

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

27 March, 2015



Introduction

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

Methods

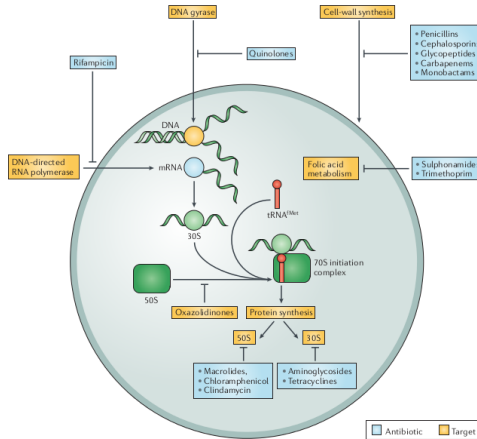
- Schema of the approach
- Input data
- Association scores

Results

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

Summary

Drug action mechanisms



Rysunek : 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

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

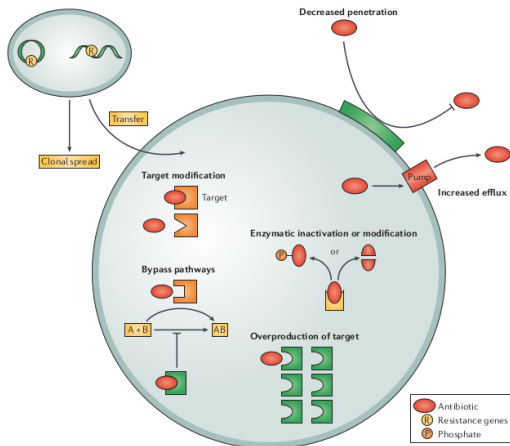
Drug resistance mechanisms I

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



Rysunek : 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):

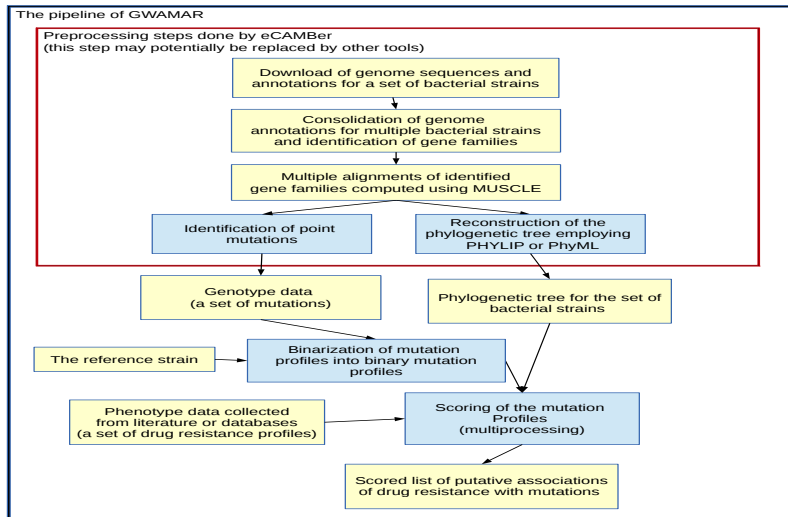
- ▶ gene gain/loss,
- ▶ amino acid point mutation.

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)

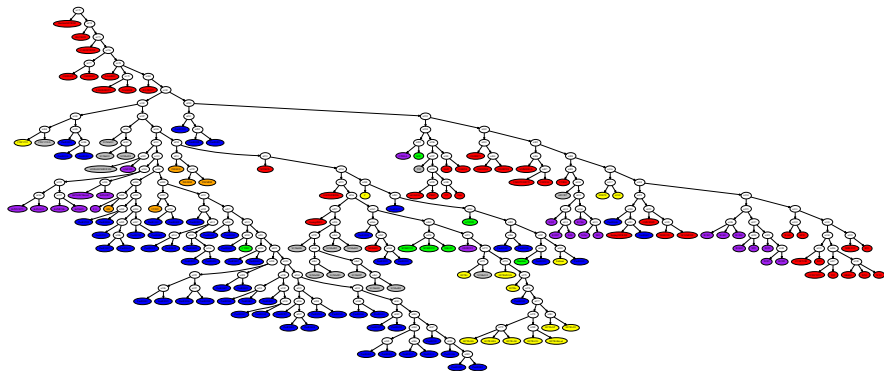
Phenotype data are represented as vectors, called **drug resistance profiles**, with possible states: 'S', 'R', 'I', '?'.
'S' stands for Susceptible, 'R' for Resistant, 'I' for Intermediate, and '?' for Unknown.

