

Bioinformatic Methods for Analyzing Mutations in Bacterial Genomes

Howard Fan

May 14, 2025



RUTGERS

What is Mutation Analysis

- Study of genetic alterations in genomes
- Focuses on point mutations, insertions, deletions, and structural changes
- Key to understanding evolution, antibiotic resistance, and pathogenicity

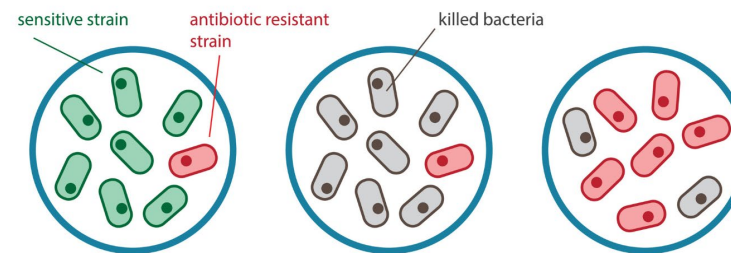


RUTGERS

Why do we care about mutations in bacteria?

Mutations can cause bacteria to be more resistant to antibiotics, leading to:

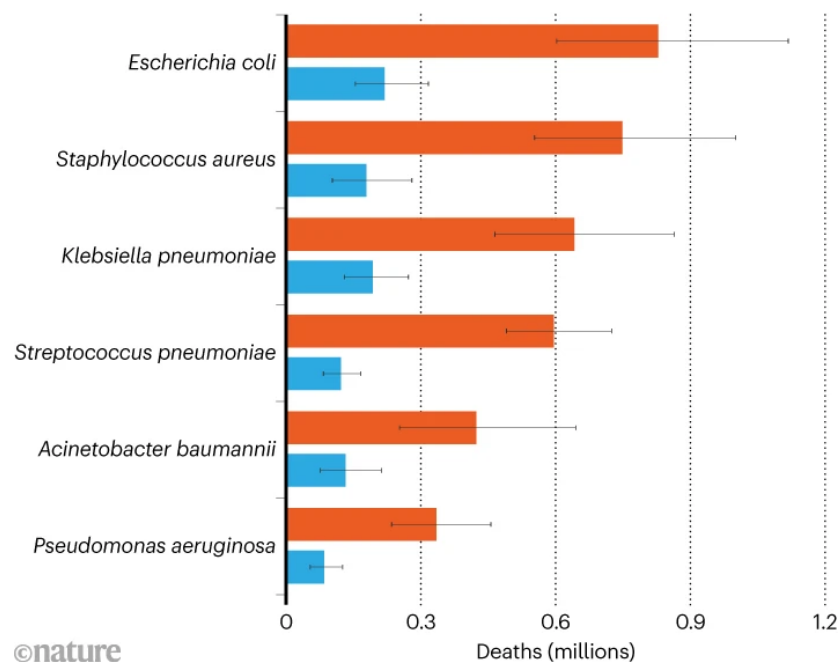
- **Reduced effectiveness of treatments**
 - Can lead to longer illnesses, more severe disease, and higher risk of death
- **Increased medical costs and burden**
 - Resistant infections often require more expensive or risky alternatives, longer hospital stays, and additional medical care.



DEADLY INFECTIONS

These 6 pathogens were responsible for almost 80% of the 1.27 million deaths attributed directly to antimicrobial resistance in 2019.

