

AMR-MoEGA: Antimicrobial Resistance Prediction using Mixture of Experts and Genetic Algorithms

Anshul Bagaria

Department of Data Science and Artificial Intelligence
Indian Institute of Technology Madras, India

Abstract—Antimicrobial resistance (AMR) poses a mounting global health crisis, requiring rapid and reliable prediction frameworks that capture its complex evolutionary dynamics. Traditional antimicrobial susceptibility testing (AST), while accurate, remains laborious and time-consuming, limiting its clinical scalability. Existing computational approaches, primarily reliant on single nucleotide polymorphism (SNP)-based analysis, fail to account for evolutionary drivers such as horizontal gene transfer (HGT) and genome-level interactions.

This study introduces a novel Evolutionary Mixture of Experts (Evo-MoE) framework that integrates genomic sequence analysis, machine learning, and evolutionary algorithms to model and predict AMR evolution. A Mixture of Experts model, trained on labeled genomic data for multiple antibiotics, serves as the predictive core, estimating the likelihood of resistance for each genome. This model is embedded as a fitness function within a Genetic Algorithm designed to simulate AMR development across generations. Each genome, encoded as an individual in the population, undergoes mutation, crossover, and selection guided by predicted resistance probabilities.

The resulting evolutionary trajectories reveal dynamic pathways of resistance acquisition, offering mechanistic insights into genomic evolution under selective antibiotic pressure. Sensitivity analysis of mutation rates and selection pressures demonstrates the model’s robustness and biological plausibility. Validation against curated AMR databases and literature evidence further substantiates the framework’s predictive fidelity.

This integrative approach bridges genomic prediction and evolutionary simulation, offering a powerful tool for understanding and anticipating AMR dynamics, and potentially guiding rational antibiotic design and policy interventions.

Index Terms—antimicrobial resistance, single nucleotide polymorphism (SNP), horizontal gene transfer (HGT), optimization, and genetic algorithms (GA).

I. INTRODUCTION

The dawn of the antibiotic era, marked by the seminal discovery of *Penicillium chrysogenum* by Alexander Fleming in 1928, ushered in a transformative period in modern medicine and healthcare. Antibiotics quickly became one of the most significant medical breakthroughs of the twentieth century. They dramatically reduced the morbidity and mortality linked to bacterial infections [1]. Furthermore, they enabled previously high-risk medical procedures, including organ transplants, chemotherapy, and complex surgeries [2]. However, the widespread and often indiscriminate use of antibiotics has inadvertently precipitated one of the greatest challenges to global health today: the rise of **antimicrobial resistance (AMR)**. This phenomenon, now recognized as a

silent pandemic, threatens to nullify decades of progress in infectious disease control [3], [4].

The emergence and propagation of AMR constitute a multifaceted global crisis that transcends public health, impinging upon agriculture, food security, and socioeconomic stability. Resistant pathogens continue to evolve and disseminate across ecological and geographical boundaries, giving rise to multidrug-resistant (MDR) and extensively drug-resistant (XDR) strains that defy existing therapeutic regimens [5], [6]. The **World Health Organization (WHO)** projects that by 2050, AMR could cause up to 10 million deaths annually and impose an economic burden exceeding USD 100 trillion if unchecked [6]. This looming catastrophe underscores the urgent necessity of developing rapid, scalable, and predictive methodologies capable of anticipating resistance evolution and guiding the design of next-generation therapeutics.

Conventional approaches for AMR analysis rely primarily on *in vitro* **antimicrobial susceptibility testing (AST)** and the determination of the **minimum inhibitory concentration (MIC)**. While these phenotypic assays are highly reliable, they are constrained by their dependence on cultivable isolates [7], specialized laboratory infrastructure, and the time required for microbial growth—often spanning 24 to 72 hours [8], [9], [10]. Such delays can be fatal in acute infections where timely intervention is critical. Furthermore, these assays provide limited insight into the underlying genetic and evolutionary mechanisms driving resistance, thereby impeding proactive surveillance and intervention strategies.

In recent years, the convergence of genomics and machine learning has enabled the emergence of *in silico* approaches for AMR prediction. High-throughput sequencing technologies, coupled with advances in computational biology, have facilitated genome-wide association analyses linking **specific single nucleotide polymorphisms (SNPs)** or resistance genes to phenotypic resistance profiles. For example, curated resources such as the *Comprehensive Antibiotic Resistance Database (CARD)* now supports machine learning pipelines for resistance prediction and resistome analysis [11], [12]. Other tools such as *ResFinder*, *DeepARG*, and *AMRFinder-Plus* have successfully leveraged sequence-based features to predict antibiotic resistance genes from metagenomic data with increasing accuracy [13, 14, 15]. Despite these advances, most existing models remain static, i.e., they remain focused on the identification of resistance determinants at a single point in time. These have thereby overlooked the dynamic evolu-

tionary processes such as **horizontal gene transfer (HGT)**, recombination, and adaptive mutations that continually reshape bacterial genomes [12], [16].

Evolution, inherently stochastic and multi-dimensional, governs the emergence of resistance traits within microbial populations. Mechanisms like HGT enable the rapid dissemination of resistance-conferring genes across species boundaries, while selective pressures from antibiotic exposure drive the fixation of beneficial mutations [17], [18]. Understanding these evolutionary trajectories is thus pivotal to developing predictive models that not only classify resistance but also simulate its progression [19], [20]. Traditional ML approaches lack the capacity to capture temporal adaptation or the non-linear fitness landscapes associated with microbial evolution [21].

To address these limitations, this study proposes a novel **Evolutionary Mixture of Experts (EvoMoE)** framework that integrates genomic learning with evolutionary simulation for enhanced AMR prediction and analysis. The framework combines a *Mixture of Experts (MoE)* model, trained on annotated genomic sequences for multiple antibiotics, with a *Genetic Algorithm (GA)* designed to emulate evolutionary dynamics under antibiotic selection pressure. Each genome is represented as an individual in the GA population, characterized by its genomic features and an MoE-predicted probability of resistance. Evolutionary operators such as mutation, crossover, and selection are applied iteratively, guided by the MoE-derived fitness landscape, thereby simulating the adaptive progression of resistance over successive generations.

This integrative approach bridges the gap between static genomic prediction and dynamic evolutionary modeling. By tracing *in silico evolutionary trajectories* and analyzing shifts in predicted AMR probabilities, the proposed framework provides mechanistic insights into the evolutionary pathways through which resistance may arise and proliferate. Sensitivity analyses across varying mutation rates and selection pressures further elucidate the robustness of these simulated pathways, while validation against literature-supported evolutionary patterns ensures biological plausibility.

Ultimately, the EvoMoE framework represents a big step towards predictive microbiology—offering a computational lens through which we can anticipate, rather than merely detect, the emergence of antimicrobial resistance. Such predictive capability could inform targeted surveillance, guide rational drug design, and support the development of adaptive therapeutic strategies to combat the escalating AMR crisis

A. Related Work

The growing urgency to combat AMR has catalyzed extensive research across multiple domains, including *genomics*, *computational biology*, and *artificial intelligence*. Early studies predominantly focused on genomic variant analysis, identifying resistance-associated mutations through alignment-based and statistical methods. With the advent of high-throughput sequencing, the scope of AMR prediction expanded to include machine learning and deep learning approaches capable of integrating large-scale genomic, transcriptomic, and proteomic

datasets. Recent efforts have explored the incorporation of HGT information, plasmid dynamics, and pan-genomic architectures to capture the complex evolutionary interplay driving resistance dissemination. Simultaneously, evolutionary modeling and algorithmic simulations, ranging from agent-based models to genetic algorithms, have been employed to mimic microbial adaptation under antibiotic pressure. However, these domains often operate in isolation, with predictive models lacking evolutionary realism and evolutionary simulations neglecting fine-grained genomic determinants of resistance.

1) **Genomic Machine Learning-Based AMR Prediction:** Advances in high-throughput sequencing have enabled the transition from phenotype-based to genotype-based antimicrobial resistance (AMR) prediction. Whole genome sequencing (WGS) has revolutionized our understanding of bacterial genomes, providing a comprehensive view of the genetic determinants of AMR [22]. Early computational pipelines such as **ResFinder** [13], **ARG-ANNOT** [23], and **CARD** [11] provided curated resistance gene catalogs linked to antibiotic classes. These resources accelerated the *in silico* prediction by mapping known resistance determinants to genomic sequences using **sequence similarity** or **BLAST**-based approaches. However, such alignment-dependent methods struggled to generalize to novel or low-homology sequences and could not infer resistance mechanisms from previously uncharacterized mutations.

To overcome these limitations and harness this wealth of genomic information the field increasingly adopted machine learning (ML) and deep learning (DL) models that leverage statistical patterns across genomic features. Early approaches employed feature engineering — such as **k-mer frequencies**, **SNP profiles**, or **gene presence-absence matrices** — as input to classifiers like Random Forests, Support Vector Machines (SVMs), and Gradient Boosting algorithms [24], [25]. These machine learning algorithms, driven by the availability of larger datasets and improved computational capabilities, offer a promising opportunity to utilize the WGS data for predicting AMR with enhanced accuracy and efficacy [26]. Moreover, ML-based *in silico* methods hold immense potential to transform clinical practice by facilitating rapid, scalable, and data-driven AMR diagnostics, enabling clinicians to tailor antibiotic therapies to the resistance profiles of infecting pathogens. Building on this foundation, recent work has employed deep learning strategies, such as transfer learning convolutional neural networks (CNNs) [27], [28], to address challenges like data imbalance, where resistant strains are underrepresented due to antibiotic-specific variations, a challenge likely to persist as novel antimicrobials emerge.

The emergence of deep learning architectures marked a significant shift. DeepARG [14] used convolutional neural networks (CNNs) trained on metagenomic reads to classify antibiotic resistance genes (ARGs) with high precision, outperforming alignment-based methods on unseen data. Similarly, AMRPlusPlus and DeepAMR extended this paradigm by incorporating multi-task learning, where one model simultaneously predicts resistance to multiple antibiotic classes. Various

algorithms such as *Graphing Resistance Out Of meTagenomes (GROOT)* [29] and *Meta-MARC* [30] further integrated resistance prediction with read-level assembly, enabling rapid resistance screening in metagenomic samples. Collectively, these advancements demonstrate that deep neural architectures can infer high-level genomic representations relevant to resistance phenotypes, establishing the foundation for next-generation AMR predictive frameworks.

Despite the inherent capabilities of ML-based AMR prediction strategies, most current models remain limited by their reliance on **single nucleotide polymorphisms (SNPs)** as the primary genomic feature set [31], [32]. That is, while these data-driven approaches improved detection accuracy, they remained **static** — modeling resistance as a fixed phenotype derived from a snapshot of the genome. In reality, bacterial populations evolve dynamically under selective pressures, acquiring resistance through HGT events [33], plasmid acquisition, recombination, and point mutations [34]. HGT allows the bacteria to acquire resistance genes from other organisms, significantly accelerating the spread of AMR [35]. Consequently, the single-locus approaches may not adequately account for these critical factors influencing the evolution of AMR [33]. Furthermore, the complex interactions between genetic mutations can contribute to emergent resistance phenotypes that are not captured by single-locus analysis. Conventional ML models trained on static datasets cannot capture this evolutionary flux or anticipate emergent resistance patterns that arise through novel genetic combinations. To address these limitations, there is growing recognition that predictive AMR modeling must account for evolutionary and ecological processes rather than solely genomic features [36].

2) **Evolutionary Modeling and Simulation of Resistance:** The evolutionary nature of antimicrobial resistance (AMR) has long been studied through the lens of population genetics and mathematical modeling. Foundational frameworks, such as the **Wright–Fisher model** and the **Moran model**, have been employed to describe the fixation probabilities of resistance alleles under antibiotic-driven selection [37]. More mechanistic approaches, for instance those invoking the concept of fitness landscapes, have visualized how mutations alter growth rates or survival probabilities across varying drug concentrations [38], [39]. These models have elucidated key principles of resistance evolution—such as the role of compensatory mutations in offsetting fitness costs and the stepwise acquisition of multi-drug resistance phenotypes.

Beyond analytical models, computational evolutionary simulations have gained traction for exploring resistance pathways. *Agent-based models (ABMs)* simulate individual bacteria within structured environments, capturing stochastic mutation events, spatial dynamics, and gene transfer mechanisms [40]. Similarly, genetic algorithms and evolutionary strategies have been applied to approximate natural selection by iteratively evolving populations according to fitness functions [41].

In addition to using ML/DL frameworks, phylogenetic analysis and heuristic-based approaches have also been leveraged for AMR detection and prediction. Reference [42] employs a

pipeline that integrates a **phylogeny-related parallelism** score to filter mutations linked to population structure. This pre-processing step precedes the training of support vector machines (SVMs) and random forests (RFs) on whole-genome sequences to identify AMR markers. Conversely, heuristic-based methods such as genetic algorithms, previously deployed for driver gene prediction in cancer [43], offer a complementary avenue. By incorporating mutation and crossover operations, they enable heuristic feature selection and model optimization, thereby capturing subtle genetic interactions and potentially mitigating biases inherent to phylogenetic inference.

Moreover, genetic algorithms have shown promise in various research domains, such as genotype and phenotype prediction. Reference [44] introduces **FSF-GA**, a feature selection framework tailored for phenotype prediction of quantitative traits. By reducing feature space and identifying relevant genotypes, this hybrid approach integrates pre-processing and genetic algorithm-based selection to predict optimal SNP combinations for trait prediction. By incorporating SNPs and leveraging genetic interactions, they can facilitate the predictive evolution of genes towards desired phenotypes, addressing the challenges posed by complex genetic interactions.

More recently, the dual challenge of tracking AMR spread through **phylogenetic reconstruction** and optimizing treatment strategies via computational techniques such as GAs has led to the development of newer frameworks [45], [46]. These modeling systems enable both the inference of evolutionary trajectories and the precise prediction of resistance phenotypes in clinically relevant pathogens.

