains with the mutation, penalty for drug-susceptible strains with the mutation, where weight $w_T(b, r, i)$ for drug resistant strains is $\frac{1}{k}$, where k denotes the size of the largest subtree with only drug resistant strains.



Weighted support for mutation m is defined as follows:

$$WS_T(b, r) = \sum_{i \in S} w_T(b, r, i) [b(i) = '1']$$

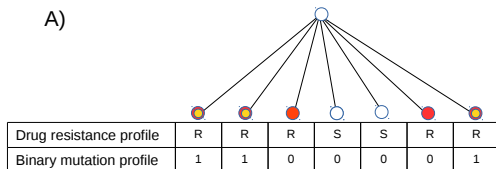
TGH score I

For a given tree T , we call a subset c of its nodes a *coloring*, if it satisfies the following two conditions:

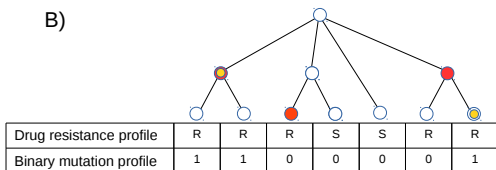
- ▶ each path from a leaf to the root contains at most one node from c ,
- ▶ each internal node in c has a sibling node which does not belong to c .

TGH score II

A)



B)



(A) an example of coloring \hat{c} induced by a given drug resistance profile (large red nodes) and coloring \bar{c} induced by a given binary mutation profile (small orange nodes) for a flat tree. In this example $|\hat{c}| = 5$, $|\bar{c}| = 3$ and $|L(\hat{c}) \cap \bar{c}| = 3$. (B) another example of colorings \hat{c} and \bar{c} induced by the same pair of profiles but for a different tree. In this example $|\hat{c}| = 3$, $|\bar{c}| = 2$ and $|L(\hat{c}) \cap \bar{c}| = 2$.

TGH score III

We define the TGH score as follows:

$$TGH_T(r, b) = -\log\left(\frac{\sum_{i=k}^n B_{T, \hat{c}}(i, n)}{V_T(n)}\right)$$

where:

$$V_T(n) = \#\{c \in C_T : |c| = n\}$$

and:

$$B_{T, \hat{c}}(k, n) = \#\{c \in C_T : |L(\hat{c}) \cap c| = k \text{ and } |c| = n\}$$

GWAMAR implements a dynamic programming approach to calculate the score. The time complexity is $O(D \cdot N^{K-1} \cdot N^2 + D \cdot N \cdot M)$.

Input datasets

We have two datasets of data for *M. tuberculosis*

- ▶ 1398 strains with 28 genes sequenced from Broad Institute (mtu_broad)
- ▶ 173 fully sequenced strains available in NCBI and PATRIC databases (mtu173)

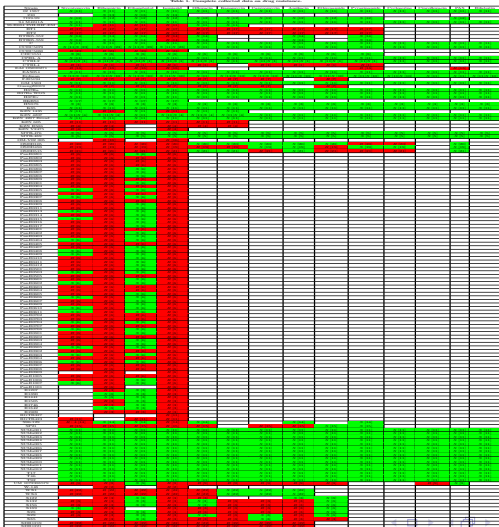
Genotype data

Point mutation profiles were determined based on gene families identified with *eCAMBer* and their multiple alignments computed with MUSCLE.

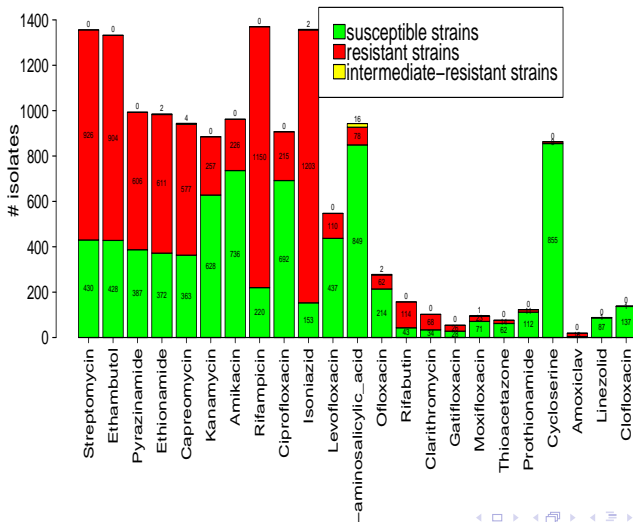
Phenotype data (drug susceptibility)