■ Associated with resistance ■ Attributable to resistance

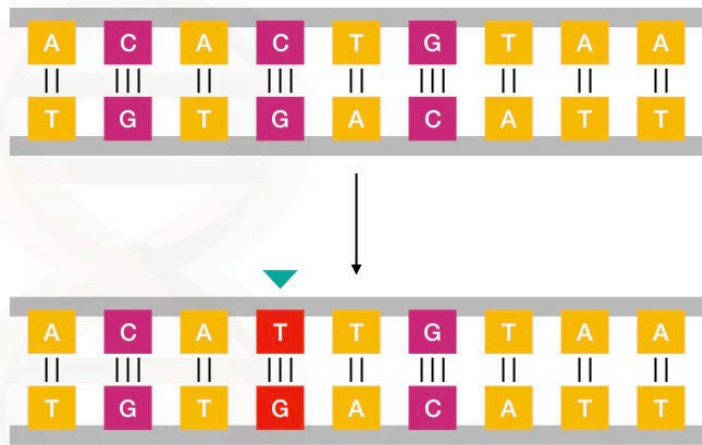




RUTGERS

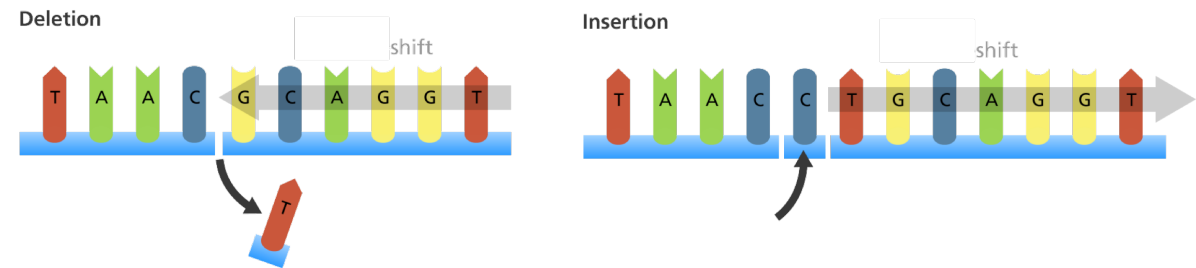
Types of Mutations

Single Nucleotide Polymorphisms



- Most common
- Also known as point mutations
- May be synonymous (silent) or non-synonymous (nonsense or missense)

Indels



- Insertion or deletion
- Often cause frameshifts if in coding regions

Identifying Variants from Bacterial Genomes

Bacterial Genome Workflow

1. Sample collection and DNA extraction
2. Sequencing (short-read vs long-read)
3. Quality control and preprocessing
4. Read alignment
5. Variant calling
6. Annotation
7. Downstream analysis



Variant Calling

We can detect variants from aligned sequences by comparing them to the reference sequence.

- SNP and Indel Callers
 - bcftools mpileup, GATK HaplotypeCaller
- Input requirements:
 - Aligned BAM files, reference genome
- Output:
 - VCF files listing variants, alleles, positions

Sequence Variants

SNV (**S**ingle **N**ucleotide **V**ariant)

Ref	A	A	G	G	G	C	T	G
Query	A	A	G	G	A	C	T	G

—

INDEL (**I**nsertion or **D**eletion)

Ref	A	A	G	G	G	C	T	G
Query	A	A	G	-----		C	T	G



RUTGERS

Variant Filtering

- Why filter?
 - Raw variant calls can include false positives
- Filtering strategies:
 - Depth of coverage (e.g., $DP > 10$)
 - Base quality (e.g., $QUAL > 30$)
 - Strand bias, allele frequency
- Tools:
 - bcftools filter, GATK VariantFiltration
 - Visual inspection with Integrative Genomics Viewer (IGV) if needed



Variant Annotation

Variant annotation tools predict the functional impact of variants and add annotations such as gene names.

- Functional impacts:
 - Synonymous, nonsynonymous, frameshift
 - Premature stop codons, splice site mutations
- Tools: SnpEff

feature type

using Sequence Ontology
transcript, motif, miRNA ...

feature ID

dependent on annotation
transcript ID, motif ID, ChipSwq peak ...

gene name

common gene name (HGNC)

biotype

Ensemble biotypes
Coding, non-coding..

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO
chr1	123456	.	C	A	.	.	ANN=A ...
chr1	234567	.	A	G,T	.	.	ANN=G ..., T ...

ANN = Annotation aka effect or consequence

putative impact

description of consequence
exon_loss_variant, stop_lost,
frameshift_variant

impact

estimation of level of impact
HIGH, LOW, MODERATE

SnpEff

Genetic variant annotation and effect prediction toolbox.