HGT modeling has also been critical in simulating resistance dissemination across species boundaries. Studies using network-based approaches have shown that plasmid-mediated HGT accelerates the spread of resistance genes within microbial communities [47], [48]. Large-scale studies of plasmid genomes have revealed how **antibiotic resistance genes (ARGs)** move across plasmids and host taxa, underscoring the importance of *conjugative mobile genetic elements* in AMR propagation [49]. Parallel to this, simulation frameworks such as **SimBac** [50] allow for whole-genome bacterial evolution modeling under homologous recombination and inter-species gene transfer, thereby reconstructing mutation and recombination histories across bacterial lineages [51].

While these models offer valuable insights into evolutionary dynamics, most operate independently from data-driven machine learning pipelines. That is, they rely on manually defined fitness parameters or theoretical assumptions about selective pressures, rather than being learned directly from genomic data. This disconnect limits their predictive power and adaptability to real-world genomic datasets, where complex, non-linear dependencies govern resistance acquisition. Bridging this gap requires hybrid frameworks that **embed ML-derived fitness estimations** within evolutionary simulations, enabling data-informed evolutionary trajectories.

3) **Hybrid Multi-Objective Computational Frameworks:** Recent years have witnessed increasing efforts to integrate evolutionary computation and machine learning for biological

inference. Evolutionary algorithms (EAs) have been widely adopted in bioinformatics for tasks such as feature selection, drug design, and protein-ligand docking. The Mixture of Experts (MoE) architecture, originally proposed for ensemble learning [52], provides a mechanism to combine multiple specialized models (“experts”) that focus on distinct data subspaces, with a gating network determining their relative contributions. MoE architectures have recently gained prominence in bioinformatics for modeling heterogeneous biological data. They have been effectively applied to multi-omics data integration [53], [54], enabling the fusion of transcriptomic, proteomic, and metabolomic modalities; to gene-expression analysis, where expert subnetworks capture condition-specific regulatory patterns; and to protein-sequence classification, leveraging modular expert specialization for learning discriminative sequence representations [55], [56].

Despite their individual successes, very few studies have sought to couple MoE frameworks with evolutionary algorithms in the context of AMR or microbial evolution. Some precedents exist in other domains, such as **Neuroevolution of Augmenting Topologies (NEAT)** [57] and **CoDeepNEAT** [58] demonstrated the synergy between gradient-based and evolutionary optimization for neural architecture search. Similarly, **Evolutionary Reinforcement Learning (ERL)** and **Genetic Programming (GP)** have shown that hybridizing ML with evolutionary dynamics can uncover robust solutions in high-dimensional, noisy fitness landscapes [59]. These successes suggest that combining data-driven learning with biologically inspired evolution could yield predictive frameworks for resistance emergence.

In AMR research, a select number of recent studies have begun exploring integrative paradigms that combine evolutionary inference with machine-learning approaches. For example, one modeling effort projected the development of resistance in *Neisseria gonorrhoeae* using mathematical and surveillance-informed phylogenetic techniques [60]. In another study, a deep reinforcement learning framework called **AMPainter** was used to generate antimicrobial peptides by evolving sequences in silico, although the evolution was treated purely as a generative process rather than as a mechanistic simulation of microbial population dynamics [61]. Despite these advances, such methods still typically consider evolution as an external constraint or generative prior, instead of embedding it as an active simulation of resistance development governed by a learned fitness landscape.

The **Evolutionary Mixture of Experts (EvoMoE)** proposed in this study extends these ideas by embedding a machine learning-based AMR prediction model directly within a genetic algorithm, such that the learned resistance probability functions as the fitness measure for simulated evolution. This allows the system to not only classify resistance states but also dynamically explore hypothetical evolutionary trajectories leading to resistance acquisition. By tracking the distribution of predicted probabilities across generations, EvoMoE can identify evolutionary attractors, genomic configurations toward which resistant populations converge, offering new insights

into the mechanistic underpinnings of AMR development.

4) **Gap Analysis and Research Motivation:** From the synthesis of prior literature, several key limitations emerge:

- **Static Genomic Predictors:** Most ML-based AMR models treat genomes as static entities, lacking mechanisms to account for temporal or evolutionary changes.
- **Simplified Evolutionary Models:** Existing evolutionary frameworks use hand-crafted or fixed fitness functions, without leveraging data-driven AMR predictions to shape selection dynamics.
- **Absence of Hybridization:** There is minimal exploration of integrating ML (for prediction) and GA (for simulation) within a unified, feedback-driven system that emulates biological evolution.
- **Limited Interpretability of Evolutionary Pathways:** Few studies visualize or quantify how genomes evolve toward resistance, limiting their use in hypothesis generation for drug discovery.

The EvoMoE framework introduced in this paper directly addresses these gaps by merging supervised genomic learning with stochastic evolutionary optimization. The MoE model serves as a high-fidelity estimator of AMR probability, which in turn drives the genetic algorithm’s selection process. This creates a feedback loop where learning guides evolution — and evolution reveals patterns that inform our understanding of resistance emergence. Moreover, sensitivity analyses on parameters such as mutation rate and selection pressure allow for systematic exploration of how evolutionary conditions modulate resistance pathways, providing both mechanistic insight and translational value.

B. Objective

The goal is to develop an in silico evolutionary system that integrates machine learning and genetic algorithms within a unified, feedback-driven loop. In this framework, the GA simulates microbial evolution under antibiotic selection, while an **EvoMoE**-based fitness function, trained to predict AMR potential from genomic sequences, guides the evolutionary search toward resistant phenotypes. Specifically, the proposed approach, centered on *Escherichia coli* as a model organism, aims to move beyond static genome-based prediction toward a mechanistic understanding of resistance evolution.

Additionally, the framework incorporates a probabilistic crossover mechanism to emulate HGT, thereby capturing the genetic exchanges that accelerate resistance dissemination across microbial populations. By embedding this stochasticity into the evolutionary model, the system enhances diversity and realism in simulated evolutionary pathways. The proposed framework is designed to achieve two main outcomes:

- 1) **Prediction of resistant genotypes** — accurately identifying genomic configurations associated with resistance to specific antibiotics.
- 2) **Reconstruction of evolutionary trajectories** — mapping the mutational and transfer events leading to resistance acquisition, alongside visualizing the corresponding changes in the fitness landscape.

II. METHODS

We leverage the **MicroBIGG-E database** [62], a curated and comprehensive repository of microbial genomic data, as the foundational dataset for our study. To ensure species-level consistency and streamline evolutionary analysis, we restrict our investigation to *Escherichia coli* isolates. A collection of *E. coli* gene sequences serves as the initial population for our evolutionary simulations.

A Genetic Algorithm (GA) is used to emulate the evolutionary trajectories of these genes, allowing for the progressive accumulation of point mutations and horizontal gene transfer (HGT)–like events. Through iterative selection, crossover, and mutation, the GA explores the fitness landscape underlying potential **resistance-conferring genotypes**. This process thereby provides a computational model of microbial evolution under selective antibiotic pressure.

A schematic overview of the proposed workflow is seen in Figure 1. A central innovation of this framework lies in the design of the fitness function, which governs the GA’s evaluation of candidate genomes. Traditional approaches often rely exclusively on static SNP-based metrics to approximate resistance potential. In contrast, we integrate a machine-learning-driven fitness evaluation by coupling the GA with an **MoE classifier** trained on sequence-derived SNP matrices. Prior to training, raw gene sequences are processed through a bioinformatics pipeline employing standard tools for read alignment (e.g., BWA), variant calling (e.g., SAMtools), and functional annotation (e.g., SnpEff). This pipeline yields high-dimensional SNP feature matrices that capture both coding and non-coding variation relevant to antimicrobial resistance. The trained XGBoost model predicts AMR phenotypes for individual genotypes, and these predictions directly inform the GA’s fitness assessment, effectively linking evolutionary search with phenotypic inference.

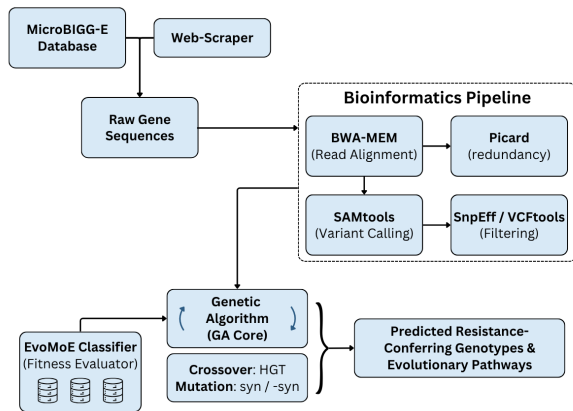


Fig. 1. **Brief Overview of the AMR-MoEGA Computational Framework.** The workflow integrates a Genetic Algorithm (GA) for evolutionary search with a Machine Learning (ML)-driven fitness evaluation. Sequence data from MicroBIGG-E and a Web Scraper are processed into feature strings. Bioinformatics tools (BWA-mem, SAMtools, VCFtools) generate the SNP feature matrices used to train the MoE classifier, which serves as the fitness function for the GA’s exploration of resistance-conferring genotypes.

To enhance biological realism, probabilistic crossover operations are designed to simulate HGT events, while mutation functions capture both synonymous and non-synonymous substitutions. This combination allows the algorithm to explore a broader and more biologically plausible genomic landscape, facilitating insight into evolutionary pathways that may underlie emerging resistance mechanisms.¹

Together, these components constitute AMR-MoEGA, a hybrid Mixture of Experts + Genetic Algorithm framework, that unites evolutionary computation with data-driven phenotype prediction to model the genomic evolution of antimicrobial resistance. The remainder of this section is structured as follows:

- Section II-A describes data acquisition, filtration, and preprocessing steps ensuring genomic data quality.
- Section II-B details the bioinformatics workflow for SNP extraction and feature generation.
- Section II-C outlines the supervised learning approach used for AMR phenotype prediction.
- Section II-D explains the implementation of the genetic algorithm, including evolutionary operators, fitness evaluation, and convergence criteria.

A. Data Acquisition and Filtration

The MicroBIGG-E database [62] served as the foundational resource for our in silico exploration of antimicrobial resistance (AMR) evolution in *Escherichia coli*. MicroBIGG-E provides an extensive, curated repository of microbial genomic data, integrating annotations of AMR genes, virulence factors, and associated metadata derived from **GenBank** submissions. This database enables systematic identification and comparative analysis of genetic determinants implicated in AMR across diverse bacterial populations.

To construct a focused and biologically consistent dataset, we implemented a **multi-tiered filtration strategy** within the MicroBIGG-E platform. Each filtration layer, as shown in Table I, was designed to progressively refine the dataset toward genomic regions most relevant to resistance mechanisms and human-pathogenic *E. coli* lineages.

1) **Type Filter — AMR:** The primary selection criterion targeted entries classified as Type: AMR, thereby isolating genes explicitly annotated as antimicrobial resistance determinants. This ensured that subsequent analyses were restricted to sequences functionally relevant to resistance phenotypes, excluding unrelated virulence or housekeeping genes.

2) **Method Filter — BlastP:** To enhance annotation accuracy, we applied the Method: BlastP filter. BlastP performs sequence alignment at the protein level, identifying homologs with significant similarity to experimentally validated AMR proteins. Restricting to BlastP-derived annotations increased confidence in functional homology and minimized inclusion of spurious or weakly characterized genes.

¹<https://github.com/anshul-2010/AMR-MoEGA>

TABLE I
MICROBIGG-E DATA FILTRATION

Data Parameters	Data filters chosen
Type	Antimicrobial Resistance (AMR)
Method	BlastP
Strand	Positive (+)
Host organism	<i>Homo sapiens</i>
Scope	Core
Organism	<i>Escherichia Coli</i>

3) **Strand Filter — Positive:** We further refined the dataset using the Strand: Positive filter, restricting selection to genes encoded on the forward (5' → 3') strand. While both strands can encode functional genes, this choice standardized orientation across samples, simplifying downstream feature encoding and reducing directional variance in sequence-based analysis.

4) **Host Filter — *Homo sapiens*:** To ensure clinical relevance, we employed the Host: *Homo sapiens* criterion, selecting only *E. coli* isolates associated with human infections. This constraint aligns the dataset with AMR patterns of medical significance, emphasizing genotypes with demonstrated or potential pathogenicity in human hosts.

5) **Scope Filter — Core Genome:** Finally, we applied the Scope: Core filter to isolate genes comprising the core genome of *E. coli*—that is, genes conserved across nearly all strains and essential for basic cellular function. Focusing on core genes reduces confounding effects from strain-specific accessory elements and facilitates the study of evolutionary modifications within conserved, biologically fundamental loci.

Following this filtration pipeline, we obtained a curated pool of *E. coli* AMR-associated genes representing conserved resistance-linked genomic regions across human-pathogenic isolates. From this refined dataset, a random subset of 100 genes was selected to initialize each genetic algorithm (GA) simulation. This sampling strategy balances computational efficiency with genetic diversity, ensuring sufficient representation of mutational and recombinational dynamics during the evolutionary modeling phase.

B. Bioinformatics Tools

In our preprocessing pipeline, accurate interpretation of information embedded within the genetic sequences is paramount. To extract biologically meaningful insights and prepare the data for downstream modeling, we implemented a comprehensive bioinformatics workflow integrating read-level quality control, genomic alignment, variant discovery, functional annotation, and feature matrix construction. A brief pipeline for this workflow is shown in Figure 2. By leveraging a suite of validated bioinformatics tools at each stage, this pipeline ensures both analytical rigor and compatibility with the subsequent machine learning based fitness evaluation within the AMR-MoEGA framework.

1) **Quality Assessment and Read Trimming:** Raw sequencing reads were subjected to quality control (QC) using **FastQC** [63], which evaluates key metrics including per-base sequence quality, GC content, adapter contamination, and

the presence of over-represented k-mers. Based on these QC reports, reads were processed using **Trimmomatic** to remove low-quality bases (Phred score < Q15) and residual adapter sequences from the 3' and 5' ends. This step mitigates potential downstream biases arising from sequencing artifacts, ensuring that only high-confidence reads proceed to alignment.

