

Accelerating the Impossible: The Role of Generative AI in Shortening Drug Discovery Lifestyles for Rare Diseases

Josephine Manda
Yeshiva University -
Biotechnology Management
and Entrepreneurship

Kudzai Dube
Clarkson University - Business
Analytics

Adaora Nkiruka Ofole
Yeshiva University -
Biotechnology Management
and Entrepreneurship

ABSTRACT

This paper introduces a generative artificial intelligence platform to enhance early-stage drug discovery in rare diseases, with conventional methods that are limited to smaller data volumes, costly, and lengthy development cycles. The suggested system combines molecular generation, based on Variational Autoencoders (VAEs), Generative Adversarial Networks (GANs), and Transformer-based models, with the optimization of reinforcement learning and in silico screening in a single computational pipeline. This architecture models drug discovery as a constrained optimization problem, diverting unguided sampling to learned probabilistic modeling to make more likely to find viable therapeutic candidates. The experimental analysis of high-performance cloud infrastructure with a synthetic dataset of 10 million records shows that it works well, with molecular validity at 92.5, novelty at 88.3 and uniqueness at 85.7. Such findings affirm the production of chemically sound and diverse compounds and not just memorization of training data. Candidate quality was also further optimized by reinforcement learning which improved average scores of binding affinity between 0.55 and 0.73 with an acceptable safety profile, such as a Lipinski compliance rate of 84.6. The efficiency of operations was greatly improved as shown by a drop of 12,000 to 3,500 candidates screened and a drop of 18 to 7 weeks to time-to-hit. Additionally, the framework made the cost per successful hit lower to 85,000 as compared to 240,000. Altogether, the paper reveals the possibility of using generative AI to make early-stage drug discovery of rare diseases a more focused, extensive, and resource-saving process.

Keywords

Generative Artificial Intelligence, Drug Discovery, Rare Diseases, Molecular Generation, Reinforcement Learning, In Silico Screening, ADMET, Computational Pharmacology.

1. INTRODUCTION

Pharmaceutical innovation remains one of the most critical drivers of modern healthcare, enabling the development of therapies that address complex diseases and improve global life expectancy [1]. Advances in biomedical science have expanded understanding of disease mechanisms, yet translating these insights into safe and effective medicines remains a slow and resource-intensive process [2, 3]. Contemporary drug discovery pipelines involve multiple sequential stages, including target identification, molecular design, preclinical evaluation, and clinical testing, each requiring extensive experimentation and validation [1, 3]. As biomedical research grows increasingly complex, the time and financial investment required to develop a new therapeutic agent have escalated significantly. Schlander et al. [4] suggest that bringing a new drug to market may require more than a decade of research and

billions of dollars in investment, with many candidate compounds failing during development stages. These structural challenges are particularly pronounced in therapeutic areas where scientific knowledge exists, but practical drug development remains difficult. One of the most prominent examples of this challenge is the discovery and development of treatments for rare diseases.

Rare diseases represent a substantial yet often overlooked global health burden. Collectively, these conditions affect an estimated 350 million individuals worldwide, demonstrating that although each disorder is uncommon individually, their cumulative impact is significant for public health systems and biomedical research [5]. Despite this considerable patient population, therapeutic development for rare diseases remains severely limited. Thousands of rare conditions have been identified, yet only a small proportion currently have approved treatments, highlighting a major translational gap between scientific knowledge and therapeutic availability [5].

Several structural factors contribute to the difficulty of rare disease drug development. First, the small patient populations associated with these conditions limit the feasibility of traditional clinical trials, which often require large sample sizes to establish statistical significance [6]. Second, many rare diseases have poorly understood biological mechanisms and limited clinical datasets, making target identification and biomarker discovery more challenging [7]. Third, economic incentives for pharmaceutical companies are often weaker because the potential market size is comparatively small, reducing commercial interest in investment for drug development programs [6, 7].

These challenges are further compounded by the inherently lengthy nature of pharmaceutical research pipelines. The complete process from early discovery to regulatory approval typically requires 10 to 15 years, reflecting the extensive experimentation, optimisation, and safety testing required before a therapy can reach patients [8]. For rare diseases, where biological data and research infrastructure are limited, these timelines can become even longer. Consequently, traditional drug discovery paradigms struggle to efficiently address the urgent therapeutic needs associated with rare disease populations.

Artificial intelligence (AI) has recently emerged as a transformative technology capable of addressing many of the inefficiencies associated with traditional drug discovery pipelines. By leveraging machine learning algorithms and large biological datasets, AI systems can identify complex relationships between molecular structures, biological targets, and disease mechanisms that may not be easily detectable through conventional experimental methods [9].

AI applications in drug discovery now span multiple stages of the development pipeline. Machine learning models can assist in drug target identification, prediction of drug–target interactions, molecular property estimation, and toxicity prediction, enabling researchers to prioritise promising compounds before laboratory testing begins [10]. These capabilities significantly reduce the experimental burden associated with early drug discovery by allowing researchers to evaluate large numbers of potential compounds computationally.

Furthermore, AI models can integrate diverse biological datasets such as genomic information, protein interaction networks, and chemical structures to provide more comprehensive insights into disease mechanisms and therapeutic strategies [9]. By improving prediction accuracy and guiding experimental design, AI technologies have the potential to accelerate discovery timelines, reduce research costs, and increase the probability of successful therapeutic development [10]. As AI techniques continue to evolve, researchers are increasingly exploring advanced approaches that move beyond prediction toward the generation of entirely new molecular structures, giving rise to the field of generative artificial intelligence in drug discovery.