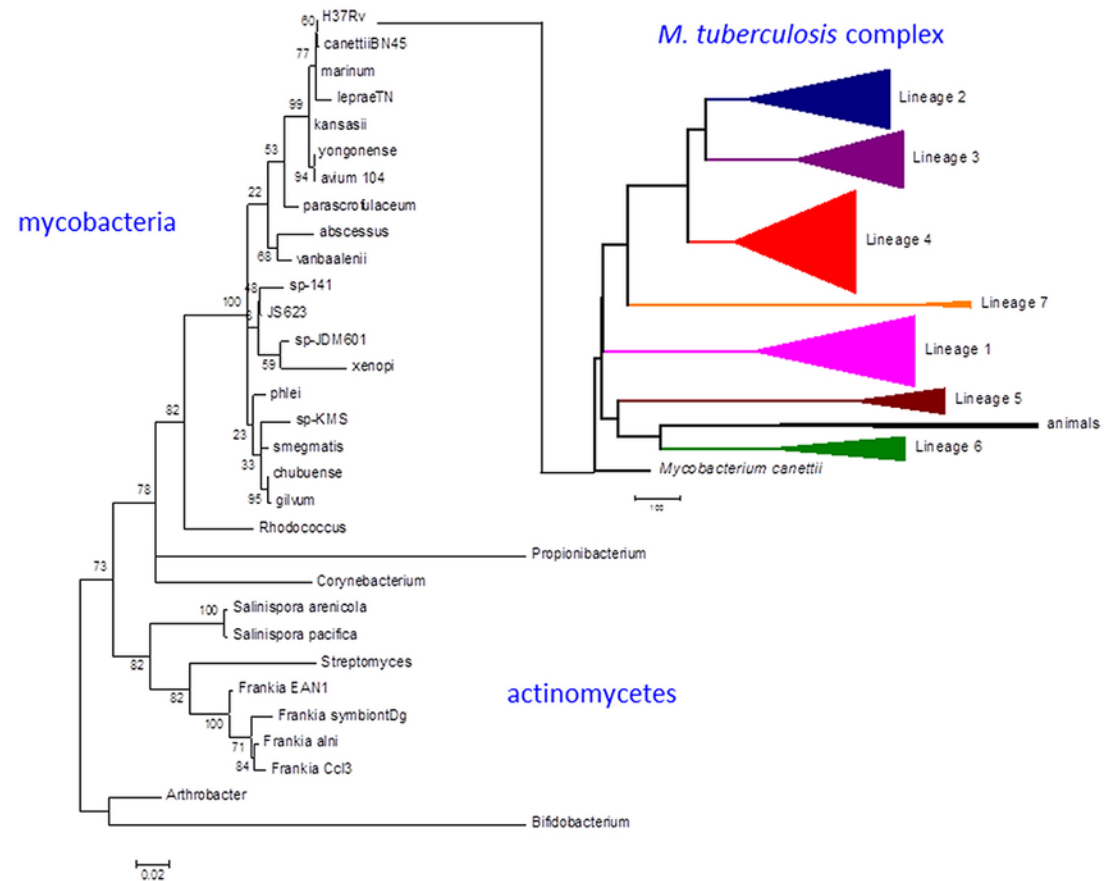
RUTGERS

Bacterial Phylogeny

Bacterial phylogeny provides insights on how bacteria evolves, spreads, and adapts.

- For example, mycobacterium tuberculosis has 7 known lineages that vary in virulence and drug resistance.

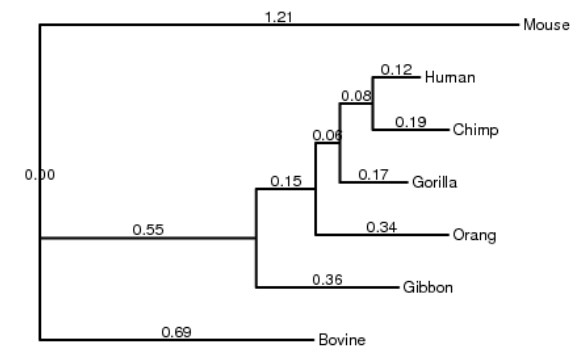
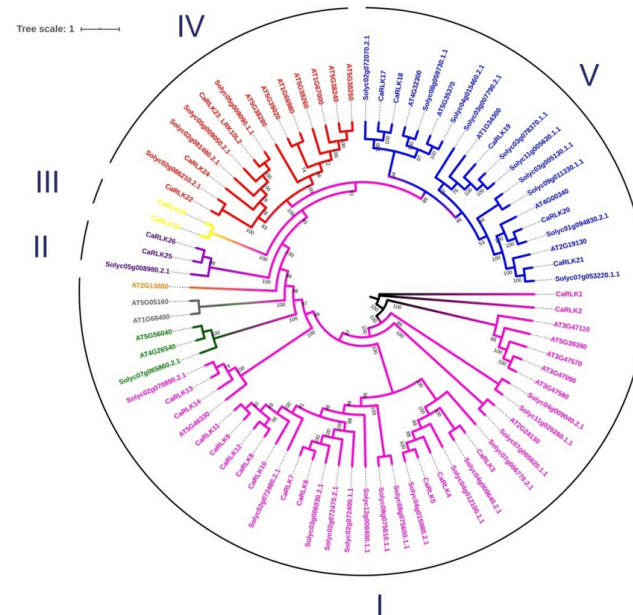
Hence, when studying bacterial variants, it is important to consider phylogeny.