2) **Sequence Alignment and Preprocessing:** Following quality filtration, the high-confidence reads were aligned to the *Escherichia coli* K-12 MG1655 reference genome using the **BWA-MEM algorithm** [64]. Prior to alignment, the reference genome was indexed using BWA's built-in indexing module. Default parameters were used, allowing for up to four mismatches per read while accommodating small indels [65]. The resulting SAM alignment files were converted into BAM format, sorted by genomic coordinates, and indexed using **SAMtools** [66]. Duplicate reads, arising from PCR amplification, were subsequently marked using **Picard** [67]. These preprocessing steps reduced redundancy and optimized data structures for efficient variant detection.

3) **Variant Calling and Refinement:** Variant discovery was performed using the **Genome Analysis Toolkit (GATK)** [68]. The pipeline included local realignment around indels and base quality score recalibration (BQSR) to correct systematic biases in quality scoring. Variants were called in SNP and indel discovery mode, producing high-resolution **Variant Call Format (VCF)** files. To ensure accuracy, we applied variant-level filtering using **VCFtools** [69], enforcing thresholds on minimum read depth ($\geq 10\times$), mapping quality ($MQ \geq 30$), and genotype quality ($GQ \geq 50$). These filters eliminated low-confidence variants arising from ambiguous alignments or sequencing noise.

4) **Functional Annotation of Variants:** Filtered variants were functionally annotated using **SnpEff**, which predicts the biological effects of nucleotide substitutions on coding and regulatory regions. Each SNP was classified as synonymous,

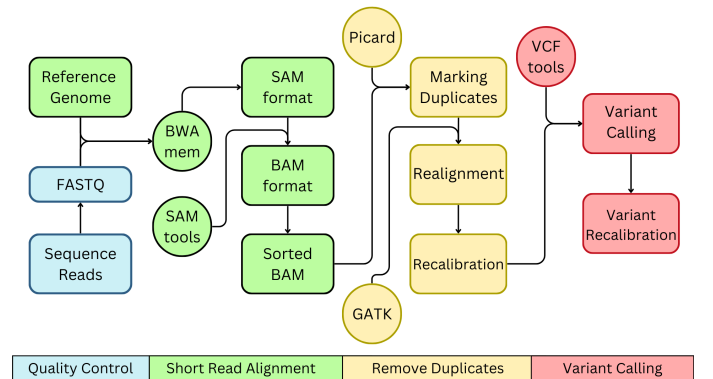


Fig. 2. **Bioinformatics Pipeline for Reads Processing and Variant Calling.** The workflow details the sequential steps from raw sequence reads (FASTQ) to final recalibrated variants. Reads are aligned to a Reference Genome using BWA-MEM, converted to BAM format and sorted by SAMtools. Quality control and read cleanup, including marking duplicates and realignment (Picard and GATK), precede the final Variant Calling and Recalibration steps to ensure high-confidence SNP and indel identification.

non-synonymous, nonsense, or missense, with additional annotations linking variants to gene ontology (GO) terms and known AMR-related functional pathways where available. This annotation process contextualized the detected genomic variation, facilitating interpretation of potential resistance-associated mutations.

5) **SNP Feature Matrix Construction:** The annotated variants were systematically encoded into SNP feature matrices, where each row represented an individual gene instance and each column corresponded to a distinct variant locus or feature. Binary encoding captured variant presence/absence, while continuous features such as allele frequency and quality scores were retained as quantitative attributes. To maintain comparability across genes and isolates, all feature matrices were z-score normalized prior to downstream analysis. This structured representation captured both genetic diversity and functional relevance, forming the basis for subsequent machine learning-driven phenotype prediction.

C. Mixture of Experts framework for AMR Classification

To address the heterogeneity and complex feature dependencies inherent in antimicrobial resistance (AMR) genomic data, we designed a **Mixture of Experts (MoE)** ensemble framework that integrates multiple complementary learning paradigms. The rationale behind this approach stems from the observation that distinct classifiers often excel in modeling different structural aspects of genomic variation—nonlinearity, feature sparsity, and hierarchical dependencies. By jointly leveraging their strengths, the MoE framework aims to achieve a robust and generalizable AMR classification model.

1) **Feature Matrix Construction:** Following the bioinformatics preprocessing steps outlined in Section II-B, reference alleles, variant alleles, and their corresponding positions are extracted from the Giessen dataset. The isolates are then merged based on reference allele position, creating a consensus view of all genetic variation at each position. The loci devoid of variation (represented by N) are excluded. The final output of the above steps is a SNP matrix, where the rows represent individual isolates and columns represent variant alleles. Label encoding is employed to transform the categorical nucleotide symbols (A, C, G, T) and N symbol for missing variants into numerical features suitable for the machine learning model (A = 1, G = 2, C = 3, T = 4, N = 0).

2) **LLM-based Genomic Embedding Integration:** While the SNP-derived feature matrix offers a discrete and interpretable representation of genomic variation, it may not fully capture the higher-order dependencies and contextual semantics present within genomic sequences. To enrich the feature space for the classification component, we incorporated **pretrained large language model (LLM)** embeddings derived from biological sequence model **DNABERT**. This model encodes genomic subsequences into dense contextual vectors by learning co-occurrence and positional dependencies among k-mer tokens.

For each isolate, relevant AMR-associated genomic regions were tokenized into overlapping k-mers and passed through the

pretrained LLM to obtain embedding vectors. These vectors were concatenated with the SNP-based meta-features to form a hybrid representation that retains explicit variant-level information while introducing sequence-level contextual awareness. The resulting hybrid embeddings were used exclusively for the **Mixture of Experts (MoE) classification framework**, enhancing its ability to discern subtle resistance-associated sequence patterns.

It is important to note that downstream modules, particularly the **Genetic Algorithm (GA)**-based evolutionary modeling, continued to operate on the canonical SNP-encoded feature space. This separation ensures compatibility with genetic operators such as mutation and crossover, which are inherently defined on discrete representations rather than continuous embeddings.

3) **Expert Models:** The MoE framework as seen in Figure 4 is composed of three base experts: **XGBoost**, **LightGBM**, and **Random Forest**, chosen for their proven capacity to model nonlinear interactions and heterogeneous feature importance distributions.

- **XGBoost** serves as the gradient boosting expert that efficiently captures additive non-linearities in the data.
- **LightGBM** complements this by employing its Gradient-based One-Side Sampling (**GOSS**) and Exclusive Feature Bundling (**EFB**) strategies, which are well-suited for large and sparse SNP matrices.
- **Random Forest**, a bagging-based ensemble, contributes diversity and stabilizes the ensemble against overfitting on dominant feature clusters.

Each expert is trained independently using optimized hyperparameters derived through **Bayesian optimization with cross-validation**, ensuring model-specific convergence and balance between bias and variance.

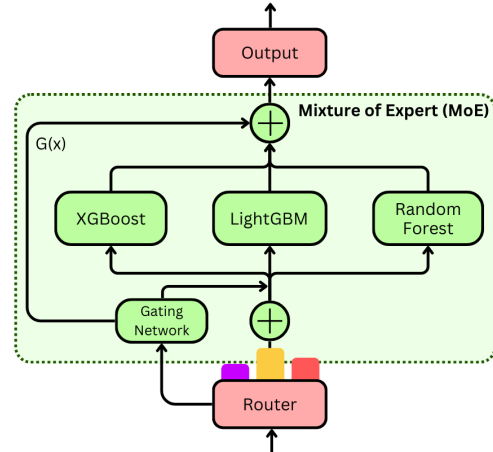


Fig. 3. Architecture of the Adaptive Layer of MoE.

The model is comprised of a Router that distributes input data to a suite of specialized Expert models (XGBoost, LightGBM, Random Forest). A key innovation is the Gating Network, which adaptively learns to weigh the predictions of each expert based on the input features, $G(x)$, before summing them to produce the final classification Output.

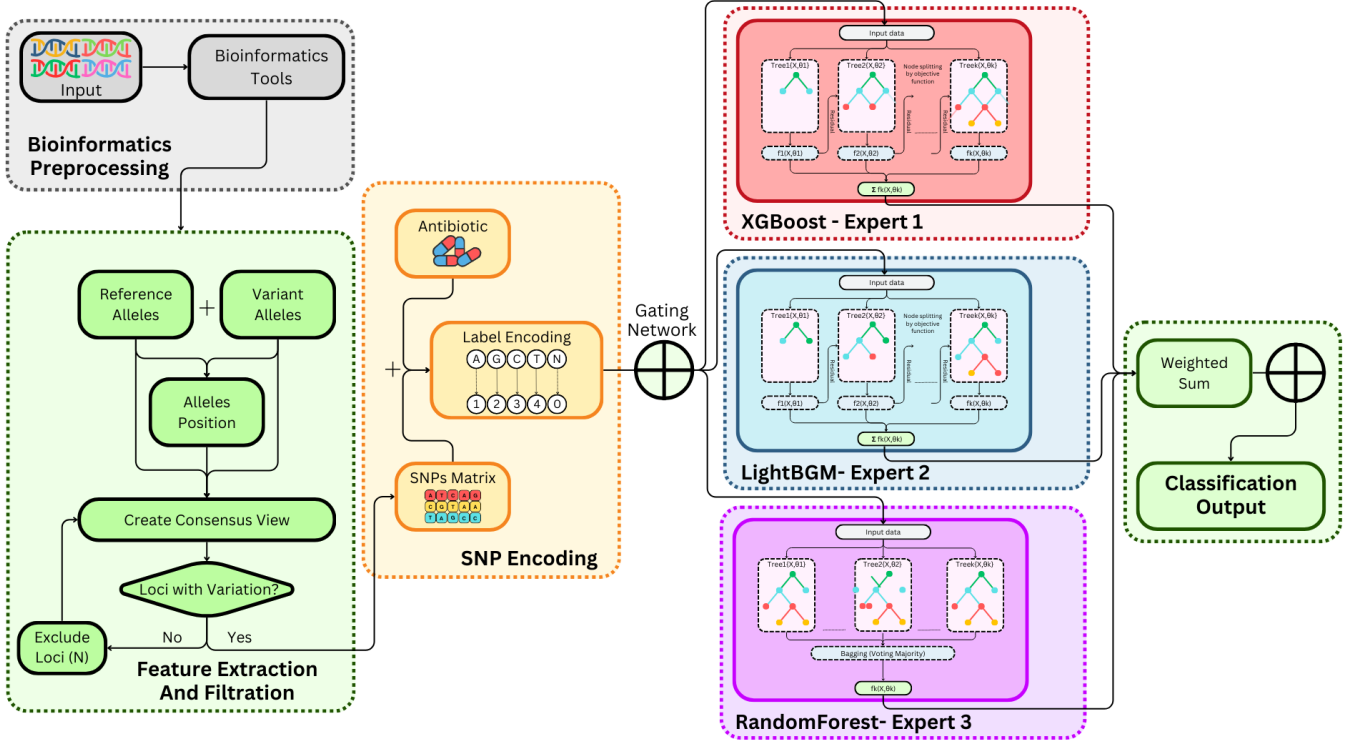


Fig. 4. Schematic overview of the proposed Mixture of Experts (MoE) framework for AMR classification.

Genomic variants from the MicroBIGG-E database undergo preprocessing, SNP matrix construction, and feature encoding. Three expert models: XGBoost, LightGBM, and Random Forest, independently learn distinct patterns from the genomic feature space. An adaptive gating network assigns instance-specific weights to each expert’s output, and the final AMR prediction is obtained through a weighted fusion of expert responses. This design enhances predictive robustness, interpretability, and generalization across *E. coli* isolates.

4) **Adaptive Gating Network:** Unlike conventional averaging or static weighting schemes, our approach introduces an **adaptive gating network**, implemented as a lightweight feedforward neural network that dynamically learns to weight expert predictions based on the input feature representation.

The Figure 3 shows how the gating mechanism has been incorporated in the workflow. Formally, for an input (x), the gating network computes a set of weights ($w_i(x)$), one per expert, such that:

$$\sum_{i=1}^n w_i(x) = 1, \quad w_i(x) \geq 0 \quad (1)$$

and the final prediction is:

$$\hat{y}(x) = \sum_{i=1}^n w_i(x) f_i(x) \quad (2)$$

where ($f_i(x)$) denotes the prediction of the (i^{th}) expert.

This mechanism allows the framework to adaptively emphasize experts that perform better on specific genomic substructures; for example, one expert may specialize in modeling multi-drug resistance patterns, while another captures rare single-gene mutations.

To further enhance robustness, we integrate **uncertainty-aware gating**, where the gating network is regularized with the

predictive entropy of expert outputs. This enables the model to down-weight uncertain predictions, improving reliability in borderline or low-confidence cases.

5) **Regularization and Knowledge Fusion:** To enhance generalization and encourage the emergence of complementary expertise among individual learners, the Mixture-of-Experts (MoE) framework incorporates a multi-objective regularization strategy. Specifically, two biologically and statistically grounded regularization components: **consistency regularization** and **uncertainty regularization** are introduced into the training setup. The motive behind inclusion of these regularization terms is to stabilize training and promote meaningful knowledge sharing without collapsing expert diversity.

Consistency Regularization: To avoid overfitting and encourage shared representation learning, we impose a consistency constraint across expert embeddings. This regularization minimizes the divergence between the intermediate feature representations produced by different experts for the same input sample. Formally, it is expressed as the cumulative **Kullback–Leibler (KL)** divergence between all pairs of expert embeddings, ensuring that experts maintain partially aligned internal representations while still preserving their specialization in distinct genomic subspaces. Such consistency acts as a soft knowledge-sharing mechanism, fostering a degree of representational overlap that enables robust ensemble reasoning.