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:

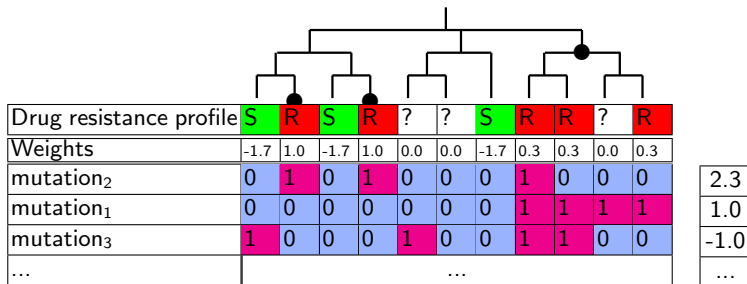
$$\text{MI}(b, r) = \sum_{x \in \{0, 1\}} \sum_{y \in \{S, I, R\}} \frac{n_x^y}{n} \cdot \log\left(\frac{n_x^y \cdot n}{n_x \cdot n^y}\right)$$

- ▶ 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:

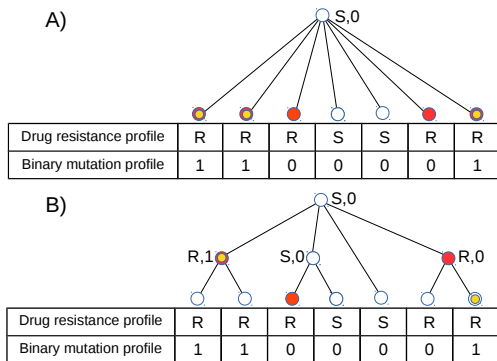
$$WS_T(b, r) = \sum_{i \in S} w_T(b, r, i) [b(i) = '1']$$

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



Rysunek : (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)