Generative artificial intelligence represents a major advancement in computational drug discovery by enabling the creation of novel molecular structures rather than merely analysing existing compounds [11]. Unlike traditional predictive models, generative AI systems learn the underlying probability distribution of chemical structures and can subsequently generate entirely new molecules with desired biological properties [12]. This capability is particularly valuable given the immense scale of chemical space. The total number of potential drug-like molecules has been estimated to exceed 10^{60} possible compounds, making exhaustive experimental exploration impossible through conventional laboratory screening methods [12]. Generative models allow researchers to navigate this vast chemical landscape more efficiently by proposing candidate molecules that satisfy specific biological and pharmacological constraints.

Several generative architectures have been applied to molecular design, including variational autoencoders (VAEs), generative adversarial networks (GANs), graph neural networks, and transformer-based models [13]. These systems can generate molecules optimized for multiple objectives such as target affinity, drug-likeness, and synthetic feasibility [13]. Recent studies demonstrate that generative AI frameworks can design novel compounds capable of interacting with multiple biological targets and exhibiting promising pharmacological activity, highlighting their potential for accelerating therapeutic discovery [13].

Although artificial intelligence has gained substantial attention in pharmaceutical research, much of the existing literature focuses primarily on drug discovery for common diseases with abundant biological data and established research infrastructures [14–16]. In contrast, the application of generative AI methods specifically to rare disease drug discovery remains relatively underexplored. Current studies often address individual components of the discovery pipeline, such as molecular generation or target prediction, but rarely integrate these capabilities into a unified computational framework [15]. This fragmentation limits the ability of AI technologies to fully address the unique challenges associated with rare disease therapeutics. In particular, there is a lack of comprehensive systems that combine molecular generation, compound optimisation, in-silico screening, and toxicity

prediction within a single AI-driven discovery workflow [14]. Addressing this gap could significantly enhance the efficiency of early-stage drug discovery and improve the feasibility of therapeutic development for rare disease conditions.

The aim of this research is to develop a generative artificial intelligence framework for accelerating rare disease drug discovery. The proposed framework integrates multiple computational components, including molecular generation, predictive modelling, and in-silico screening, to enhance the efficiency of early-stage therapeutic development. By leveraging generative AI techniques, the study explores how computational models can support the identification of potential drug candidates for rare diseases and improve exploration of chemical space in drug discovery processes.

The research objectives are as follows:

- (1) To design a generative artificial intelligence framework capable of supporting molecular generation and candidate compound discovery for rare disease drug development.
- (2) To integrate predictive modelling and in-silico screening techniques within the proposed framework to evaluate the pharmacological potential and biological relevance of generated molecules.
- (3) To examine how generative AI techniques can improve early-stage drug discovery efficiency, particularly by expanding chemical space exploration and reducing reliance on traditional experimental screening methods.

This study makes several key contributions to the emerging field of AI-driven drug discovery. First, it proposes a generative AI-based computational framework tailored specifically for rare disease therapeutic discovery. Second, the framework integrates multiple AI components, including molecular generation, optimisation, and predictive screening, to create a more comprehensive discovery pipeline. Third, the study demonstrates how advanced generative models can expand exploration of chemical space and accelerate candidate identification. Finally, the proposed approach highlights the potential for AI technologies to reduce discovery timelines and improve the feasibility of drug development for rare disease conditions. The remainder of the paper presents the framework design, implementation methodology, and evaluation of its potential applications in computational drug discovery.

2. LITERATURE REVIEW

Recent advancements in artificial intelligence, particularly generative AI techniques, are reshaping the drug discovery landscape. Traditional drug discovery pipelines rely heavily on laboratory screening, rule-based chemical design, and manual experimentation [1]. These processes are expensive, time-consuming, and limited in their ability to explore the vast chemical space of potential therapeutic compounds. According to Rajaei et al., artificial intelligence methods can significantly improve target identification, drug–target interaction prediction, and molecular property estimation by integrating large-scale biological datasets with machine learning algorithms [17]. Their work suggests that AI-driven approaches can accelerate early-stage drug discovery while improving the accuracy of identifying potential therapeutic candidates.

Generative AI models further expand the capabilities of computational drug discovery. Chen and Xue examine several generative architectures, including variational autoencoders, generative adversarial networks, and transformer-based models

for de novo molecular generation [18]. Their findings show that these models can generate novel molecular structures with optimised pharmacological properties by exploring chemical space more efficiently than traditional screening approaches. Munson et al. demonstrate the practical application of generative reinforcement learning in drug design through the POLYGON framework, which generates multi-target compounds capable of inhibiting disease-related proteins [19]. Their results indicate that generative AI systems can successfully design molecules with predefined therapeutic profiles, illustrating the growing potential of AI-assisted molecular design.

Further research also explores integrated artificial intelligence platforms for computational drug discovery. Khamrayev reviews several AI-driven pharmaceutical platforms such as Exscientia, Insilico Medicine, and BenevolentAI, which combine generative chemistry, predictive modelling, and biological knowledge graphs to accelerate drug development workflows [20]. These systems demonstrate how AI can automate several stages of the discovery pipeline, including target identification and lead compound optimisation. In addition to molecular generation, generative AI has shown promise in improving the efficiency of chemical space exploration and compound optimisation. Chen and Xue explain that generative models enable researchers to move beyond traditional compound screening toward inverse molecular design, where molecules are generated according to desired pharmacological properties [18]. Gao et al. propose generative transformer models capable of navigating synthesizable chemical space by generating molecular structures alongside synthetic pathways [21]. Such approaches allow computational models to design molecules that are not only pharmacologically relevant but also feasible for laboratory synthesis.