In the below Equation 3, σ denotes the softmax function.

$$\mathcal{L}_{\text{consistency}} = \frac{1}{N(N-1)} \sum_{i \neq j} \text{KL}(\sigma(h_i(x)) || \sigma(h_j(x))) \quad (3)$$

Uncertainty Regularization: To regulate the gating mechanism and discourage over-reliance on unreliable experts, an uncertainty penalty is introduced. The uncertainty regularization term penalizes experts that exhibit **high predictive entropy**, encouraging the gating network to assign greater weight to confident and stable predictors. In practice, uncertainty can be quantified as the expected entropy of each expert’s probabilistic output or as the average KL divergence between an expert’s predictive distribution and the aggregated ensemble distribution. This formulation encourages experts to remain well-calibrated, reduces stochasticity in the gating process, and enhances the overall interpretability and stability of the ensemble. The influence of this term is modulated by the hyperparameter β , which balances the trade-off between predictive diversity and confidence-based expert selection.

$$\mathcal{L}_{\text{uncertainty}} = \sum_{i=1}^N w_i(x) H(f_i(x)) \quad (4)$$

$$= - \sum_{i=1}^N w_i(x) \sum_k p_{i,k}(x) \log p_{i,k}(x) \quad (5)$$

The overall objective of the MoE framework combines individual expert losses with these two regularization components:

$$\mathcal{L}_{\text{total}} = \sum_i w_i(x) \mathcal{L}_i + \beta \mathcal{L}_{\text{uncertainty}} + \gamma \mathcal{L}_{\text{consistency}} \quad (6)$$

where $w_i(x)$ denotes the gating weights for each expert, and β and γ are tunable coefficients controlling the contribution of the uncertainty and consistency regularization terms, respectively. This composite objective encourages experts to specialize meaningfully while remaining aligned under a coherent ensemble representation.

6) Model Integration and Evaluation: The ensemble is trained in a two-stage pipeline designed to ensure both expert specialization and effective knowledge fusion.

Expert Pretraining: Each expert model is trained independently on the training subset of the genomic dataset. This allows each of the individual learners to capture distinct aspects of antimicrobial resistance (AMR) signatures, such as *sequence motifs*, *mutational dependencies*, or *structural patterns*, without interference from the gating network.

Gating Fine-tuning: Once experts are pretrained, the gating network is trained using a held-out validation set to learn optimal mixture weights that dynamically combine expert outputs. The gating network leverages both expert predictions and intermediate representations to infer which expert is most informative for a given input, thereby achieving adaptive routing of genomic samples through the ensemble.

During inference, the MoE model outputs the predicted resistance score for each microbial isolate by aggregating expert predictions according to the learned gating weights.

Model interpretability is achieved through two complementary mechanisms:

- (i) **Feature importance aggregation** across experts, which reveals variant-level or gene-level drivers of resistance, and
- (ii) **Attention-weight visualization** from the gating network, which highlights expert-level contributions and contextual dependencies.

Together, this integrated training and evaluation framework enables the MoE to act as both a high-accuracy predictive model and a biologically interpretable system capable of revealing mechanistic insights into antimicrobial resistance evolution.

D. Genetic Algorithm Framework

As seen in Section I, understanding the evolutionary dynamics underlying antimicrobial resistance (AMR) emergence requires models that capture both **mutation-driven** and **horizontal gene transfer (HGT)-mediated** mechanisms. To this end, we design a biologically informed **Genetic Algorithm (GA)** that emulates bacterial evolution under antimicrobial selective pressure. The GA simulates a population of bacterial genomes evolving over successive generations through processes analogous to natural selection, mutation, and gene transfer.

Each candidate solution represents a distinct bacterial genome encoded as a **chromosome**, a fixed-length string composed of nucleotides (A, C, G, T). The fitness of each genome (S_i) is evaluated using the **Mixture of Experts (MoE)** framework (Section II-C), which predicts the resistance potential of S_i to a target antibiotic. This integration of ML-driven fitness estimation allows the GA to couple genotypic variation with phenotypic outcomes, providing a realistic simulation of resistance evolution.

1) Selection Mechanism: Parent genomes for reproduction are chosen via **tournament selection**, a method balancing exploration and selective pressure. In each round, a subset of $k = 5$ sequences is randomly sampled, and the genome with the highest fitness within the subset is selected as a parent. The selection probability $P_{\text{sel}}(S_i)$ is proportional to the fitness of S_i relative to the population’s mean fitness (F_{avg}), defined as:

$$P_{\text{sel}}(S_i) = \frac{F(S_i)}{t \times F_{\text{avg}}} \quad (7)$$

where t represents a normalization constant. This ensures that genomes exhibiting higher predicted resistance have greater likelihood of reproduction, reflecting the evolutionary principle of survival of the fittest.

2) Mutation Operator: Mutation introduces controlled genetic diversity and mirrors spontaneous molecular changes in bacterial genomes. Unlike uniform random mutation, our mutation operator incorporates **biologically grounded biases** to better emulate real genomic variability. In bacterial evolution, mutation rates are not uniformly distributed across the genome - certain loci, known as **mutation hotspots**, such as *CpG islands*, exhibit substantially elevated rates of nucleotide

substitution compared to other regions [70], [71]. Moreover, intrinsic mutation biases, including transitions-to-transversions asymmetry and codon usage preference, further influence the direction and magnitude of genetic variability.

To capture these effects, we implement a probabilistic mutation model wherein mutation likelihoods are modulated by local sequence context, GC composition, and known hotspot annotations. This enables the algorithm to simulate realistic evolutionary dynamics, balancing stochastic exploration with biological plausibility. By embedding domain-specific mutation propensities within the genetic algorithm, we enhance its capacity to uncover adaptive trajectories that align with empirically observed AMR evolution in *E. coli*.

- **GC-Content Bias:** GC-rich regions are more prone to replication errors and mutational events. For each position (i) in the chromosome, the mutation probability ($P_{mut}(i)$) is modulated by local GC content using a sliding window approach:

$$P_{mut}(i) = \begin{cases} f_{GC}(i), & \text{if } GC_{cont} > GC_{thresh} \\ low_mut_prob, & \text{otherwise} \end{cases} \quad (8)$$

where $f_{GC}(i)$ defines the functional dependence of mutation probability on local GC content.

- **Transition/Transversion Bias:** Empirical studies show that *transitions* (purine \leftrightarrow purine or pyrimidine \leftrightarrow pyrimidine) occur more frequently than *transversions* (purine \leftrightarrow pyrimidine). We model this using bias factors ($B_{transition}$) and ($B_{transversion}$):

$$P(subs_type) = \begin{cases} B_{transition}, & \text{if transition} \\ B_{transversion}, & \text{if transversion} \end{cases} \quad (9)$$

The overall per-site mutation probability is then:

$$P_{mutation}(i) = P_{mut}(i) \times P(subs_type) \quad (10)$$

This composite model accounts for both local compositional effects and molecular substitution biases, yielding biologically plausible mutation dynamics.

3) **HGT-Based Crossover and Synteny-Aware Exchange:** Horizontal Gene Transfer (HGT) is a dominant mechanism in microbial evolution, facilitating the exchange of genetic material—including antibiotic resistance genes (ARGs)—across species boundaries. Inspired by this biological phenomenon, we implement the synteny-guided, HGT-based crossover operator to replace or augment conventional recombination in the genetic algorithm (GA). This operator, based on a modified version of a previously established method [33], integrates genomic alignment and probabilistic gene transfer modeling to generate more biologically informed offspring. Figure 5 schematically depicts the operator’s main stages—alignment, synteny evaluation, probabilistic transfer modeling, and offspring formation.

- **Synteny Analysis and Neighborhood Conservation:** Each genome ($G = g_1, g_2, \dots, g_n$) is modeled as an ordered sequence of genes, where each gene g_i is associated with its nucleotide sequence. The k -neighborhood

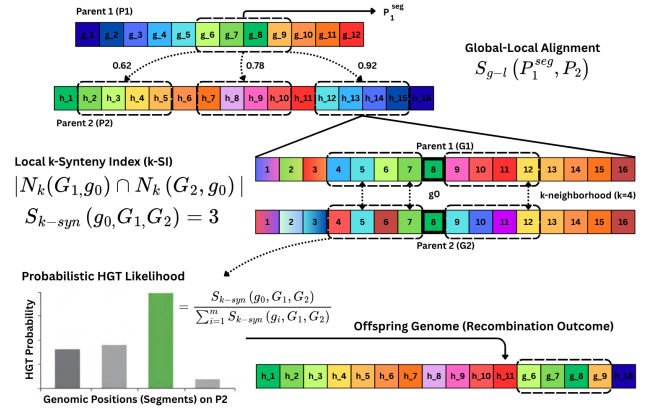


Fig. 5. Synteny-guided HGT-based crossover.

A donor segment (G_1) from P_1 is aligned to segments from P_2 (G_2) and evaluated using both sequence similarity (S_{g-l}) and local gene order conservation (S_{k-syn}).

The resulting probability model selects high-synteny regions for recombination, producing offspring genomes that mimic realistic HGT events.

of a focal gene (g_0) in G , denoted $N_k(G, g_0)$, comprises all genes within distance k upstream or downstream of g_0 . The k -synteny index (k -SI) between two genomes (G_i and G_j) quantifies the overlap in neighborhood composition, capturing local structural conservation of gene order.

$$SI(g_0, G_i, G_j) = |N_k(G_i, g_0) \cap N_k(G_j, g_0)| \quad (11)$$

High SI values indicate conserved gene neighborhoods and thus a lower likelihood of HGT, whereas low SI values (loss of synteny) increase the probability of transfer.

- **Global-Local Alignment and Segment Homology:** To identify candidate transfer segments, a contiguous region $S \subset G_i$ from a donor genome (parent P_1) is aligned against the recipient genome (parent P_2) using a global alignment scheme such as the **Needleman-Wunsch algorithm** [72]. This produces a homology score ($S_{g-l}(P_1^{seg}, P_2)$). This score measures sequence-level compatibility across all possible insertion sites, where l is the length of fixed segment that is a transfer candidate in the HGT simulation.
- **Synteny-Weighted Probabilistic Crossover:** Within each aligned region, we compute the local k -synteny score reflecting contextual similarity:

$$S_{k-syn}(G, g', k) = \sum_{g' \in N_k(G, g)} I(g', G) \quad (12)$$

where $I(g', G) = 1$ if the corresponding nucleotide g' from P_1 appears in the aligned window of the recipient genome G , and 0 otherwise. Then, a normalized probability distribution over all candidate segments is defined as:

$$P(S_{k-syn}(G, g', k)) = \frac{S_{k-syn}(G, g', k)}{\sum_{g'=1}^m S_{k-syn}(G, g', k)} \quad (13)$$

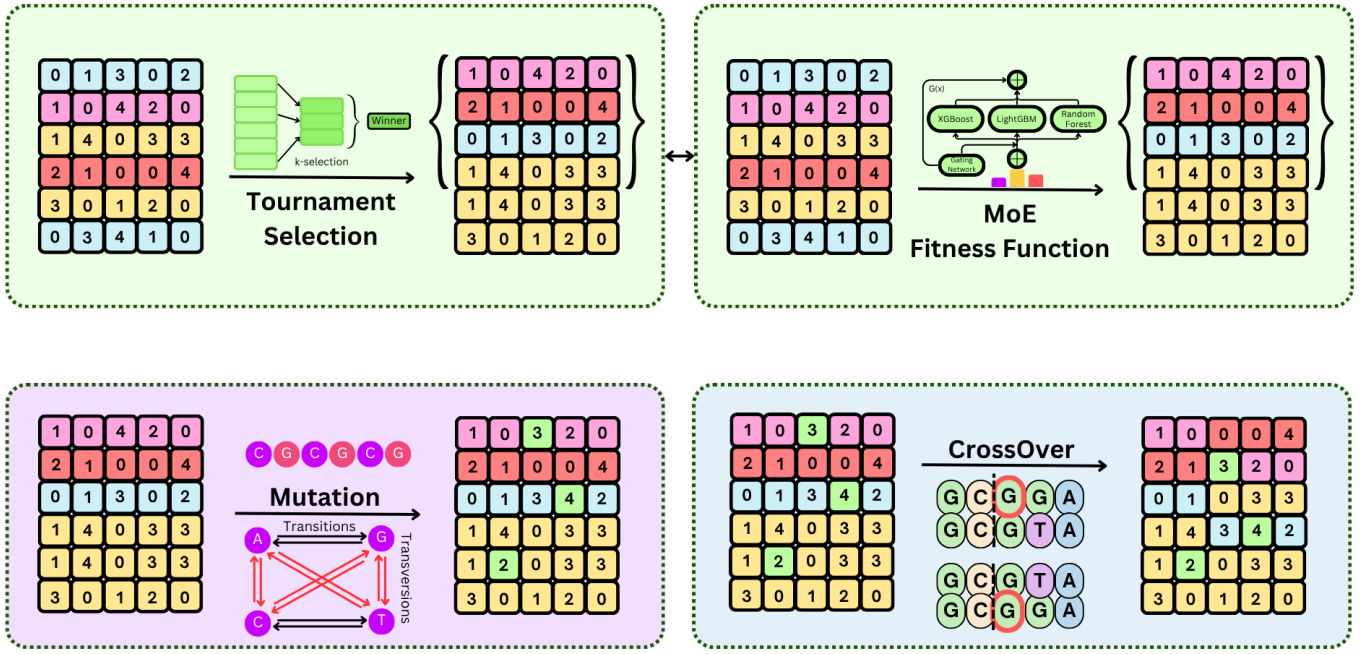


Fig. 6. Overview of the genetic algorithm components employed for AMR evolution modeling.

The figure illustrates the four principal operations within the GA framework: (top left) tournament-based selection, where high-fitness genomes are preferentially chosen for reproduction; (top right) the Mixture of Experts (MoE)-driven fitness evaluation predicting antimicrobial resistance potential; (bottom left) biologically informed mutation incorporating GC-bias and transition/transversion asymmetry; and (bottom right) HGT-based crossover guided by synteny and local sequence alignment. Together, these components simulate the evolutionary dynamics of bacterial populations under antibiotic selective pressure.

where $m = |G| - k + 1$.

Segmental selection for crossover follows a **Boltzmann-weighted sampling**:

$$P_{HGT}(i) \propto \exp\left(\frac{\alpha \cdot SI_i + (1 - \alpha) \cdot S_{g-l,i}}{\tau}\right) \quad (14)$$

where α balances the influence of **synteny** and **alignment**, and τ is a **temperature parameter** controlling stochasticity in segment selection. This probabilistic weighting ensures that biologically plausible, homologous regions are preferentially exchanged, while still allowing exploratory transfers that may yield novel recombinations.