Bacterial Phylogeny: Tree Building Tools

Tool	Purpose	Features
IQ-TREE	Maximum likelihood phylogenetic inference	Fast, model selection, bootstrapping
RAxML/FastTree	Tree building from alignments	FastTree for large trees, RAxML for complex models



Newick Format

Statistical and Bioinformatic Methods for Studying Genetic Variants

Statistical Association Analysis

Regression models can be used to identify significant variants associated with traits of interest.

- For binary traits (e.g., disease presence, mortality):

- Logistic regression

$$\log \left(\frac{P(\text{trait} = 1)}{P(\text{trait} = 0)} \right) = \beta_0 + \beta_1 \cdot \text{SNP} + \beta_2 \cdot \text{Covariates}$$

- For quantitative traits (e.g., weight, blood pressure, WBC count)

- Linear regression

$$\text{Trait} = \beta_0 + \beta_1 \cdot \text{SNP} + \beta_2 \cdot \text{Covariates} + \epsilon$$

The results typically include a p-value and effect size.



Multiple Testing Correction

Because we are often testing many SNPs (possibly millions) at once, we need to control the false positive rate.

- Bonferroni Correction
 - Set p-value cut-off as $\alpha = \frac{0.05}{\text{number of SNPs}}$
- False Discovery Rate (FDR)
 - Benjamini-Hochberg procedure
 - Controls the expected proportion of false positives among all significant results
- Permutation Testing
 - Calculate empirical p-values based on phenotype permutation

Identifying Significant Variants: Bioinformatic Tools

Tool	Best For	Features
PLINK	General GWAS analysis	Fast; supports linear/logistic regression
GEMMA	Mixed models	Linear mixed models to correct for population structure and relatedness
SAIGE	Case-control studies with imbalance	Scalable generalized mixed models; controls for unbalanced case-control ratios and relatedness
RVTESTS	Rare variant analysis	Single variant and burden tests; useful for rare disease and sequencing data

Machine Learning in Genomic Studies

Machine Learning in Genome-Wide Association Studies

Silke Szymczak,^{1*} Joanna M. Biernacka,² Heather J. Cordell,³ Oscar González-Recio,⁴ Inke R. König,¹ Heping Zhang,⁵ and Yan V. Sun⁶

¹Institut für Medizinische Biometrie und Statistik, Universität zu Lübeck, Lübeck, Germany

²Department of Health Sciences Research, Mayo Clinic, Rochester, Minnesota

³Institute of Human Genetics, International Centre for Life, Newcastle University, Central Parkway, Newcastle upon Tyne, United Kingdom

⁴Department of Dairy Science, University of Wisconsin-Madison, Madison, Wisconsin

⁵Public Health, Yale University School of Medicine, New Haven, Connecticut

⁶School of Public Health, University of Michigan, Ann Arbor, Michigan

scientific reports

OPEN

Leveraging large-scale *Mycobacterium tuberculosis* whole genome sequence data to characterise drug-resistant mutations using machine learning and statistical approaches

Siddharth Sanjay Pruthi^{1,2}, Nina Billows¹, Joseph Thorpe¹, Susana Camp¹, Jody E. Phelan¹, Fady Mohareb² & Taane G. Clark^{1,3,4,5}

Machine learning and feature extraction for rapid antimicrobial resistance prediction of *Acinetobacter baumannii* from whole-genome sequencing data

Yue Gao^{1,2}, Henan Li², Chunjiang Zhao², Shuguang Li², Guankun Yin² and Hui Wang^{1,2*}

Article | Published: 08 December 2022

A generalizable deep learning framework for inferring fine-scale germline mutation rate maps

Yiyuan Fang, Shuyi Deng & Cai Li

communications biology

ARTICLE

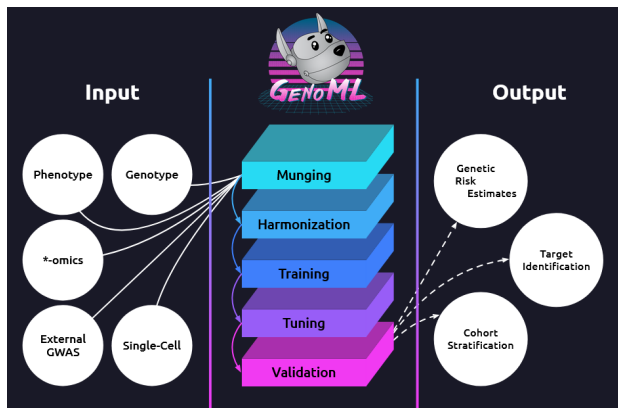
<https://doi.org/10.1038/s42003-021-02622-z>