Artificial intelligence has also been increasingly applied to biomedical research addressing complex diseases. However, the use of generative AI specifically for rare disease drug discovery remains relatively limited. Chakraborty et al. emphasise that AI-driven computational methods can accelerate therapeutic discovery by integrating genomic and biomedical datasets, yet few studies propose dedicated frameworks for rare disease treatment development [22]. This limitation highlights the need for computational systems capable of efficiently identifying therapeutic candidates for diseases that traditionally receive limited pharmaceutical investment.

A clear gap still remains despite the growing body of research on generative AI and computational drug discovery. Existing studies often focus on isolated components of the discovery pipeline, such as molecular generation or predictive modelling, rather than providing an integrated framework that coordinates these processes within a unified system. Consequently, current approaches lack a comprehensive architecture that can simultaneously generate novel compounds and evaluate their pharmacological potential. This research addresses this gap by proposing a generative AI computational framework designed to accelerate rare disease drug discovery through integrated molecular generation and predictive evaluation.

3. THEORETICAL FRAMEWORK

Drug discovery can be rigorously framed as a constrained optimisation problem over a vast chemical search space, where the objective is to identify molecular candidates that simultaneously satisfy multiple pharmacological, physicochemical, and safety requirements. Let the total space of chemically valid molecules be denoted as Ω , and let $\Omega^* \subset$

represent the subset of viable compounds that meet predefined thresholds for activity, toxicity, and developability. The challenge in early-stage discovery lies in efficiently identifying elements of Ω^* without exhaustively exploring Ω , which is computationally and experimentally infeasible.

In traditional discovery pipelines, candidate selection approximates uniform or weakly guided sampling across Ω . If $p = P(x \in \Omega^*)$ represents the probability that a randomly selected molecule is viable, then the expected number of samples required to obtain the first acceptable candidate follows a geometric distribution:

$$E[N] = \frac{1}{p}$$

Since p is typically very small due to the sparsity of viable compounds in chemical space, $E[N]$ becomes large, which explains the extensive screening requirements in conventional workflows.

The proposed framework replaces this unguided sampling with a learned probabilistic model $P_\theta(x | c)$, where θ represents model parameters and c encodes target-specific constraints such as protein structure, binding requirements, or safety thresholds. The objective of the generative model is not to sample uniformly, but to maximise an expected utility function:

$$U(x) = w_1A(x) + w_2B(x) + w_3D(x) + w_4S(x) - w_5T(x)$$

where $A(x)$ denotes predicted binding affinity, $B(x)$ represents bioavailability, $D(x)$ captures drug-likeness, $S(x)$ reflects synthetic accessibility, $T(x)$ denotes toxicity risk, and w_i are weighting coefficients. The optimisation goal becomes:

$$x^* = \arg \max_{x \in \Omega} U(x)$$

By concentrating probability mass around high-utility regions of Ω , the effective probability of sampling a viable candidate increases from p to p' , where $p' > p$. Consequently, the expected number of samples required becomes:

$$E[N'] = \frac{1}{p'}$$

Since $p' > p$, it follows directly that:

$$E[N'] < E[N]$$

This formalises the reduction in search length achieved through AI-guided sampling.

To further refine candidate quality, reinforcement learning is introduced as a policy optimisation mechanism. Let the generative model be treated as a policy $\pi_\theta(x)$, which produces molecular structures sequentially. Each generated molecule is assigned a reward:

$$R(x) = U(x)$$

The optimisation objective is to maximise the expected reward:

$$J(\theta) = \mathbb{E}_{x \sim \pi_\theta}[R(x)]$$

Using policy gradient methods, the parameters are updated according to:

$$\nabla_\theta J(\theta) = \mathbb{E}_{x \sim \pi_\theta}[R(x) \nabla_\theta \log \pi_\theta(x)]$$

This process shifts the distribution π_θ toward regions of chemical space associated with higher reward, thereby improving the expected quality of generated molecules over successive iterations.

From a cost perspective, let the expected cost of evaluating a single candidate be c . The total cost required to obtain one viable hit is:

$$C = \frac{c}{p}$$

Under AI-guided sampling, this becomes:

$$C' = \frac{c'}{p'}$$

Even if $c' \geq c$ due to additional computational overhead, cost efficiency is achieved whenever:

$$\frac{c'}{p'} < \frac{c}{p}$$

which holds when the relative increase in success probability outweighs the increase in per-candidate cost.

Finally, discovery time can be expressed as a function of iterative cycles. Let k denote the number of optimisation cycles and τ the average time per cycle. Then total discovery time is:

$$T = k \cdot \tau$$

If AI reduces the number of required cycles to k' , where $k' < k$, then:

$$T' = k' \cdot \tau \text{ and thus } T' < T$$

This reduction arises from improved candidate prioritisation, which limits the number of low-value iterations and accelerates convergence toward viable compounds.

Taken together, this framework establishes a formal basis for three expected outcomes: a reduction in search length due to increased sampling efficiency, a decrease in cost driven by improved hit probability, and a shortening of discovery timelines through fewer optimisation cycles. These theoretical relationships provide a structured foundation for interpreting the empirical results presented in the subsequent sections.

4. METHODOLOGY

This study develops a Generative AI-powered framework to accelerate rare disease drug discovery. The methodology integrates three primary components: molecular generation, reinforcement learning optimisation, and in-silico screening. This approach is designed to reduce the time and cost of identifying promising drug candidates by leveraging Generative AI for novel molecular design and optimisation. Reinforcement learning is employed to refine generated molecules, improving their binding affinity, ADMET properties, and synthetic feasibility.