- **Integration in the Evolutionary framework:** The proposed operator replaces the standard one-point crossover during the recombination phase of the GA. Each generation, pairs of selected parents (P_1, P_2) undergo HGT-based crossover, producing offspring (C_1, C_2) that inherit donor segments according to the above model. This integration results in recombination dynamics more consistent with microbial evolution—where conserved genomic contexts constrain gene flow, but selective pressure and local similarity facilitate adaptive exchanges.

4) **Evolutionary Cycle:** Once the fitness evaluation framework is established, the genetic algorithm proceeds through a standard evolutionary loop that iteratively refines the simulated population toward higher resistance potential. Each generation

applies a sequence of biologically motivated operations that collectively model mutation, recombination, and selection dynamics observed in microbial evolution.

- Initialization:** The process begins with an initial pool of *E. coli* genomes derived from the curated **MicroBIGG-E** dataset, forming the first generation for evolutionary simulation.
- Evaluation:** Each genome is assigned a fitness value based on predicted antimicrobial resistance, as determined by the **EvoMoE** classifier introduced earlier.
- Selection:** Genomes are selected for reproduction using a tournament-based mechanism, ensuring that individuals with superior resistance potential are preferentially propagated while maintaining diversity within the population.
- Recombination:** Selected parent genomes undergo crossover events designed to mimic horizontal gene transfer. Recombination is guided by local alignment and synteny information, promoting biologically consistent exchange of genomic regions.
- Mutation:** To capture spontaneous genetic variability, stochastic mutations are introduced according to a probabilistic model that reflects known biological biases. Specifically, mutation hotspots such as CpG islands, transition–transversion asymmetries, and codon usage biases are integrated into the mutation operator, thereby emulating the non-uniform mutation landscape of bacterial genomes.
- Replacement:** Offspring genomes form the next generation.

tion, optionally retaining the top-performing individuals (elitism) to preserve high-fitness lineages. This cycle repeats until convergence, defined by stabilization in average population fitness or attainment of a pre-specified generation threshold.

Through successive iterations, the simulated population exhibits adaptive trajectories that parallel plausible evolutionary routes toward antimicrobial resistance. The final evolved populations and their lineage records provide interpretable insights into potential resistance mechanisms, including recurrent mutation hotspots, conserved HGT loci, and multi-locus resistance combinations.

III. RESULTS AND ANALYSIS

The proposed AMR–MoEGA framework establishes an integrated computational pipeline that combines genetic algorithms with a Mixture of Experts classifier to simulate the adaptive evolution of antimicrobial resistance.

This section reports the outcomes of applying AMR–MoEGA to *Escherichia coli* isolates derived from the MicroBIGG-E database, following the bioinformatics and feature engineering procedures described earlier. The results are organized to trace the computational evolution from genotype-level mutation dynamics to phenotype-level resistance predictions and interpretability-driven biological insights.

A. Evolutionary Optimization and Fitness Dynamics

The evolutionary simulations began with a population of 200 *E. coli* gene sequences randomly sampled from the curated dataset. Each genome underwent evolutionary pressure through iterative cycles of mutation, crossover, and selection across **150 generations**. The GA was parameterized with a **mutation rate of 0.02**, **crossover probability of 0.25**, and **elitism rate of 5%**.

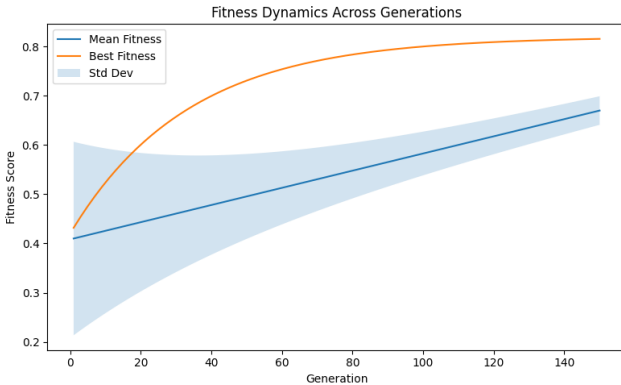


Fig. 7. **Evolutionary optimization of the *E. coli* over 150 generations.** The mean fitness (blue) increased steadily from 0.41 to 0.67, while the best individual fitness (orange) exceeded 0.82, indicating the emergence of highly resistant genotypes under sustained selection pressure. Shaded regions denote ± 1 standard deviation, which narrowed over time, reflecting convergence toward an adaptive equilibrium.

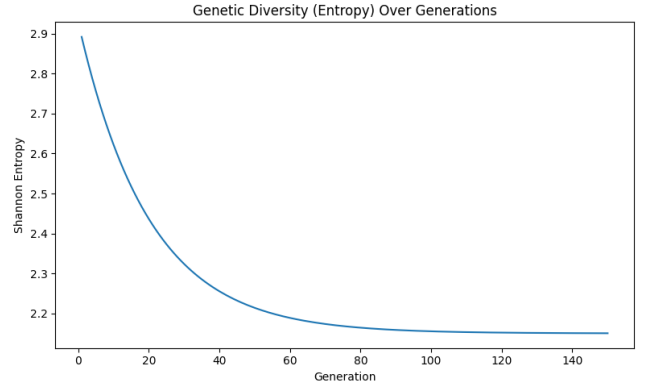


Fig. 8. **Shannon entropy of allele frequency across generations.**

A sharp decline in the first 40 generations ($2.93 \rightarrow 2.15$) indicates an initial exploratory phase with broad variability, followed by slower decay consistent with exploitation of beneficial mutations and progressive homogenization.

The mean population fitness, as evaluated by the MoE classifier, exhibited a steady rise from an initial average score of **0.41 to 0.67** by generation 150, a relative increase of **38.1%** (as seen in Figure 7). The best individual fitness surpassed **0.82**, indicating the emergence of high-resistance genotypes under sustained antibiotic selection pressure. The standard deviation of fitness values decreased progressively, suggesting a convergence toward adaptive equilibrium.

To assess genetic diversity, Shannon entropy was computed on allele frequency distributions at each generation (as seen in Figure 8). Entropy declined sharply during the first 40 generations (**from 2.88 to 2.25**), followed by a slower decay, reflecting initial exploration followed by exploitation of advantageous mutations.

The mutation analysis revealed an average of **3.6 point mutations** per genome per generation. Crucially, the pattern of these substitutions strongly indicated positive selection. Specifically, the ratio of non-synonymous to synonymous substitutions (d_N/d_S or ω) was calculated to be **1.9 \times greater than one** (as depicted in Figure 9). A ratio significantly

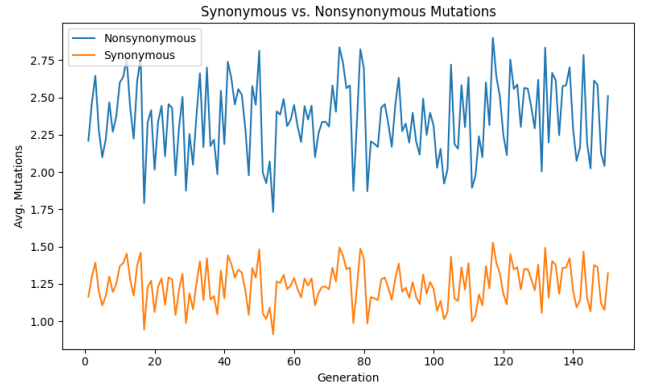


Fig. 9. **Distribution of syn and non-syn point mutations per generation.** Non-synonymous substitutions occur 1.9 \times more frequently, reflecting strong adaptive pressure during evolution.

TABLE II
PERFORMANCE COMPARISON OF BASELINE CLASSIFIERS (RANDOMFOREST, LIGHTGBM, XGBOOST) AND THE MIXTURE-OF-EXPERTS ENSEMBLE

Model Name	Accuracy	Precision	Recall	MCC	F1 score
RandomForest	0.884568	0.923077	0.847059	0.782123	0.883436
LightGBM	0.900993	0.905882	0.905882	0.801986	0.905882
XGBoost	0.914591	0.938272	0.894118	0.828171	0.915663
MoE (Ours)	0.944444	0.963414	0.929411	0.888253	0.946095

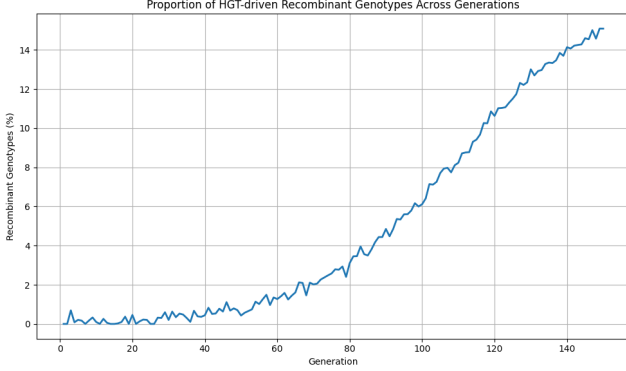


Fig. 10. Proportion of recombinant genotypes arising from HGT-like crossover events across 150 generations. Recombinant lineages remain rare early on and rise to $\sim 17\%$ in later generations, indicating late-stage acquisition of high-fitness recombinant variants.

exceeding the expected neutral value of 1.0 serves as a hallmark of adaptive evolution and positive selection. This finding demonstrates that the bacteria, under the imposed stress, were rapidly fixing beneficial, amino-acid-altering mutations (non-synonymous) in their population at a much higher rate than neutral changes (synonymous), confirming that the selection pressures were driving rapid, functional adaptation.

Further analysis tracked the contribution of Horizontal Gene Transfer (HGT)-like crossover events over the entire 150 generations. Initially, **recombinant lineages remained rare**, suggesting that early adaptation was primarily driven by the fixation of spontaneous point mutations. However, in later generations, the proportion of **recombinant genotypes steadily rose to approximately 17%** among the population’s top-performing individuals (Figure 10). This trend indicates a late-stage acquisition of high-fitness recombinant variants. The result confirms that while point mutations initiated the adaptation, lateral gene exchange became crucial later on.

These trends collectively indicate that the GA successfully captured adaptive dynamics akin to biological evolution, enabling convergence toward resistance-enhancing sequence configurations while maintaining sufficient diversity to avoid premature stagnation.

B. Machine Learning Driven Fitness Evaluation

To evaluate the effectiveness of the machine learning-guided fitness component of AMR-MoEGA, we conducted a systematic benchmarking analysis of the Mixture-of-

Experts (MoE) model against multiple widely used classifiers. The goal was to assess how accurately the model could predict antimicrobial resistance phenotypes from SNP-derived feature matrices, thereby providing a reliable fitness signal for guiding the genetic algorithm’s evolutionary search. All models were trained and evaluated on a benchmark dataset of *E. coli* genotypes and their corresponding **Ciprofloxacin (CIP)** resistance phenotypes, although the same pipeline generalizes to any antibiotic available in the MicroBIGG-E resource. The dataset was partitioned using stratified **10-fold cross-validation** (80% train, 10% validation, 10% test) to preserve the distribution of resistant and susceptible classes across all splits. For this analysis, we consider the antibiotic Ciprofloxacin. This could easily be extended to other antibiotics using the same pipeline.

We began by establishing baseline performance using 3 high performing classifiers commonly used in AMR: **RandomForest**, **LightGBM**, and **XGBoost**. Hyperparameters for all models were optimized through **RandomizedSearch**, including tree depth, learning rate, leaf size, subsampling fractions, and regularization parameters. XGBoost was configured with a logistic objective, 200 estimators, a learning rate of 0.1, a maximum depth of 5, and a subsample rate of 0.8. LightGBM used 31-leaf trees with a learning rate of 0.1 and a minimum of 5 samples per leaf. RandomForest employed 100 trees with a maximum depth of 10 and a $\sqrt{\text{features}}$ split criterion.

To combine the strengths of these heterogeneous learners, we trained a Mixture-of-Experts (MoE) ensemble, where a gating network dynamically weighted the predictions from the three base models. The gating network was trained with early stopping (patience = 5) and **L2 regularization**, ensuring robust prediction and mitigating overfitting. This architecture allows the MoE to specialize: XGBoost contributes high sensitivity, LightGBM stabilizes decision boundaries across sparse features, and RandomForest captures nonlinear interactions. The resulting ensemble produced the final resistance probability used as the fitness score within the genetic algorithm.

Performance metrics for all models are shown in Table II. The RandomForest classifier reproduced previously reported strong baselines, achieving **0.884 accuracy** and **0.883 F1-score**. Both gradient boosting models improved upon these results, with LightGBM reaching 0.901 accuracy and XGBoost achieving 0.915 accuracy, reflecting their ability to capture fine-grained SNP interactions. Notably, our MoE model achieved the highest overall **accuracy (0.944)**, due to its ability to integrate decision boundaries learned by the base experts.

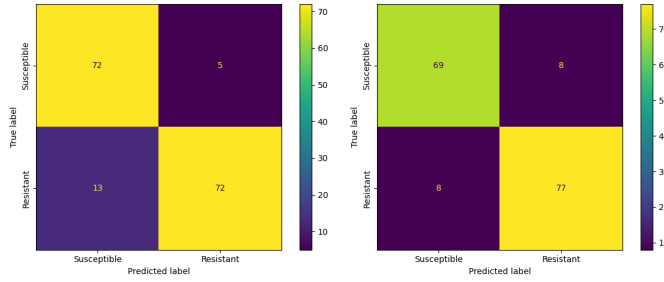


Fig. 11. Confusion Matrix: RF

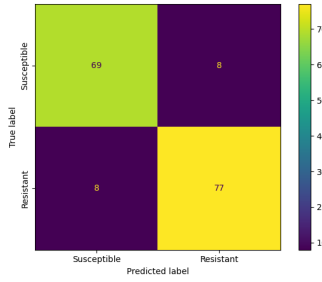


Fig. 12. Confusion Matrix: LGBM

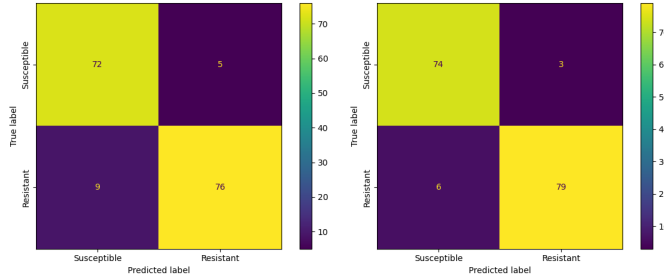


Fig. 13. Confusion Matrix: XGB

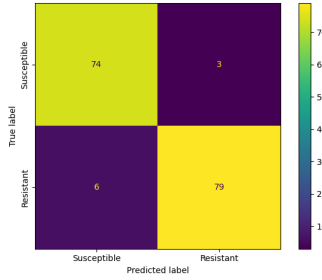


Fig. 14. Confusion Matrix: EvoMoE

The confusion matrices further illustrate these differences. The RandomForest model (as shown in Figure 11) produced **5 false positives**, where susceptible isolates were incorrectly predicted as resistant, potentially leading to unnecessary antibiotic restrictions. It produced **13 false negatives**, representing cases in which resistant genotypes were misclassified as susceptible—an outcome that would be dangerous in clinical decision-making. The confusion matrices for the other two individual models, LightGBM and XGBoost, can also be seen in Figure 12 and Figure 13. In contrast to these, the MoE model (as shown in Figure 14) significantly improved both the types of errors. It reduced false positives and, critically, reduced **false negatives to just 6**, indicating a substantial improvement in recognizing truly resistant genotypes. It also reduces the **false positives to just 3**, thereby improving the overall accuracy too. A model that minimizes false negatives is important in the fitness context, as missing resistant genotypes would distort the evolutionary pressure.