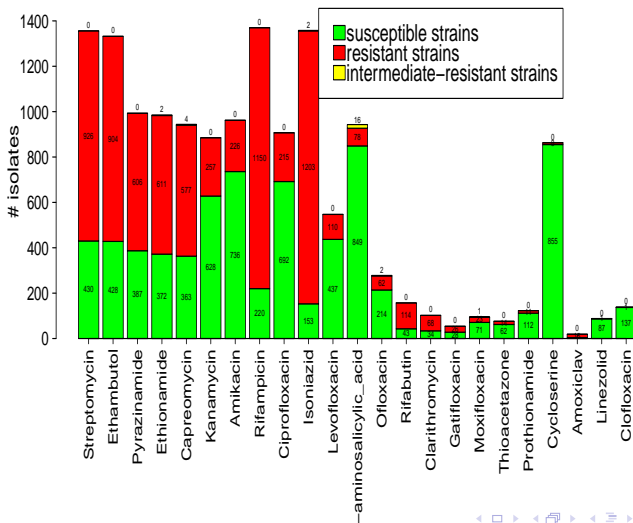
Genotype data

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

Phenotype data for the mtu_broad dataset

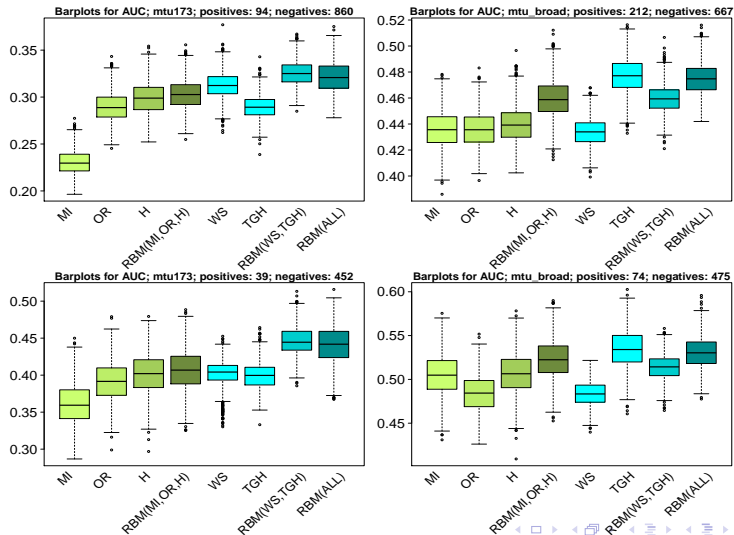


Gold standard associations

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

drug name	gene name	positions
Fluoroquinolones	gyrA gyrB	90,91,94,102,126 538
Ethambutol	embB	306,406,497
Isoniazid	ahpC fabG1-inhA kasA katG	-46,-39,21 -15,-8 269 315
Rifampicin	rpoB	432,435,441,445,450,452
Streptomycin	rpsL rrs	43,88 492,513,514,517,907
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	D94H ₁ A ₅ N ₂ Y ₂ G ₁₂	Y	Y	14.184
Isoniazid	Rv1908c	katG	S315N ₁ G ₂ T ₇₅	Y	Y	9.045
Rifampicin	Rv0667	rpoB	S450L ₇₁	Y	Y	8.602
Streptomycin	Rv0682	rpsL	K43R ₁₅	Y	Y	8.323
Ethambutol	Rv3795	embB	M306L ₁ I ₃₂ V ₁₈	Y	Y	8.250
Isoniazid	Rv1483	fabG1	C-15T ₃₀	Y	Y	5.845
Rifampicin	Rv0667	rpoB	D435Y ₂ F ₅ V ₁₁ G ₃ A ₁	Y	Y	5.040
Streptomycin	Rv0682	rpsL	K88R ₅ M ₁	Y	Y	4.164
Ethambutol	Rv3795	embB	E504G ₁ D ₁	N	N	3.331
Pyrazinamide	Rv2043c	pncA	H51P ₁	Y	Y	2.708
Pyrazinamide	Rv2043c	pncA	W68L ₁	Y	Y	2.708
Rifampicin	Rv0667	rpoB	H445D ₈ Y ₂ R ₁	Y	Y	2.530
Streptomycin	Rvnr01	rrs	G1108C ₂	N	N	1.717
Ethambutol	Rv3795	embB	D869G ₁	N	N	1.688
Ethambutol	Rv3795	embB	A505T ₁	N	N	1.688
Ethambutol	Rv3795	embB	D1024N ₁	Y	N	1.688
Fluoroquinolones	Rv0005	gyrB	N538T ₁	Y	Y	1.685
Fluoroquinolones	Rv0006	gyrA	S91P ₁	Y	Y	1.685
Fluoroquinolones	Rv0005	gyrB	T539I ₁	N	N	1.685
Streptomycin	Rvnr01	rrs	A1401G ₁₇	Y	N	1.288
Ethambutol	Rv3795	embB	Y334H ₂	Y	N	1.054
Ethambutol	Rv3795	embB	Q497R ₂	Y	Y	1.054
Rifampicin	Rv0667	rpoB	E250G ₃	N	N	1.047
Fluoroquinolones	Rv0006	gyrA	A90V ₆ G ₃	Y	Y	1.035
Streptomycin	Rvnr01	rrs	C517T ₃₃	Y	Y	0.915