- ▶ publications issued together with the fully sequenced genomes;
- ▶ other publications found by searching of related literature;
- ▶ drug resistance profiles for separate drugs are combined into: Rifampicin, Isoniazid, Fluoroquinolones, Ethambutol, Pyrazinamide, Streptomycin

Collected phenotype data for the mtu173 dataset



Phenotype data for the mtu_broad dataset

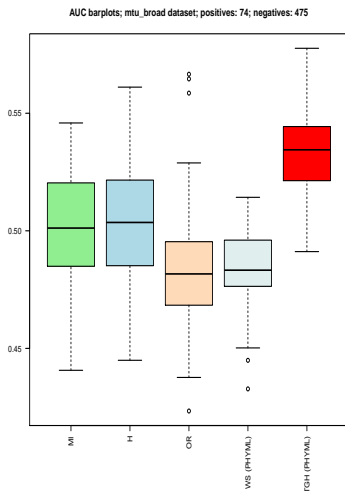
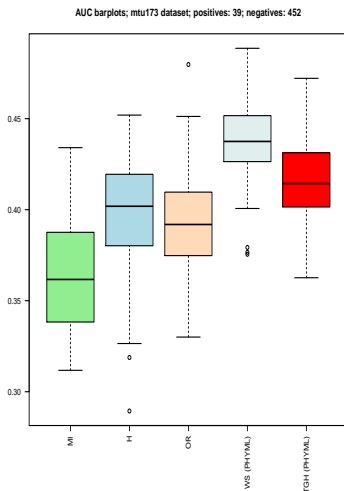


Gold standard associations

We retrieve the gold standard associations from the TBDreamDB database for:
Rifampicin, Isoniazid, Fluoroquinolones, Ethambutol, Pyrazinamide, Streptomycin.

drug name	gene name	positions
Fluoroquinolones	gyrA	90,91,94,102,126
	gyrB	538
Ethambutol	embB	306,406,497
Isoniazid	ahpC	-46,-39,21
	fabG1-inhA	-15,-8
	kasA	269
	katG	315
Isoniazid	ahpC	-46,-39,21
Rifampicin	rpoB	432,435,441,445,450,452
Streptomycin	rpsL	43,88
	rrs	492,513,514,517,907
Pyrazinamide	pncA	-11,7,10(60 in total)

Comparison on *mtu173* and *mtu_broad* datasets



Top-scoring mutations on the *mtu173* dataset

drug name	gene id	gene name	mutation	all	h.c.	TGH
Fluoroquinolones	Rv0006	gyrA	D94G/A/H/N/Y	Y	Y	14.1474612056
Rifampicin	Rv0667	rpoB	S450L	Y	Y	13.7718163314
Isoniazid	Rv1908c	katG	S315T/G/N	Y	Y	11.3299511445
Streptomycin	Rv0682	rpsL	K43R	Y	Y	8.48207709024
Rifampicin	Rv0667	rpoB	D435F/V/Y/G/A	Y	Y	6.27518366477
Ethambutol	Rv3795	embB	M306I/V/L	Y	Y	5.86457315376
Isoniazid	Rv1483	fabG1	C-15T	Y	Y	5.54308155286
Streptomycin	Rv0682	rpsL	K88R/M	Y	Y	4.2132021894
Ethambutol	Rv3795	embB	E504G/D	N	N	3.33044011637
Rifampicin	Rv0667	rpoB	H445D/Y/R	Y	Y	2.71070582544
Pyrazinamide	Rv2043c	pncA	H51P	Y	Y	2.7080502011
Pyrazinamide	Rv2043c	pncA	W68L	Y	Y	2.7080502011
Fluoroquinolones	Rv0006	gyrA	A90V/G	Y	Y	2.50721373928
Fluoroquinolones	Rv0005	gyrB	N538T	Y	Y	1.82454929205
Fluoroquinolones	Rv0006	gyrA	S91P	Y	Y	1.82454929205
Fluoroquinolones	Rv0005	gyrB	T539I	N	N	1.82454929205
Ethambutol	Rv3795	embB	D869G	N	N	1.67147330335
Ethambutol	Rv3795	embB	A505T	N	N	1.67147330335
Ethambutol	Rv3795	embB	D1024N	Y	N	1.67147330335
Rifampicin	Rv0667	rpoB	A692T	N	N	1.65151195989
Isoniazid	Rv1483	fabG1	T-8A/C	Y	Y	1.36702851766
Streptomycin	Rvnr01	rrs	A1401G	Y	N	1.31458509629
Rifampicin	Rv0667	rpoB	I1106T	N	N	1.11729579263
Rifampicin	Rv0667	rpoB	E250G	N	N	1.11729579263
Rifampicin	Rv0667	rpoB	L452P	Y	Y	1.11729579263