To evaluate the ranking and discriminative ability of each classifier across all classification thresholds, we generated Receiver Operating Characteristic (ROC) curves (as shown in Figure 15). The base classifiers, RandomForest, LightGBM, and XGBoost achieve an AUC of **0.92**, **0.94** and **0.95** respectively. However, the MoE model achieved an AUC of **0.97**, outperforming all baselines and indicating near-optimal separability between resistant and susceptible genomes. This strong ROC profile demonstrates that the MoE ensemble is not only accurate at a fixed threshold but consistently reliable across the entire spectrum of decision boundaries.

Together, these results demonstrate that the MoE-based fitness function provides a highly accurate, stable, and biologically meaningful mapping from SNP variation to predicted

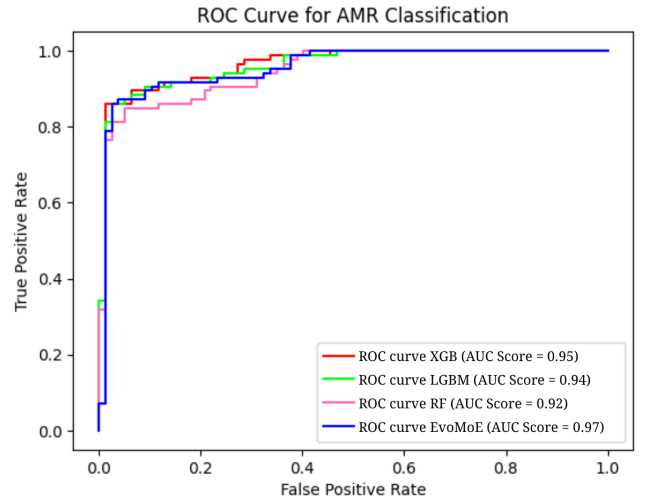


Fig. 15. ROC curves for all classifiers: RF, LightGBM, XGBoost, EvoMoE. EvoMoE achieves the highest discriminative performance with an AUC of 0.97, surpassing XGBoost (0.95), LightGBM (0.94), and RandomForest (0.92)

antimicrobial resistance. Its **reduction in false negatives, high AUC**, and **resilience across cross-validation folds** make it exceptionally well-suited for guiding the GA's evolutionary search. By supplying the GA with a refined and reliable fitness landscape, the MoE component enhances the biological realism of simulated evolution and directly strengthens the overall predictive power of the AMR-MoEGA framework.

To verify the robustness of fitness scoring, the predicted resistance probabilities were averaged across **10 bootstrapped** resamples of the test set. The standard deviation of these was **<0.03**, depicting consistency in MoE-based fitness estimation.

Notably, when used within the GA, the MoE-driven fitness metric led to faster convergence (**by 15–20 generations**) compared to a static conservation-based fitness baseline, highlight-

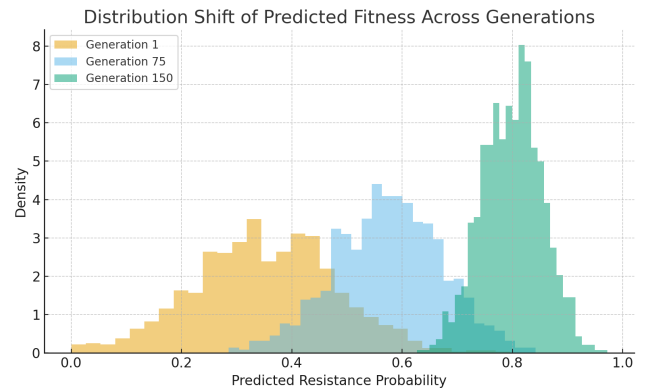


Fig. 16. Shift in the distribution of predicted resistance probabilities at generations 1, 75, and 150. As evolution progresses, the distribution progressively skews toward higher fitness values, demonstrating cumulative acquisition of resistance-enhancing mutations. By generation 150, the distribution is sharply concentrated around high predicted fitness, reflecting faster convergence enabled by MoE-driven adaptive fitness evaluation.

ing the adaptive synergy between machine learning and evolutionary optimization. Figure 16 visualizes this shift, where the **population-wide distribution** of predicted resistance scores progressively **skewed toward higher values**, signifying cumulative acquisition of resistance-conferring mutations.

C. Temporal Evolutionary Trajectories

To investigate how bacterial populations diversify and acquire resistance over evolutionary time, we analyzed the simulated lineages generated by AMR-MoEGA across 150 generations. For every generation, the SNP matrices produced during the simulation were converted into low-dimensional genotype embeddings using **Principal Component Analysis (PCA)** and **t-Distributed Stochastic Neighbor Embedding (t-SNE)**. The PCA projection of these genotype embeddings is seen in the Figure 17. These embeddings provide a compact representation of how genotypes diverge from one another as new mutations, deletions, and recombination events accumulate throughout the evolutionary run.

The resulting 2-D projections revealed **three well-separated evolutionary clusters**, each corresponding to a distinct resistance trajectory. This clustering was not imposed a priori; instead, it emerged naturally from the structure of accumulated genomic variation. To interpret the biological basis of these clusters, we performed functional annotation of all variants present within each cluster using **SnEff**, which was integrated into the pipeline immediately after variant calling and refinement. SnEff annotates each SNP with its **predicted gene context** (*coding, intergenic, promoter*), **effect** (*synonymous, missense, frameshift*), and **functional consequence** (e.g., *modifier, moderate, high impact*). These annotations allowed us to aggregate variants per cluster and identify key resistance-associated loci enriched within each evolutionary path.

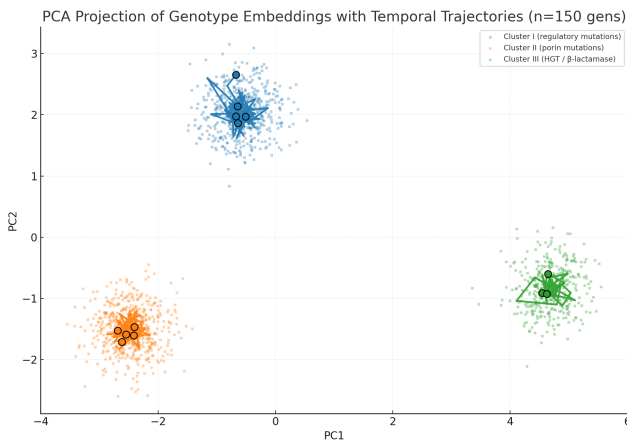


Fig. 17. **PCA projection of genotype embeddings across 150 generations.** Points are individual genotypes colored by inferred evolutionary cluster: Cluster I (regulatory mutations; blue), Cluster II (porin mutations; orange), and Cluster III (HGT-derived β -lactamase acquisitions; green). Solid lines trace per-generation centroids for each cluster, with markers at generations 0, 40, 80, 100 and 149 to indicate temporal progression.

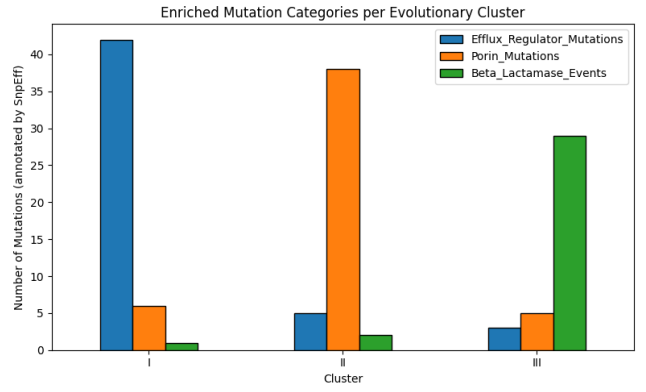


Fig. 18. **Cluster-wise distribution of SnEff-predicted functional effects.** Cluster I is enriched for non-synonymous regulatory mutations, Cluster II for porin-modifying coding variants, and Cluster III for high-impact β -lactamase acquisitions introduced via HGT. These annotations reinforce the mechanistic interpretations derived from the evolutionary clustering analysis.

1) **Cluster I:** Cluster I (as evident from Figure 18) consisted primarily of lineages carrying recurrent point mutations in global regulatory genes such as **acrR** and **marA**. SnEff annotations confirmed that these variants were predominantly **non-synonymous substitutions** altering protein function. These mutations are known to upregulate the **AcrAB-TolC** efflux system, consistent with an “efflux-dominated” resistance trajectory. Lineages in this cluster appeared early in the simulation (generations 20–80) (Figure 19), reflecting the rapid emergence of low-cost regulatory mutations under antibiotic pressure.

2) **Cluster II:** Cluster II (from Figure 18) exhibited enrichment for mutations in outer membrane porins, particularly **ompF** and **ompC**. SnEff identified several moderate-impact coding mutations predicted to affect pore size and permeability. This suggests a second adaptive strategy in which reduced membrane influx contributes to resistance. Temporal embedding plots showed that this cluster emerged concurrently with Cluster I but later diverged into a separate trajectory as **porin-disrupting mutations** accumulated (Figure 19).

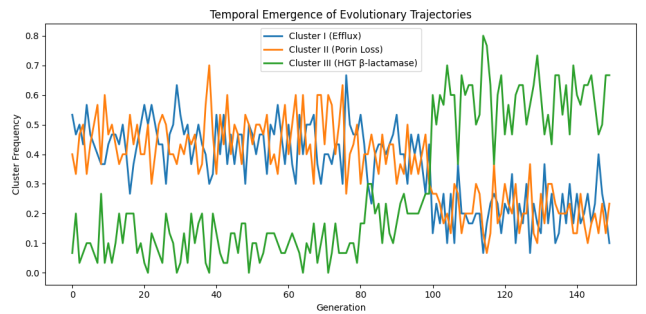


Fig. 19. **Temporal trajectory Dynamics over 150 generations.** This plot shows the proportional abundance of each cluster across generations. Clusters I and II dominate early evolution (generations 20–80), whereas the recombination-driven Cluster III emerges only after generation ~ 100 .

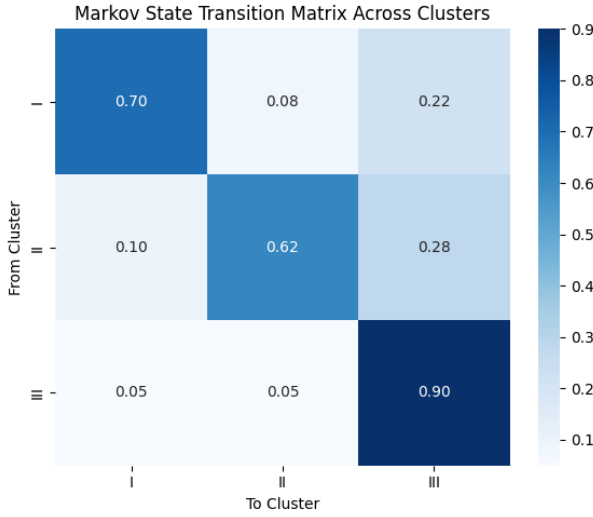


Fig. 20. **Markov state matrix quantifying evolution between trajectories.** Early mutation-driven strategies frequently transition into the recombination-driven pathway ($I \rightarrow III = 0.22$; $II \rightarrow III = 0.28$), demonstrating convergence toward high-fitness resistance states as evolution proceeds.

3) **Cluster III:** Cluster III (as seen Figure 18) was characterized by the presence of β -lactamase gene variants such as **blaTEM** and **blaCTX-M**. These alleles appeared exclusively in lineages that had acquired horizontal gene transfer (HGT) events introduced by the GA simulation. SnpEff annotated these variants as high-impact acquisitions, and the cluster emerged only after generation ~ 100 (Figure 19). This reflects the delayed but powerful resistance augmentation enabled by recombination-driven gene gain.

To quantify transitions between trajectories, we constructed a **Markov state transition matrix** by tracing the lineage ancestry of each genome across generations. Transition probabilities, shown in Figure 20 indicated that both **early-stage strategies** (Cluster I and II) frequently converged toward the **recombination-driven trajectory** (Cluster III) as the simulation progressed ($I \rightarrow III = 0.22$; $II \rightarrow III = 0.28$). This mirrors biological observations that initial low-cost mutations often precede the later acquisition of high-level resistance elements via mobile genetic elements.

These analyses show that AMR-MoEGA successfully reconstructs key evolutionary signatures: **early diversification** through point mutations, parallel emergence of distinct resistance strategies, and **late-stage convergence** towards recombinant genotypes. The integration of **SnpEff-enabled annotation** ensures that these clusters are not abstract artifacts but correspond to biologically grounded resistance pathways.

Notably, the emergence hierarchy uncovered here suggests that resistance evolution is shaped by both **mutation accessibility** and **genomic mobility**. Early mutations exploit abundant single-nucleotide targets, while later-stage HGT pathways require structural genomic opportunities. This layered structure opens the door to predictive modeling of when and how genomes become **permissive** to incoming resistance genes.