OPEN

GenNet framework: interpretable deep learning for predicting phenotypes from genetic data

Arno van Hilten^{1,2}, Steven A. Kushner², Manfred Kayser³, M. Arfan Ikram⁴, Hieab H. H. Adams^{1,5}, Caroline C. W. Klaver^{4,6}, Wiro J. Niessen^{1,7,8} & Gennady V. Roshchupkin^{1,4,8,9}

- Machine learning has become an essential tool in genomic studies due to its ability to process and extract patterns from massive, complex datasets.
- Open-source projects and data sharing have made incorporation of machine learning easier.





Model Training Workflow

1) Preparing the data

Sample	SNP1	SNP2	SNP3	SNP4	SNP5	...
DRR034340	1	0	0	0	0	...
DRR034341	0	0	0	1	0	...
DRR034342	0	1	0	0	0	...
DRR034343	1	0	0	1	0	...
...

51,229 samples x 9,643 SNPs

Test set
(20%)

Training set
(80%)

2) Training the model

Tune
Parameters

- Random Forest
- Max_depth, min_samples_leaf, min_samples_split, n_estimators
- Gradient-Boosted Trees
- N_estimators, learning_rate, max_depth

Perform cross
validation

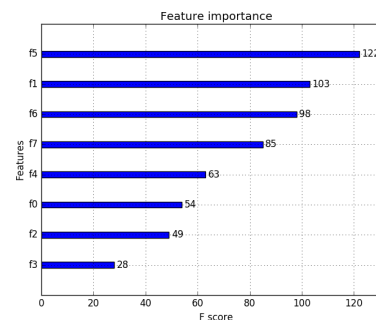
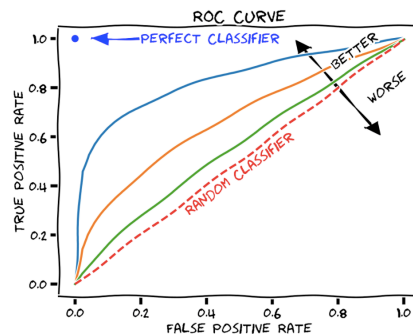
- Stratified K-Fold cross-validator
- Test model with multiple subsets of the training data using varying thresholds

Find optimal
threshold

- Test thresholds between 0.3 and 0.8
- Choose threshold with best accuracy score