Top-scoring mutations on the *mtu_broad* dataset

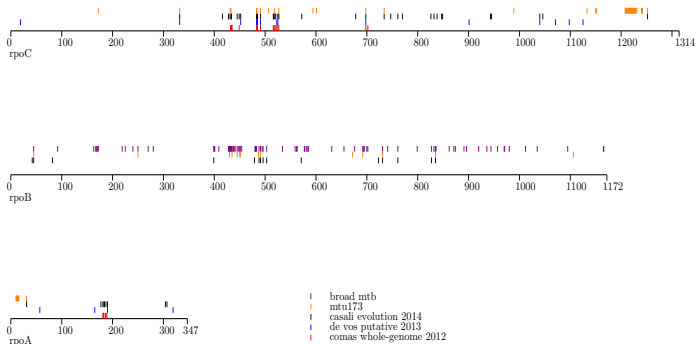
drug name	gene id	gene name	mutation	all	h.c.	TGH
Fluoroquinolones	Rv0006	gyrA	D94A/G/N/Y/H	Y	Y	138.414507992
Ethambutol	Rv3795	embB	M306V/I/L/T	Y	Y	75.677195766
Rifampicin	Rv0667	rpoB	S450L/W	Y	Y	73.8586969047
Fluoroquinolones	Rv0006	gyrA	A90G/V	Y	Y	41.3671629169
Isoniazid	Rv1908c	katG	S315T/G/N/I/R	Y	Y	31.2319396693
Streptomycin	Rv0682	rpsL	K88T/R/Q/M	Y	Y	22.7910550626
Ethambutol	Rv3795	embB	Q497P/R/K/H	Y	Y	18.8920918506
Streptomycin	Rv0682	rpsL	K43R	Y	Y	15.0916682608
Fluoroquinolones	Rv0005	gyrB	N538K/T/D/S	Y	Y	14.6924139234
Rifampicin	Rv0667	rpoB	H445P/D/R/Y/L/N/Q	Y	Y	13.0077710882
Pyrazinamide	Rv2043c	pncA	T135P/A	Y	N	9.80925396944
Pyrazinamide	Rv2043c	pncA	T-11C/G	Y	Y	9.01421480751
Ethambutol	Rv3795	embB	G406S/D/A/C	Y	Y	8.05127512926
Fluoroquinolones	Rv0006	gyrA	S91P	Y	Y	7.62877200647
Streptomycin	Rvnr01	rrs	A1401G	Y	N	7.62190832659
Rifampicin	Rv0667	rpoB	D435Y/V/H/G/A/N	Y	Y	7.54687538111
Streptomycin	Rvnr01	rrs	A514C	Y	Y	7.04416442854
Pyrazinamide	Rv2043c	pncA	Q10R/P/H	Y	Y	6.79992327959
Fluoroquinolones	Rv0006	gyrA	D89G/N	Y	N	6.3023070068
Ethambutol	Rv3795	embB	D328Y/G/H	Y	N	5.93350968238
Pyrazinamide	Rv2043c	pncA	V139G/A/M/L	Y	Y	5.67144832937
Isoniazid	Rv1483	fabG1	C-15T	Y	Y	5.5292894594
Pyrazinamide	Rv2043c	pncA	T76P/I	Y	Y	5.22287436396
Pyrazinamide	Rv2043c	pncA	D12G/A/E/N	Y	Y	5.10653842681
Ethambutol	Rv3795	embB	D354A/G	N	N	5.02061345733

Putative compensatory mutations

Recent publications reporting putative compensatory mutations in *M. tuberculosis*:

- ▶ Whole-genome sequencing of rifampicin-resistant *Mycobacterium tuberculosis* strains identifies compensatory mutations in RNA polymerase genes; *Nature Genetics*; 2012
- ▶ Putative Compensatory Mutations in the *rpoC* Gene of Rifampin-Resistant *Mycobacterium tuberculosis* Are Associated with Ongoing Transmission; *Antimicrobial Agents and Chemotherapy*; 2013
- ▶ Evolution and transmission of drug-resistant tuberculosis in a Russian population; *Nature Genetics*; 2014

Putative compensatory mutations



Interestingly, several mutations identified by GWAMAR that has also been reported in at least one of the papers.

- ▶ *rpoA*: G31S/A
- ▶ *rpoB*: P45S/L, L731P, E761D, R827C, H835P/R
- ▶ *rpoC*: G332R/S, V431M, G433C/S, V483G/A, W484G, I491T/V, L527V, N698K, A734V

Summary

- ▶ The fast growing number of fully sequenced bacterial strains enables us to develop and test new methods to identifying drug resistance associated genes and mutations.
- ▶ We developed and implemented GWAMAR – a new framework for detection of drug resistance-associated mutations. This software is available at the project website: <http://bioputer.mimuw.edu.pl/gwamar/>.
- ▶ We proposed a new association score, called TGH, which employ phylogenetic information. It outperforms the standard tree-ignorant scores, but is more computationally expensive.
- ▶ Applying our approach we identified some novel putative drug resistance-associated mutations.
- ▶ Future possible direction of research may include: classification of mutations into primary and secondary, grouping of mutations which are close together, incorporation of PPI networks.

Thank you

Thank you!

You are welcome to give comments or ask questions.