Top-scoring mutations on the *mtu_broad* dataset

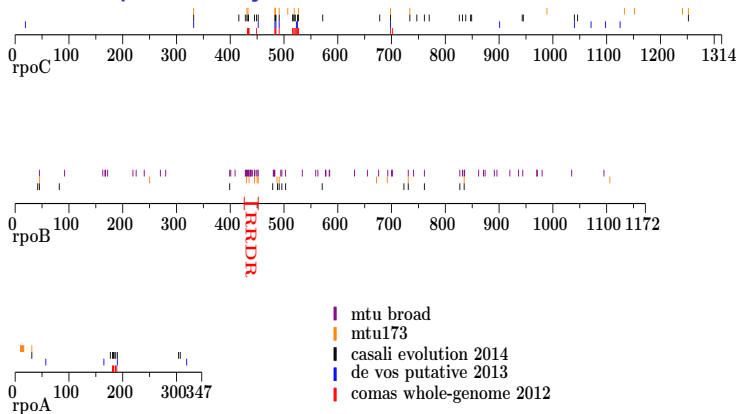
drug name	gene id	gene name	mutation	all	h.c.	TGH
Fluoroquinolones	Rv0006	gyrA	D94Y ₆ H ₂ A ₂₆ G ₇₈ N ₁₄	Y	Y	128.323
Rifampicin	Rv0667	rpoB	S450L ₇₄₃ W ₂₂	Y	Y	72.284
Ethambutol	Rv3795	embB	M306T ₁ L ₁₆ V ₂₉₀ I ₃₁₃	Y	Y	70.217
Fluoroquinolones	Rv0006	gyrA	A90G ₂ V ₄₆	Y	Y	41.699
Streptomycin	Rv0682	rpsL	K43R ₂₂₈	Y	Y	30.012
Isoniazid	Rv1908c	katG	S315T ₈₉₅ G ₂ I ₃ R ₃ N ₂₇	Y	Y	27.966
Ethambutol	Rv3795	embB	Q497H ₅ K ₁₈ P ₁₀ R ₄₃	Y	Y	17.081
Streptomycin	Rv0682	rpsL	K88Q ₁ R ₂₈ T ₃₂ M ₇	Y	Y	16.327
Fluoroquinolones	Rv0005	gyrB	N538K ₁ S ₁ T ₉ D ₂	Y	Y	12.605
Rifampicin	Rv0667	rpoB	H445P ₂ Q ₂ L ₂₇ Y ₅₃ R ₄₂ D ₂₅ N ₇	Y	Y	12.252
Streptomycin	Rvnr01	rrs	A140I _{G254}	Y	N	9.509
Streptomycin	Rvnr01	rrs	A514C ₉₀	Y	Y	8.940
Pyrazinamide	Rv2043c	pncA	T135A ₁ P ₂₂	Y	N	8.814
Fluoroquinolones	Rv0006	gyrA	S91P ₉	Y	Y	7.557
Rifampicin	Rv0667	rpoB	D435H ₁ N ₂ A ₂ Y ₂₇ G ₃ V ₁₄₀	Y	Y	7.480
Ethambutol	Rv3795	embB	G406C ₃ A ₆₈ D ₅₂ S ₄₃	Y	Y	7.057
Pyrazinamide	Rv2043c	pncA	T-11G ₃ C ₂₄	Y	Y	6.766
Fluoroquinolones	Rv0006	gyrA	D89G ₂ N ₄	Y	N	6.253
Pyrazinamide	Rv2043c	pncA	L120P ₂₀ R ₅	Y	N	6.146
Streptomycin	Rvnr01	rrs	C517T ₂₆	Y	Y	5.169
Pyrazinamide	Rv2043c	pncA	Q10H ₃ R ₁₀ P ₁₂	Y	Y	5.053
Pyrazinamide	Rv2043c	pncA	V139M ₃ G ₂ A ₇ L ₁	Y	Y	5.053
Ethambutol	Rv3795	embB	D328G ₅ H ₁ Y ₉	Y	N	5.032
Streptomycin	Rvnr01	rrs	A908C ₇ G ₁	Y	N	4.779
Pyrazinamide	Rv2043c	pncA	D12E ₁ G ₅ N ₁ A ₁₂	Y	Y	4.725

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/figures/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.