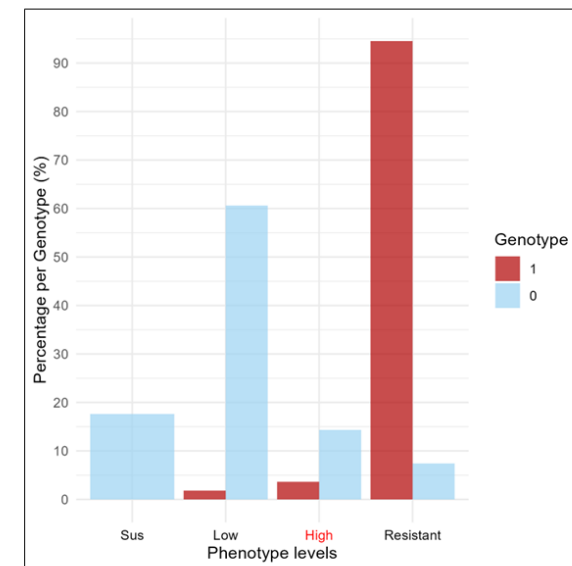
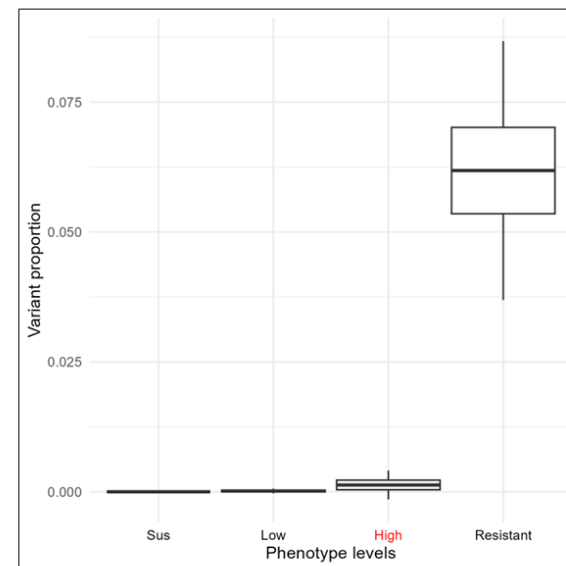
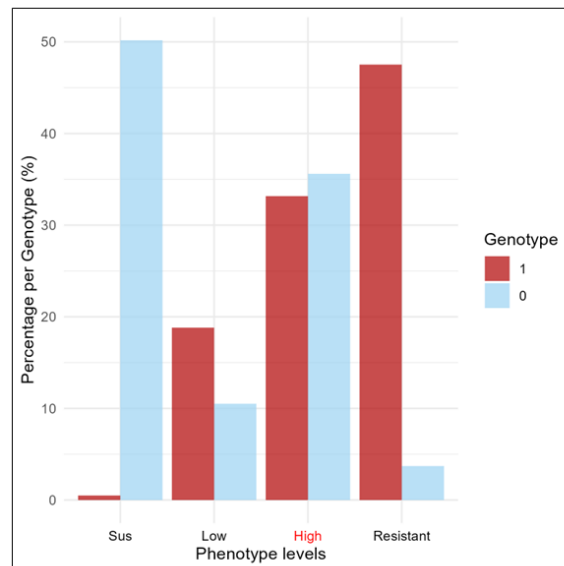
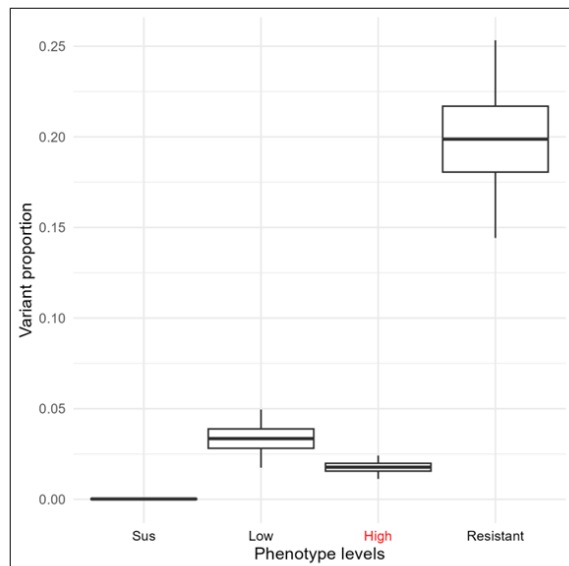
Trained
model

3) Evaluating the model



Genomic Associations with MIC Phenotypes

Drug	Position (gene)	Variant*	Suscept. Coeff.	Low Coeff.	Resist. Coeff.	Resist. P-value	Resist. RR
RIF	763,555 (<i>rpoB</i>)	230 C > T	– 0.358	0.474	1.913	$< 3 \times 10^{-7}$	6.773
INH	1,673,423 (<i>inhA</i>)	– 779G > T**	– 4.843	1.050	4.192	$< 10^{-6}$	66.168
EMB	4,248,002 (<i>embB</i>)	1489 C > A	– 4.584	– 1.930	2.164	$< 3 \times 10^{-10}$	8.707
AMK	1,473,246 (<i>Rv1313c</i>)	– 3741T > C **	– 6.422	– 4.975	2.320	$< 2 \times 10^{-14}$	10.174
KAN	1,473,246 (<i>Rv1313c</i>)	– 3741T > C **	– 3.684	– 3.063	3.915	$< 10^{-6}$	50.165
ETH	1,674,263 (<i>inhA</i>)	62T > C	– 14.794	0.380	3.457	$< 2 \times 10^{-6}$	31.708
LEV	7581 (<i>gyrA</i>)	280G > A	– 2.933	– 3.774	2.746	$< 10^{-6}$	15.575
MOX	7581 (<i>gyrA</i>)	280G > A	– 3.340	– 3.694	2.221	$< 10^{-6}$	9.215



Summary

- Genetic variants in bacteria can result in antibiotic resistance, which can hinder clinical treatments and increase risk of mortality.
- Variant calling (GATK, bcftools) and annotation (SnEff) tools are used to identify and interpret genetic variants in query sequences.
- Statistical association tests, bioinformatic tools, and machine learning help to identify key variants associated with traits of interest.

Thank you for listening!