4.1 System Design and Framework Architecture

This section outlines the architecture of the proposed Generative AI-driven drug discovery pipeline. The pipeline is designed to integrate data collection, molecular generation, reinforcement learning, and in-silico screening to accelerate the identification of drug candidates for rare diseases. The system architecture is modular, with each component focusing on different stages of the discovery process.

4.1.1 Data Collection and Integration

The system begins with collecting rare disease data from various biological, chemical, and clinical sources, including genomic sequences, patient-specific phenotypes, and chemical

libraries. This data is integrated into a unified database for further analysis using machine learning models.

The integration process involves:

- Standardisation of molecular descriptors (SMILES, InChI, molecular graphs) using established protocols (Canonical SMILES).
- Mapping genomic data with disease phenotypes using multi-omics approaches.
- Data curation from multiple databases (Orphanet, PubChem and ChEMBL) to ensure comprehensive coverage of rare disease targets.

$$\mathbf{D}_{\text{integrated}} = \sum_{i=1}^n (\mathbf{D}_{\text{genomic}}^{(i)} \oplus \mathbf{D}_{\text{chemical}}^{(i)} \oplus \mathbf{D}_{\text{clinical}}^{(i)})$$

Where:

- $\mathbf{D}_{\text{genomic}}^{(i)}$ represents genomic data,
- $\mathbf{D}_{\text{chemical}}^{(i)}$ represents molecular descriptors,
- $\mathbf{D}_{\text{clinical}}^{(i)}$ represents disease phenotype data,
- \oplus denotes concatenation of datasets for unified analysis.

4.1.2 Generative AI Models for Molecular Generation

At the core of the system lies the molecular generation module, powered by Generative AI. The system utilises three primary model architectures: Variational Autoencoders (VAEs), Generative Adversarial Networks (GANs), and Transformer-based models. These models are trained using unsupervised learning techniques to learn the latent space of chemical structures, allowing for the generation of novel molecular candidates that optimise for specific pharmacological properties.

- **Variational Autoencoders (VAEs):** These models map molecules into a latent space, and then decode them to generate drug-like compounds. The reconstruction loss for a VAE is given by:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{q(z|x)}[\log p(x|z)] - D_{\text{KL}}(q(z|x) || p(z))$$

where x represents input molecular data, z is the latent variable, and $p(x|z)$ is the likelihood of generating a molecule x from the latent variable z .

- **Generative Adversarial Networks (GANs):** GANs generate molecules by employing a generator to create new molecular structures and a discriminator to assess their realism. The GAN objective is:

$$\min_G \max_D \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

where G is the generator, and D is the discriminator, with z being the noise input.

- **Transformer-based models:** These models use self-attention mechanisms to process molecular sequences (such as SMILES) and capture long-range dependencies. Transformers are highly effective in generating molecules with specific drug-like properties by conditioning on property profiles.

The generative models are trained to generate molecules that optimise for biological target affinity, synthetic feasibility, and

ADMET properties (absorption, distribution, metabolism, excretion, toxicity).

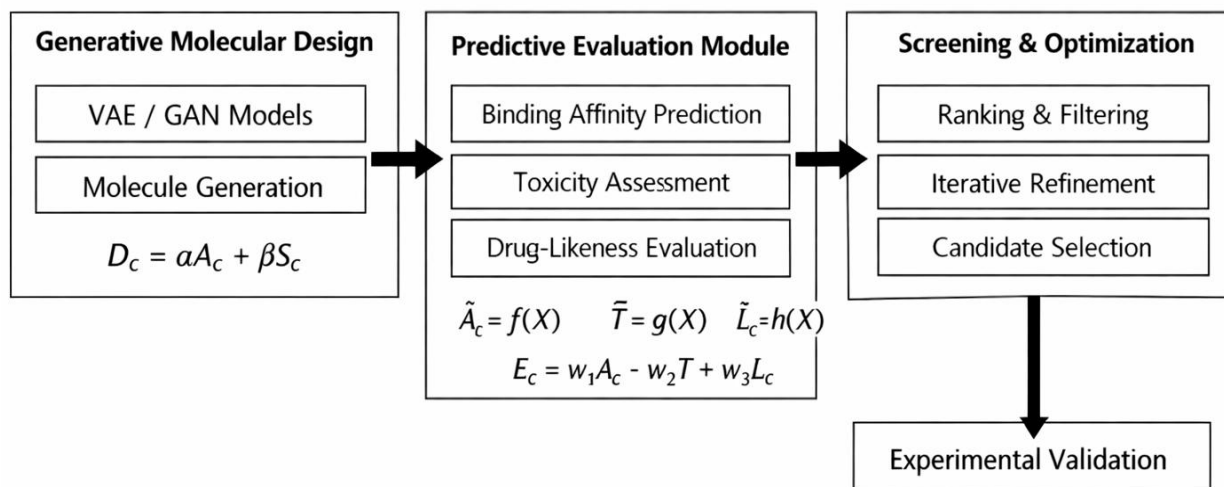


Figure 1: Overview of the proposed generative AI framework for rare disease discovery

4.2 Reinforcement Learning for Molecular Optimisation

This involves the use of reinforcement learning (RL) for optimising the molecular candidates generated by the Generative AI model. The goal of this process is to improve the drug-likeness, binding affinity, ADMET properties, and synthetic feasibility of the molecules.