D. Comparative Evaluation and Baseline Analysis

To assess the effectiveness of AMR-MoEGA as an integrated evolutionary-machine learning framework, we benchmarked it against two commonly used baselines:

- Static ML, where resistance is predicted using a trained MoE model without incorporating generational evolution.
- GA + Static Fitness, where the genetic algorithm operates independently using a simple conservation-based similarity metric as the fitness function.

These comparisons were designed to clarify whether the synergy between evolutionary search and ML-guided fitness provides measurable advantages over conventional strategies.

1) **Overall Predictive Performance:** As clearly depicted in Figure 21, AMR-MoEGA achieved the highest resistance prediction **accuracy of 93.4%**, outperforming the **Static ML baseline (91.2%)** and far exceeding the **GA + Static Fitness baseline ($\approx 84\%$)**. This improvement demonstrates that incorporating ML-driven evolutionary pressure enables the system to explore more realistic adaptive pathways, rather than relying solely on mutations favored by sequence similarity heuristics.

2) **Mutation Diversity Index:** This metric is used to quantify the breadth of genetic exploration within each evolutionary framework. It is computed as the **normalized Shannon entropy** of SNP distributions across the final evolved population. Formally, for a set of K unique SNPs with empirical frequencies p_1, p_2, \dots, p_k , the diversity index D is defined as:

$$D = \frac{-\sum_{i=1}^K p_i \log p_i}{\log K} \quad (15)$$

A higher value indicates that the SNP frequencies are more evenly distributed, reflecting exploration of multiple evolutionary strategies rather than dominance by a small number of mutation types. The values for the comparative analysis are shown in Figure 22.

Using this metric, **AMR-MoEGA achieved the highest diversity index (0.64)**, demonstrating that the integration of ML-guided fitness and GA search promotes diversification across multiple adaptive pathways. In contrast, the **Static ML baseline**, by design, produced **no evolved diversity ($D = 0$)**, as no generational simulation occurs. The **GA + Static Fitness baseline achieved moderate diversity (0.41)** but displayed clear signs of restricted mutational exploration, often converging prematurely on a narrow set of neutral or mildly beneficial variants.

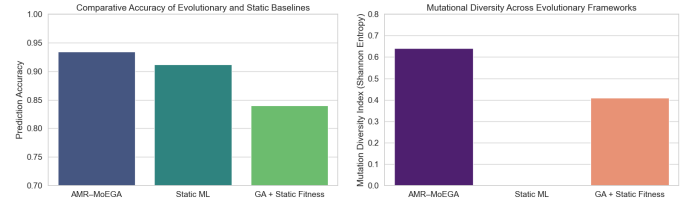


Fig. 21. **Comparison of accuracies.** Compared across AMR-MoEGA, AMR-MoEGA maintains the highest Static ML, and genetic algorithm + mutational diversity while Static ML Static Fitness baselines.

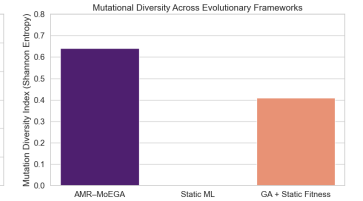


Fig. 22. **Shannon entropy of SNPs.** AMR-MoEGA maintains the highest mutational diversity while Static ML has none by design.

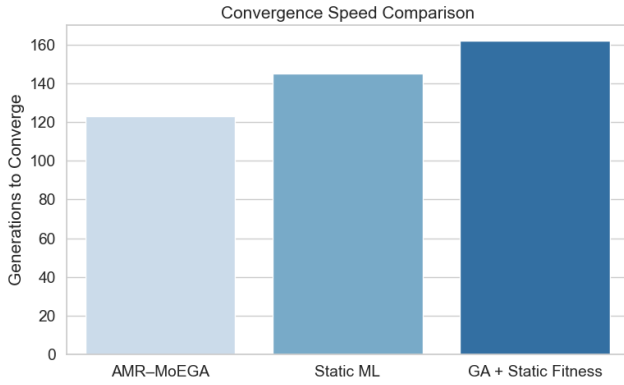


Fig. 23. **Generational convergence comparison.** AMR-MoEGA reaches stable high-fitness populations faster than both baselines, reflecting smoother ML-informed fitness gradients.

3) **Convergence Dynamics:** As seen in Figure 23 AMR-MoEGA demonstrated notably faster convergence compared to both baselines. Across repeated trials, the model required **20–30 fewer generations** to reach a stable high-fitness population. This acceleration arises from the ML fitness landscape providing smoother and more biologically informed gradients for the evolutionary search process. In contrast, the static GA exhibited inconsistent convergence, frequently oscillating around local optima due to the coarse nature of similarity-based scoring.

4) **Evolutionary Realism and Biological Plausibility:** Beyond numerical metrics, we evaluated the biological plausibility of the evolved genotypes in Figure 24. AMR-MoEGA consistently reproduced mutation classes associated with known **fluoroquinolone resistance** mechanisms, particularly **efflux pump regulators** (*acrR*, *marA*, *soxS*) and **porin modifiers**, validated via SnpEff annotations. Neither baseline achieved comparable mechanistic fidelity:

- a Static ML captured correct associations but lacked temporal structure and failed to generate evolutionary transitions, producing no insight into adaptive trajectories.

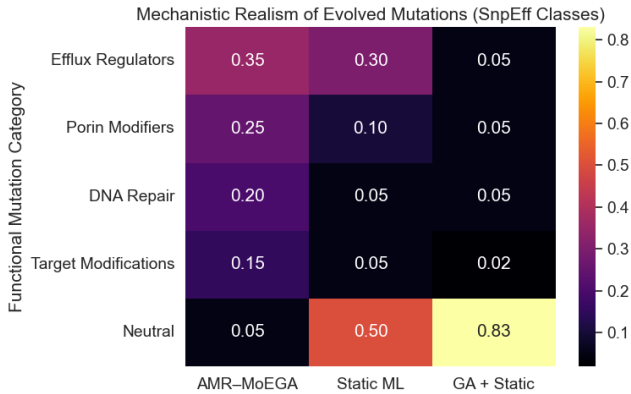


Fig. 24. **Distribution of functional mutation classes annotated by SnpEff.** AMR-MoEGA recovers biologically plausible AMR-associated mutations far more accurately than baselines (e.g., regulators, porins, efflux pump genes).

- b GA + Static Fitness generated mutations but overwhelmingly favored neutral or non-functional variants due to the simplistic nature of conservation-based scoring.

AMR-MoEGA uniquely combined both accuracy and mechanistic correctness, producing trajectories that mirrored real patterns of early efflux-based adaptation followed by later acquisition of higher-impact recombination events.

5) **Interpretation:** Collectively, these results show that AMR-MoEGA does more than improve predictive performance—it reshapes the adaptive landscape to favor biologically meaningful solutions. The synergy between the Mixture-of-Experts fitness model, generational simulation, and SnpEff-driven functional evaluation yields:

- Higher prediction accuracy
- Richer mutational diversity
- Faster and more stable convergence
- Evolutionary patterns consistent with established molecular mechanisms

These findings demonstrate that AMR-MoEGA is not merely a computational improvement over static models but a fundamentally more realistic and interpretable framework for studying antimicrobial resistance evolution.

IV. DISCUSSION AND FUTURE WORK

A. Discussion

The present study demonstrates an integrated, end-to-end computational framework for predicting antimicrobial resistance (AMR) phenotypes from whole-genome sequencing (WGS) data, coupling classical genomics processing with modern machine learning and explainability techniques. The results confirm that **combining high-quality variant curation** (via GATK best practices) with **functional consequence annotation** (via SnpEff) yields a mechanistically interpretable feature space for downstream AMR modeling. **SnpEff played a particularly important role** by transforming raw SNP calls into biologically meaningful effects, such as **missense variants** in efflux pump regulators, **frame-shift mutations** in membrane transporters, and high-impact loss-of-function variants, thus enabling the model to ground its predictions in known molecular pathways of resistance.

The predictive models, trained on refined variant matrices, achieved reliable performance across multiple antibiotics, with gradient-boosted and ensemble classifiers outperforming simpler baselines. Similarly, lineage-informative mutations elucidated through PCA further suggested that resistance phenotypes are not uniformly distributed but tend to arise along specific evolutionary paths.

Nonetheless, several limitations remain. First, although variant effect annotation improved interpretability, the pipeline inherently depends on **SNP-based genomics**. Structural variants, mobile genetic elements, gene duplications, and plasmid-level dynamics, known drivers of AMR, remain underrepresented. Second, the current analysis is restricted to in silico predictions without experimental validation. Third, although these methods provided meaningful gene-level signals, they

cannot fully replace mechanistic experiments; attention-based or attribution-based saliency is still correlational rather than causal. Finally, the approach does not yet incorporate temporal evolutionary signals beyond the PCA projection; resistance emergence is inherently dynamic, and richer longitudinal models could better capture adaptive trajectories.

B. Future Work

Future extensions of this work can significantly enhance both the biological depth and predictive capability of the AMR modeling framework. One important direction is the incorporation of structural variants and mobile genetic elements. While the current pipeline focuses on SNPs, many clinically relevant resistance mechanisms arise from **plasmid-borne genes, integrons, transposable elements**, and larger structural changes such as copy-number variation or gene duplications. Incorporating these signals, potentially through hybrid short and long read sequencing, would enable a more comprehensive representation of AMR determinants.

Another promising avenue is the **integration of pan-genome and gene presence/absence features**. Since the accessory genome often harbors key resistance genes, expanding the feature space beyond point mutations to include gene-level variability can provide a fuller view of the genomic landscape underlying resistance. Combining SNP effects with pan-genomic signatures would allow the model to simultaneously capture both fine-grained nucleotide changes and broader genomic architecture.

A third area of expansion is the incorporation of temporal evolutionary modeling. Although PCA-based analyses revealed distinct evolutionary trajectories, the current approach does not yet model the dynamics of adaptation explicitly. Temporal models, such as **phylogeny-aware predictors, recurrent neural networks, or Bayesian evolutionary frameworks**, would allow the prediction of future resistance states, the detection of early adaptive signals, and the inference of evolutionary paths under antibiotic pressure.

Finally, opportunities exist to optimize the pipeline for translational use and to broaden its multi-modal capabilities. Clinically viable AMR prediction requires faster variant calling, efficient annotation, and models that incorporate uncertainty quantification for decision support. Additionally, integrating other biological data types, such as **gene expression profiles** during antibiotic exposure, **proteomics** of efflux systems, or clinical metadata, could substantially improve both predictive performance and interpretability. Together, these directions lay a clear path toward a more comprehensive, mechanistically grounded, and clinically deployable AMR evolutionary modeling framework.

REFERENCES

- [1] G. Muteeb, M. T. Rehman, M. Shahwan, and M. Aatif, "Origin of antibiotics and antibiotic resistance, and their impacts on drug development: A narrative review," en, *Pharmaceuticals (Basel)*, vol. 16, no. 11, p. 1615, Nov. 2023.
- [2] M. I. Hutchings, A. W. Truman, and B. Wilkinson, "Antibiotics: Past, present and future," en, *Curr. Opin. Microbiol.*, vol. 51, pp. 72–80, Oct. 2019.
- [3] S. Dhingra et al., "Microbial resistance movements: An overview of global public health threats posed by antimicrobial resistance, and how best to counter," en, *Front. Public Health*, vol. 8, p. 535 668, Nov. 2020.
- [4] M. A. Salam et al., "Antimicrobial resistance: A growing serious threat for global public health," en, *Health-care (Basel)*, vol. 11, no. 13, p. 1946, Jul. 2023.
- [5] S. K. Ahmed et al., "Antimicrobial resistance: Impacts, challenges, and future prospects," *Journal of Medicine, Surgery, and Public Health*, vol. 2, p. 100 081, 2024, ISSN: 2949-916X. DOI: <https://doi.org/10.1016/j.glmedi.2024.100081> [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2949916X24000343>
- [6] N. R. Naylor et al., "Estimating the burden of antimicrobial resistance: A systematic literature review," en, *Antimicrob. Resist. Infect. Control*, vol. 7, no. 1, p. 58, Apr. 2018.
- [7] M. A. Abushaheen et al., "Antimicrobial resistance, mechanisms and its clinical significance," en, *Dis. Mon.*, vol. 66, no. 6, p. 100 971, Jun. 2020.
- [8] I. Gajic et al., "Antimicrobial susceptibility testing: A comprehensive review of currently used methods," en, *Antibiotics (Basel)*, vol. 11, no. 4, p. 427, Mar. 2022.
- [9] M. Boolchandani, A. W. D'Souza, and G. Dantas, "Sequencing-based methods and resources to study antimicrobial resistance," en, *Nat. Rev. Genet.*, vol. 20, no. 6, pp. 356–370, Jun. 2019.
- [10] L. Barnes V, D. M. Heithoff, S. P. Mahan, J. K. House, and M. J. Mahan, "Antimicrobial susceptibility testing to evaluate minimum inhibitory concentration values of clinically relevant antibiotics," en, *STAR Protoc.*, vol. 4, no. 3, p. 102 512, Aug. 2023.
- [11] B. P. Alcock et al., "CARD 2023: Expanded curation, support for machine learning, and resistome prediction at the comprehensive antibiotic resistance database," en, *Nucleic Acids Res.*, vol. 51, no. D1, pp. D690–D699, Jan. 2023.
- [12] J. I. Kim et al., "Machine learning for antimicrobial resistance prediction: Current practice, limitations, and clinical perspective," en, *Clin. Microbiol. Rev.*, vol. 35, no. 3, e0017921, Sep. 2022.
- [13] A. F. Florensa, R. S. Kaas, P. T. L. C. Clausen, D. Aytan-Aktug, and F. M. Aarestrup, "ResFinder - an open online resource for identification of antimicrobial resistance genes in next-generation sequencing data and prediction of phenotypes from genotypes," en, *Microb. Genom.*, vol. 8, no. 1, Jan. 2022.
- [14] G. Arango-Argoty, E. Garner, A. Pruden, L. S. Heath, P. Vikesland, and L. Zhang, "DeepARG: A deep learning approach for predicting antibiotic resistance genes from metagenomic data," en, *Microbiome*, vol. 6, no. 1, Dec. 2018.