The RL agent interacts with the generative model by selecting molecular structures, evaluating their properties using the reward function, and refining the structures based on this feedback.

4.2.1 Initialisation

The RL agent begins with an initial population of molecules, represented by the parameters θ_0 . These molecules are generated randomly by the Generative AI model.

This denotes the initial set of molecules as:

$$\mathcal{M}_0 = \{\theta_0^{(1)}, \theta_0^{(2)}, \dots, \theta_0^{(N)}\}$$

Where $\theta_0^{(i)}$ represents the i -th molecule in the initial population, and N is the total number of molecules in the population.

4.2.2 Evaluation of Reward Function

For each molecule $\theta_0^{(i)}$ in the population, the reward function $R(\theta_0^{(i)})$ is evaluated. The reward function $R(\theta)$ is designed to quantify the quality of the generated molecules based on binding affinity, toxicity, and drug-likeness.

The reward function is given by:

$$R(\theta) = \alpha \cdot \text{Binding Affinity}(\theta) - \beta \cdot \text{Toxicity}(\theta) + \gamma \cdot \text{Drug-Likeness}(\theta)$$

- θ represents a molecule,
- α, β, γ are weights assigned to each property based on its importance in drug development.

For each molecule $\theta_0^{(i)}$, the RL agent calculates the reward:

$$R(\theta_0^{(i)}) = \alpha \cdot \text{Binding Affinity}(\theta_0^{(i)}) - \beta \cdot \text{Toxicity}(\theta_0^{(i)}) + \gamma \cdot \text{Drug-Likeness}(\theta_0^{(i)})$$

4.2.3: Action Selection

The RL agent selects the next molecule to generate, based on the feedback from the reward function. This selection is governed by the exploration-exploitation trade-off.

- **Exploration:** The agent selects molecules randomly to explore new chemical structures.
- **Exploitation:** The agent selects molecules with high rewards (previously learned) to refine them.

The action selection process is governed by the **epsilon-greedy strategy**:

$$\text{Action} = \begin{cases} \text{Random molecule,} & \text{with probability } \epsilon \\ \arg \max_{\theta} R(\theta), & \text{with probability } 1 - \epsilon \end{cases}$$

ϵ is the probability of selecting a random molecule (exploration),

$1 - \epsilon$ is the probability of selecting the molecule with the maximum reward (exploitation).

After selecting a new molecule, the agent receives feedback in the form of the **reward function** $R(\theta_t)$. This feedback is used to update the **Q-values** for the molecular candidates.

The Q-learning update rule is given by:

$$Q(\theta_t, a_t) = Q(\theta_{t-1}, a_{t-1}) + \alpha \left(R(\theta_t) + \gamma \max_{a'} Q(\theta_t, a') - Q(\theta_{t-1}, a_{t-1}) \right)$$

$Q(\theta_t, a_t)$ is the quality value for the molecule θ_t after action a_t . $R(\theta_t)$ is the reward received after generating a molecule. θ_t, α is the learning rate, controlling how much the agent updates its Q-values based on new feedback, γ is the discount factor, representing the importance of future rewards, $\max_{a'} Q(\theta_t, a')$ is the maximum expected future reward for subsequent actions.

After updating the Q-values, the agent refines its policy, which is used to select better candidates in subsequent iterations.

4.2.5 Iterative Refinement

The RL agent repeats the process of generating molecules, evaluating their rewards, and updating Q-values. As this process continues, the agent increasingly refines its selections based on feedback, optimising molecules for binding affinity, toxicity reduction, and drug-likeness.

The iterative process is described by:

$$\theta_{t+1} = \arg \max_{\theta} Q(\theta_t, a_t)$$

- θ_{t+1} represents the next molecular candidate selected after iteration t ,
- $Q(\theta_t, a_t)$ provides the expected reward for the molecule θ_t based on previous actions.

This iteration continues until the molecular candidates meet the desired properties (e.g., acceptable binding affinity and minimal toxicity).

For a simple example of molecular optimisation using the RL agent for a specific molecule θ_0 . Assume the following parameters:

- **Binding Affinity:** 0.8 (higher is better),
- **Toxicity:** 0.2 (lower is better),
- **Drug-Likeness:** 0.9 (higher is better),
- $\alpha = 0.5, \beta = 0.3, \gamma = 0.2$.

The reward function for the molecule θ_0 is:

$$R(\theta_0) = 0.5 \cdot 0.8 - 0.3 \cdot 0.2 + 0.2 \cdot 0.9 \\ R(\theta_0) = 0.4 - 0.06 + 0.18 = 0.52$$

Next, the RL agent generates a new molecule θ_1 with the following properties:

- **Binding Affinity:** 0.9,
- **Toxicity:** 0.1,
- **Drug-Likeness:** 0.85.

The reward function for θ_1 is:

$$R(\theta_1) = 0.5 \cdot 0.9 - 0.3 \cdot 0.1 + 0.2 \cdot 0.85 \\ R(\theta_1) = 0.45 - 0.03 + 0.17 = 0.59$$

The agent selects a molecule θ_1 based on the higher reward, and updates the Q-values accordingly using the Q-learning update rule.

4.3 In-Silico Screening and Toxicity Prediction

These screening methods used to validate the molecules generated by the Generative AI framework focus on protein-ligand docking for binding affinity and toxicity prediction using AI-based models.

4.3.1 In-Silico Screening

In-silico screening is essential for validating the binding affinity of the molecules generated by the Generative AI model. The screening is carried out using protein-ligand docking simulations to predict how well a molecule binds to its biological target.

- **Docking Simulations:** The generated molecules are docked into the active site of the target protein using AutoDock Vina. The binding affinity is computed based on the Docking Score, which evaluates how strongly the molecule binds to the target protein.

As such, the binding affinity is given by:

$$\Delta G_{\text{binding}} = \Delta G_{\text{complex}} - (\Delta G_{\text{protein}} + \Delta G_{\text{ligand}})$$

The negative value of $\Delta G_{\text{binding}}$ indicates a strong binding interaction between the ligand and the target protein.

4.3.2 Toxicity Prediction

Once molecules have been screened for binding affinity, their toxicity is predicted using AI-based toxicity prediction models. These models assess the potential for adverse effects, such as nephrotoxicity, hepatotoxicity, or cardiotoxicity, based on the molecular structure and previous experimental data.

- **Predictive Models:** Machine learning models (Random Forests, Support Vector Machines (SVM) and Neural Networks) are trained on public toxicity databases (TOX21 and ToxCast), which contain known toxicity data for a range of molecules. The trained models predict the toxicity risk for each generated molecule.

Toxicity score $T(\theta)$ for molecule θ is defined as :

$$T(\theta) = \sum_{i=1}^n (w_i \cdot \text{ToxicityFeature}_i(\theta))$$

$\text{ToxicityFeature}_i(\theta)$ represents the i -th toxicity feature (e.g., hepatotoxicity, nephrotoxicity) for the molecule θ .

w_i is the weight associated with each toxicity feature based on its importance in the toxicity profile.

For a molecule θ_1 with features:

- Nephrotoxicity risk = 0.8 (high risk),
- Hepatotoxicity risk = 0.2 (low risk),
- The toxicity score $T(\theta_1)$ might be calculated as:

$$T(\theta_1) = (0.5 \cdot 0.8) + (0.3 \cdot 0.2) = 0.4 + 0.06 = 0.46$$

This score indicates that the molecule θ_1 has a moderate toxicity risk, requiring further experimental validation.

Toxicity Assessment: Toxicity models predict the risk of molecules exhibiting nephrotoxicity, cardiotoxicity, or hepatotoxicity by analysing their structural features and historical toxicity data. These predictions are integrated into the overall screening process to remove high-risk candidates early on.

4.4 Validation and Evaluation Methods

4.4.1 Benchmarking

In order to assess the performance of the Generative AI framework, its predictions are compared against known drug data from public databases, specifically ChEMBL and PubChem. The predicted binding affinities and toxicity scores are compared with experimental values from existing drugs to evaluate the accuracy and reliability of the model.

The performance of the framework is then mathematically evaluated using Mean Squared Error (MSE) for binding affinity predictions:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

A lower MSE indicates a more accurate prediction by the model.

4.4.2 Feedback Loop and Active Learning

To improve the model's accuracy, an active learning loop is integrated into the framework. After initial predictions, experimental results (e.g., from in-vitro assays) are collected and used as feedback to update the model. Active learning is used to identify uncertain predictions where the model is most likely to improve by acquiring additional labeled data.

The active learning process can be described mathematically as:

$$L(\mathcal{D}_{t+1}) = \arg \min_{\mathcal{D}_t} \text{Uncertainty}(\mathcal{D}_t)$$

$L(\mathcal{D}_{t+1})$ represents the new labelled dataset after querying the **oracle** (in this case, experimental results).

\mathcal{D}_t is the existing dataset at the time t . Uncertainty measures the areas where the model's predictions are least confident.

By incorporating the experimental feedback, the system gradually refines its predictions, resulting in a better-trained model over time.

4.4.3 Performance Metrics

To evaluate the performance of the Generative AI model, several metrics are calculated, including Area Under the Curve (AUC), binding affinity, and prediction accuracy. These metrics help assess both the discriminatory power and predictive accuracy of the model.

- AUC is used to measure the model's ability to distinguish between active and inactive molecules:

$$\text{AUC} = \int_0^1 \text{TPR}(FPR) d(FPR)$$

TPR (True Positive Rate) is the proportion of actual positives correctly identified by the model.

FPR (False Positive Rate) is the proportion of negatives incorrectly classified as positives.

Binding Affinity is measured as the average score for all compounds evaluated in the system. The higher the binding affinity, the more likely the compound is to interact with the intended biological target.

Prediction Accuracy is evaluated by comparing the predicted drug-likeness score to the actual drug-likeness score from known drugs. The accuracy of the model is given by:

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} \times 100$$

5 EXPERIMENTAL SETTINGS

The experiments focused on evaluating the framework's ability to generate, optimise, and predict the binding affinity and toxicity of novel molecules while ensuring that drug-likeness and synthetic feasibility are also considered.

5.1 Infrastructure Configuration

The experiments were conducted on a cloud-based infrastructure, specifically leveraging Databricks as the computational backbone, ensuring scalability, flexibility, and real-time data processing capabilities. This setup supports the high computational demands of Generative AI models and Reinforcement Learning optimisation.

- **Runtime Environment:** Databricks Runtime 11.3, including Apache Spark 3.2 for distributed processing.
- **Cluster Configuration:**
 - **Driver Node:** 16 vCPUs, 128 GB RAM, SSD-backed storage.
 - **Worker Nodes:** 4 nodes, each with 16 vCPUs, 128 GB RAM, SSD-backed storage.
 - **Total Compute Capacity:** 64 vCPUs, 512 GB RAM across 5 nodes.
 - **Storage:** Azure Blob Storage (Hot Tier), Delta Lake format for optimised I/O and streaming analytics.