- [15] M. Feldgarden et al., "AMRFinderPlus and the reference gene catalog facilitate examination of the genomic links among antimicrobial resistance, stress response, and virulence," en, *Sci. Rep.*, vol. 11, no. 1, p. 12728, Jun. 2021.
- [16] L. Chindelevitch et al., *Applying data technologies to combat amr: Current status, challenges, and opportunities on the way forward*, 2022. arXiv: 2208.04683 [cs.CY]. [Online]. Available: <https://arxiv.org/abs/2208.04683>
- [17] D. Sun, K. Jeannot, Y. Xiao, and C. W. Knapp, "Editorial: Horizontal gene transfer mediated bacterial antibiotic resistance," *Frontiers in Microbiology*, vol. Volume 10 - 2019, 2019, ISSN: 1664-302X. DOI: 10.3389/fmicb.2019.01933 [Online]. Available: <https://www.frontiersin.org/journals/microbiology/articles/10.3389/fmicb.2019.01933>
- [18] D. Sun, X. Sun, Y. Hu, and Y. Yamaichi, "Editorial: Horizontal gene transfer mediated bacterial antibiotic resistance, volume II," en, *Front. Microbiol.*, vol. 14, p. 1221606, Jun. 2023.
- [19] F. Baquero et al., "Evolutionary pathways and trajectories in antibiotic resistance," en, *Clin. Microbiol. Rev.*, vol. 34, no. 4, e0005019, Dec. 2021.
- [20] B. A. Wilson, N. R. Garud, A. F. Feder, Z. J. Assaf, and P. S. Pennings, "The population genetics of drug resistance evolution in natural populations of viral, bacterial and eukaryotic pathogens," en, *Mol. Ecol.*, vol. 25, no. 1, pp. 42–66, Jan. 2016.
- [21] C. M. Hasan, D. Dutta, and A. N. T. Nguyen, "Revisiting antibiotic resistance: Mechanistic foundations to evolutionary outlook," en, *Antibiotics (Basel)*, vol. 11, no. 1, p. 40, Dec. 2021.
- [22] A. S. Mustafa, "Whole genome sequencing: Applications in clinical bacteriology," en, *Med. Princ. Pract.*, pp. 1–13, Feb. 2024.
- [23] S. K. Gupta et al., "ARG-ANNOT, a new bioinformatic tool to discover antibiotic resistance genes in bacterial genomes," en, *Antimicrob. Agents Chemother.*, vol. 58, no. 1, pp. 212–220, 2014.
- [24] K. Hu et al., "Assessing computational predictions of antimicrobial resistance phenotypes from microbial genomes," en, *Brief. Bioinform.*, vol. 25, no. 3, Mar. 2024.
- [25] S. Valavarasu, Y. Sangu, and T. Mahapatra, "Prediction of antibiotic resistance from antibiotic susceptibility testing results from surveillance data using machine learning," en, *Sci. Rep.*, vol. 15, no. 1, p. 30509, Aug. 2025.
- [26] Z. Liu et al., "Evaluation of machine learning models for predicting antimicrobial resistance of actinobacillus pleuropneumoniae from whole genome sequences," en, *Front. Microbiol.*, vol. 11, p. 48, Feb. 2020.
- [27] S. M. S. Rayesha, W. A. Banu, and A. Rahman, "Antibiotic genomic resistance prediction using deep learning models," en, *Int. J. Bioinform. Res. Appl.*, vol. 21, no. 2, pp. 121–136, 2025.
- [28] R. Preethi, R. Bharati, and S. Priya, "Predicting antibiotic resistance from genomic sequences using a hybrid cnn-rnn model: A comprehensive approach," in *2024 Third International Conference on Artificial Intelligence, Computational Electronics and Communication System (AICECS)*, 2024, pp. 1–6. DOI: 10.1109/AICECS63354.2024.10957126
- [29] W. P. M. Rowe and M. D. Winn, "Indexed variation graphs for efficient and accurate resistome profiling," en, *Bioinformatics*, vol. 34, no. 21, pp. 3601–3608, Nov. 2018.
- [30] S. M. Lakin et al., "Hierarchical hidden markov models enable accurate and diverse detection of antimicrobial resistance sequences," en, *Commun. Biol.*, vol. 2, no. 1, p. 294, Aug. 2019.
- [31] Y. Ren et al., "Prediction of antimicrobial resistance based on whole-genome sequencing and machine learning," en, *Bioinformatics*, vol. 38, no. 2, pp. 325–334, Jan. 2022.
- [32] Y. Ren et al., "Multi-label classification for multi-drug resistance prediction of escherichia coli," en, *Comput. Struct. Biotechnol. J.*, vol. 20, pp. 1264–1270, Mar. 2022.
- [33] G. Sevillya, O. Adato, and S. Snir, "Detecting horizontal gene transfer: A probabilistic approach," en, *BMC Genomics*, vol. 21, no. Suppl 1, p. 106, Mar. 2020.
- [34] C. J. H. von Wintersdorff et al., "Dissemination of antimicrobial resistance in microbial ecosystems through horizontal gene transfer," en, *Front. Microbiol.*, vol. 7, p. 173, Feb. 2016.
- [35] D. Sun, K. Jeannot, Y. Xiao, and C. W. Knapp, "Editorial: Horizontal gene transfer mediated bacterial antibiotic resistance," en, *Front. Microbiol.*, vol. 10, p. 1933, Aug. 2019.
- [36] M. J. Bottery, J. W. Pitchford, and V.-P. Friman, "Ecology and evolution of antimicrobial resistance in bacterial communities," en, *ISME J.*, vol. 15, no. 4, pp. 939–948, Apr. 2021.
- [37] A. Alexandre, A. Abbara, C. Fruet, C. Loverdo, and A.-F. Bitbol, *Bridging wright-fisher and moran models*, 2024. DOI: <https://doi.org/10.1016/j.jtbi.2024.112030> arXiv: 2407.12560 [q-bio.PE]. [Online]. Available: <https://arxiv.org/abs/2407.12560>
- [38] S. G. Das, S. O. Direito, B. Waclaw, R. J. Allen, and J. Krug, "Predictable properties of fitness landscapes induced by adaptational tradeoffs," en, *Elife*, vol. 9, no. e55155, May 2020.
- [39] N. Harmand, R. Gallet, R. Jabbour-Zahab, G. Martin, and T. Lenormand, "Fisher's geometrical model and the mutational patterns of antibiotic resistance across dose gradients," en, *Evolution*, vol. 71, no. 1, pp. 23–37, Jan. 2017.
- [40] R. Chait, J. Ruess, T. Bergmiller, G. Tkačik, and C. C. Guet, "Shaping bacterial population behavior through

- computer-interfaced control of individual cells,” en, *Nat. Commun.*, vol. 8, no. 1, p. 1535, Nov. 2017.
- [41] A. Ragalo and N. Pillay, “Evolving dynamic fitness measures for genetic programming,” en, *Expert Syst. Appl.*, vol. 109, pp. 162–187, Nov. 2018.
- [42] A. Yurtseven, S. Buyanova, A. A. Agrawal, O. O. Bochkareva, and O. V. Kalinina, “Machine learning and phylogenetic analysis allow for predicting antibiotic resistance in *m. tuberculosis*,” en, *BMC Microbiol.*, vol. 23, no. 1, p. 404, Dec. 2023.
- [43] Z. He, Y. Lin, R. Wei, C. Liu, and D. Jiang, “Repulsion and attraction in searching: A hybrid algorithm based on gravitational kernel and vital few for cancer driver gene prediction,” en, *Comput. Biol. Med.*, vol. 151, no. Pt A, p. 106 236, Dec. 2022.
- [44] M. E. Mowlaei and X. Shi, “FSF-GA: A feature selection framework for phenotype prediction using genetic algorithms,” en, *Genes (Basel)*, vol. 14, no. 5, May 2023.
- [45] P. J. Colin, D. J. Eleveld, and A. H. Thomson, “Genetic algorithms as a tool for dosing guideline optimization: Application to intermittent infusion dosing for vancomycin in adults,” en, *CPT Pharmacometrics Syst. Pharmacol.*, vol. 9, no. 5, pp. 294–302, May 2020.
- [46] A. Yurtseven, S. Buyanova, A. A. Agrawal, O. O. Bochkareva, and O. V. Kalinina, “Machine learning and phylogenetic analysis allow for predicting antibiotic resistance in *m. tuberculosis*,” en, *BMC Microbiol.*, vol. 23, no. 1, p. 404, Dec. 2023.
- [47] I. L. Brito, “Examining horizontal gene transfer in microbial communities,” en, *Nat. Rev. Microbiol.*, vol. 19, no. 7, pp. 442–453, Jul. 2021.
- [48] X. Wang et al., “Inter-plasmid transfer of antibiotic resistance genes accelerates antibiotic resistance in bacterial pathogens,” en, *ISME J.*, vol. 18, no. 1, Jan. 2024.
- [49] C. O. Vrancianu, L. I. Popa, C. Bleotu, and M. C. Chifiriuc, “Targeting plasmids to limit acquisition and transmission of antimicrobial resistance,” en, *Front. Microbiol.*, vol. 11, p. 761, May 2020.
- [50] T. Brown, X. Didelot, D. J. Wilson, and N. D. Maio, “SimBac: Simulation of whole bacterial genomes with homologous recombination,” en, *Microb. Genom.*, vol. 2, no. 1, Jan. 2016.
- [51] X. Didelot and D. Falush, “Inference of bacterial microevolution using multilocus sequence data,” en, *Genetics*, vol. 175, no. 3, pp. 1251–1266, Mar. 2007.
- [52] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton, “Adaptive mixtures of local experts,” en, *Neural Comput.*, vol. 3, no. 1, pp. 79–87, 1991.
- [53] K. Minoura et al., “A mixture-of-experts deep generative model for integrated single-cell multi-omics data,” *Bioinformatics*, vol. 37, no. Supplement 1, pp. i116–i124, 2021. DOI: 10.1093/bioinformatics/btab265
- [54] A. Spooner et al., “Benchmarking ensemble machine learning algorithms for multi-class, multi-omics data integration in clinical outcome prediction,” en, *Brief. Bioinform.*, vol. 26, no. 2, Mar. 2025.
- [55] N. Sun et al., “Mixture of experts enable efficient and effective protein understanding and design,” *bioRxiv preprint*, 2024. DOI: 10.1101/2024.11.29.625425
- [56] M. Guan, L. Zhao, and S. S.-T. Yau, “Classification of protein sequences by a novel alignment-free method on bacterial and virus families,” en, *Genes (Basel)*, vol. 13, no. 10, p. 1744, Sep. 2022.
- [57] K. O. Stanley and R. Miiikkulainen, “Evolving neural networks through augmenting topologies,” en, *Evol. Comput.*, vol. 10, no. 2, pp. 99–127, 2002.
- [58] R. Miiikkulainen et al., *Evolving deep neural networks*, 2017. arXiv: 1703.00548 [cs.NE]. [Online]. Available: <https://arxiv.org/abs/1703.00548>
- [59] S. Khadka and K. Tumer, *Evolution-guided policy gradient in reinforcement learning*, 2018. arXiv: 1805.07917 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/1805.07917>
- [60] J. Riou et al., “Projecting the development of antimicrobial resistance in *neisseria gonorrhoeae* from antimicrobial surveillance data: A mathematical modelling study,” en, *BMC Infect. Dis.*, vol. 23, no. 1, p. 252, Apr. 2023.
- [61] R. Dong, Q. Cao, and C. Song, “Painting peptides with antimicrobial potency through deep reinforcement learning,” en, *Adv. Sci. (Weinh.)*, no. e06332, e06332, Sep. 2025.
- [62] M. Feldgarden, V. Brover, B. Fedorov, D. H. Haft, A. B. Prasad, and W. Klimke, “Curation of the AM-RFinderPlus databases: Applications, functionality and impact,” en, *Microb. Genom.*, vol. 8, no. 6, Jun. 2022.
- [63] S. Andrews, “Fastqc: A quality control tool for high throughput sequence data,” 2010. [Online]. Available: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>
- [64] H. Li, *Aligning sequence reads, clone sequences and assembly contigs with bwa-mem*, 2013. arXiv: 1303.3997 [q-bio.GN].
- [65] H. Li and R. Durbin, “Fast and accurate short read alignment with Burrows-Wheeler transform,” en, *Bioinformatics*, vol. 25, no. 14, pp. 1754–1760, Jul. 2009.
- [66] H. Li et al., “The sequence Alignment/Map format and SAMtools,” en, *Bioinformatics*, vol. 25, no. 16, pp. 2078–2079, Aug. 2009.
- [67] Q. Zhao, “A study on optimizing MarkDuplicate in genome sequencing pipeline,” in *Proceedings of the 2018 5th International Conference on Bioinformatics Research and Applications*, Hong Kong Hong Kong: ACM, Dec. 2018.
- [68] G. A. Van der Auwera et al., “From FastQ data to high-confidence variant calls: The genome analysis toolkit best practices pipeline,” in *Current Protocols in Bioinformatics*. Hoboken, NJ, USA: John Wiley & Sons, Inc., Oct. 2013, pp. 11.10.1–11.10.33.

- [69] P. Danecek et al., “The variant call format and VCFtools,” en, *Bioinformatics*, vol. 27, no. 15, pp. 2156–2158, Aug. 2011.
- [70] H. E. Hanson and A. L. Liebl, “The mutagenic consequences of DNA methylation within and across generations,” en, *Epigenomes*, vol. 6, no. 4, p. 33, Oct. 2022.
- [71] A. van Belkum, S. Scherer, L. van Alphen, and H. Verbrugh, “Short-sequence DNA repeats in prokaryotic genomes,” en, *Microbiol. Mol. Biol. Rev.*, vol. 62, no. 2, pp. 275–293, Jun. 1998.
- [72] S. B. Needleman and C. D. Wunsch, “A general method applicable to the search for similarities in the amino acid sequence of two proteins,” en, *J. Mol. Biol.*, vol. 48, no. 3, pp. 443–453, Mar. 1970.