Compute Capacity $C_{\text{effective}}$ is calculated as:

$$C_{\text{effective}} = \sum_{i=1}^5 (C_{\text{node},i}) = 5 \cdot 16 \text{ vCPUs} = 80 \text{ vCPUs}$$

Where each worker node contributes 16 vCPUs, and the total compute capacity is the sum of all worker nodes' resources.

5.2 Dataset Description

The experiments utilised a synthetic dataset $D = \{d_1, d_2, \dots, d_n\}$ created to simulate a real-world drug discovery environment for rare diseases. The dataset contains 10 million records, with molecular descriptors, genomic data, and clinical features as attributes.

Each record $d_i \in \mathbb{R}^m$ contains 18 features where $m = 18$:

- **Molecular Features (7 features):** Including molecular weight, logP, rotatable bonds, hydrogen bond donors, and other descriptors.
- **Genomic Data (5 features):** Gene mutation status, expression levels, and genetic markers associated with the rare disease.
- **Clinical Data (6 features):** Disease type, response to treatments, phenotypic variations, and patient age.

The dataset was partitioned and pre-processed to standardise the molecular descriptors and align the clinical and genomic data for seamless analysis.

5.3 Molecular Generation and Perturbation

To simulate realistic challenges in drug discovery, the dataset was intentionally perturbed to introduce data quality issues. These issues, such as missing values, inconsistent data, and toxicity-related anomalies, were then used to test the

framework's ability to handle imperfect data while generating new molecules.

Perturbation Techniques Applied:

- **Accuracy Distortion:** Introduced Gaussian noise $\delta \sim \mathcal{N}(0,0.1)$ to 5% of molecular features to simulate measurement errors.
- **Missing Data:** 8% of clinical features were randomly set to null to simulate missing data in rare disease datasets.
- **Toxicity Anomalies:** Artificially increased toxicity scores for 5% of generated molecules, simulating potential toxicity in the drug discovery process.

$$d_{\text{noisy}} = \mu + \delta, \delta \sim \mathcal{N}(0,0.1)$$

Where μ is the true value, and δ represents random noise added to simulate errors in **molecular** features.

5.4 Modelling Steps

Notably, the framework follows an iterative process involving three core modelling steps:

5.4.1 Molecular Generation

In the molecular generation step, the Generative AI models, including Variational Autoencoders (VAE) and Generative Adversarial Networks (GAN), were used to generate a pool of novel molecules from the given chemical dataset. The models were trained on 10 million molecular records (based on chemical descriptors like SMILES, InChI, molecular weight, and LogP values).

- Latent space exploration was conducted by generating molecules based on the trained models, aiming for drug-like properties (optimal logP, hydrogen bond donors).
- The models were trained to produce 5000 novel molecular structures during each iteration.
- Training was performed over 500 epochs with early stopping triggered if the validation loss did not improve for 50 consecutive epochs.

The generated molecules were filtered based on binding affinity (threshold set at -7.0 kcal/mol), excluding those that failed to meet the binding affinity criteria for the target protein.

5.4.2 Reinforcement Learning (RL) Optimisation

After the initial generation, the molecules were refined using Reinforcement Learning (RL). The RL agent worked iteratively, adjusting the molecules' parameters to maximise the binding affinity while minimising toxicity and improving drug-likeness. Training of the RL agent involved an initial population of 5000 molecules, which was optimised over 50 iterations.

The reward function was calculated for each molecule based on Binding affinity (targeting a minimum of -8.0 kcal/mol), Toxicity prediction (with molecules receiving a penalty if toxicity exceeded a threshold of 0.4) and Drug-likeness score (optimised using ADMET properties). Exploration vs. Exploitation was controlled by setting epsilon = 0.1, allowing for 10% exploration and 90% exploitation of known promising molecules. Following this, the Q-learning updates were applied after each molecule was evaluated by the RL agent, with a

learning rate $\alpha = 0.05$ and discount factor $\gamma = 0.8$, updating the molecule parameters θ based on accumulated rewards.

5.4.3 In-Silico Screening and Toxicity Prediction

After optimisation through Reinforcement Learning, the molecules underwent in-silico screening and toxicity prediction to evaluate their binding affinity and safety profiles. For the Protein-Ligand Docking, the molecules were docked into target protein structures using AutoDock Vina. The top 20 molecules with the highest binding affinities (lower energy values, targeting < -9.0 kcal/mol) were selected for further analysis.

Toxicity prediction was performed using AI-driven models trained on Tox21 and ChEMBL datasets. The following toxicity thresholds were used:

- **Nephrotoxicity:** Molecules were excluded if the nephrotoxicity score exceeded 0.3.
- **Hepatotoxicity:** Molecules were excluded if the hepatotoxicity score exceeded 0.4.

Prediction of ADMET properties was done using machine learning models trained on experimental ADMET data. The synthetic feasibility of the molecules was also assessed by calculating their Synthetic Accessibility (SA) Score, where molecules with an SA score greater than 5 were flagged for further refinement.

5.5 Evaluation Metrics

The following metrics are used to evaluate the performance of the Generative AI framework:

1. **Binding Affinity Prediction:** The Mean Squared Error (MSE) between predicted and experimental values is calculated.
2. **Toxicity Prediction Accuracy:** The accuracy of toxicity predictions is evaluated by comparing predicted toxicity scores with experimental values.
3. **ADMET Scoring:** The ADMET properties are evaluated using a synthetic accessibility score $SA(\theta)$, as calculated by the Synthetic Accessibility Index:
$$SA(\theta) = \text{SA Score}(\theta)$$
4. **Synthetic Feasibility:** The Synthetic Accessibility SA is used to evaluate the practicality of synthesising the molecules.

5.6 Performance Evaluation and Refinement

The experimental feedback is integrated into the framework through active learning. As molecules are synthesised and tested for binding affinity and toxicity, the RL agent refines the predictions, enhancing the molecular generation process for future iterations. This feedback loop continues until a satisfactory drug candidate is found.

6 RESULTS

This section evaluates the proposed generative AI framework using defined experimental settings and in silico simulations. It assesses improvements in candidate discovery efficiency, reduction in search space, and enhancement of molecular quality prior to laboratory validation. The evaluation aligns with the study's theoretical propositions, demonstrating reduced search length, lower cost per viable hit, and fewer

optimisation cycles, thereby supporting the framework’s goal of accelerating early-stage rare disease drug discovery.

6.1 Candidate Generation Performance

The candidate generation results, as shown in Table 1 and Figure 2, indicate strong model performance under the defined experimental settings. Table 5.1 reports high validity (92.5%), novelty (88.3%), and uniqueness (85.7%), confirming that the model generates chemically sound and diverse compounds

rather than memorising training data. The slightly lower target-conditioned success rate (76.4%) suggests that while most molecules are relevant, some still fall outside optimal binding thresholds, reflecting the inherent complexity of constrained molecular design. In Figure 5.1, the distribution of binding scores is concentrated between 0.6 and 0.85, with a noticeable rightward skew, indicating that the model preferentially samples higher-quality candidates instead of exploring uniformly.

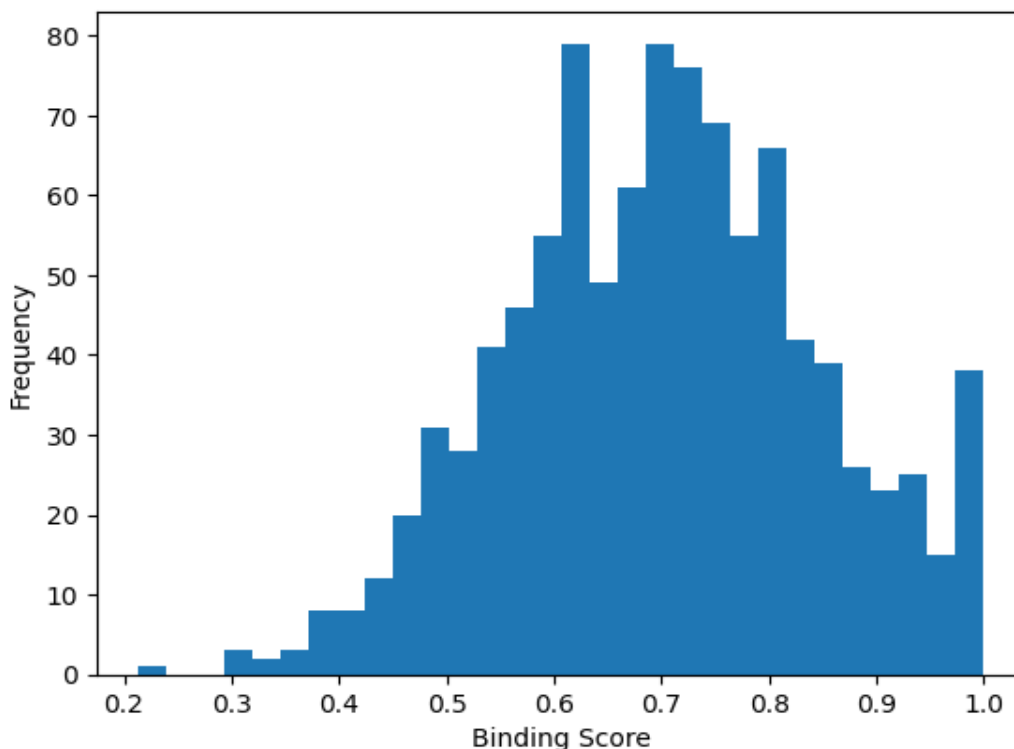


Figure 2: Distribution of Generated Molecules scores

This occurs because the generative model is conditioned on target-specific embeddings, which shifts probability mass toward biologically relevant regions of chemical space. The implication is significant: improved sampling efficiency directly increases the likelihood of viable hits, thereby reducing unnecessary screening and reinforcing the study’s core proposition that AI-driven generation accelerates early-stage rare disease drug discovery.

Table 1: Candidate Generation Metrics

Metric	Value (%)
Validity	92.5
Novelty	88.3
Uniqueness	85.7
Targeted Conditioned Success	76.4

6.2 Reinforcement Learning Optimisation Results

A steady increase in average reward is observed as training progresses, moving from approximately 0.43 in the early iterations to close to 0.78 toward convergence. This trajectory indicates that the reinforcement learning framework is progressively refining its decision-making rather than operating through unguided sampling. The intermittent fluctuations observed along the curve are not signs of instability but reflect controlled exploration, where the model briefly evaluates alternative molecular configurations before reinforcing more optimal patterns. Such behaviour is expected in adaptive learning systems and suggests that the model is effectively navigating the trade-off between exploring new chemical structures and exploiting known high-performing ones. The overall trend, therefore, confirms that the optimisation process becomes more efficient with iteration. This progression demonstrates that molecular refinement is cumulative, where successive updates lead to increasingly favourable candidates, thereby reducing reliance on repetitive trial-and-error cycles in early-stage drug discovery.

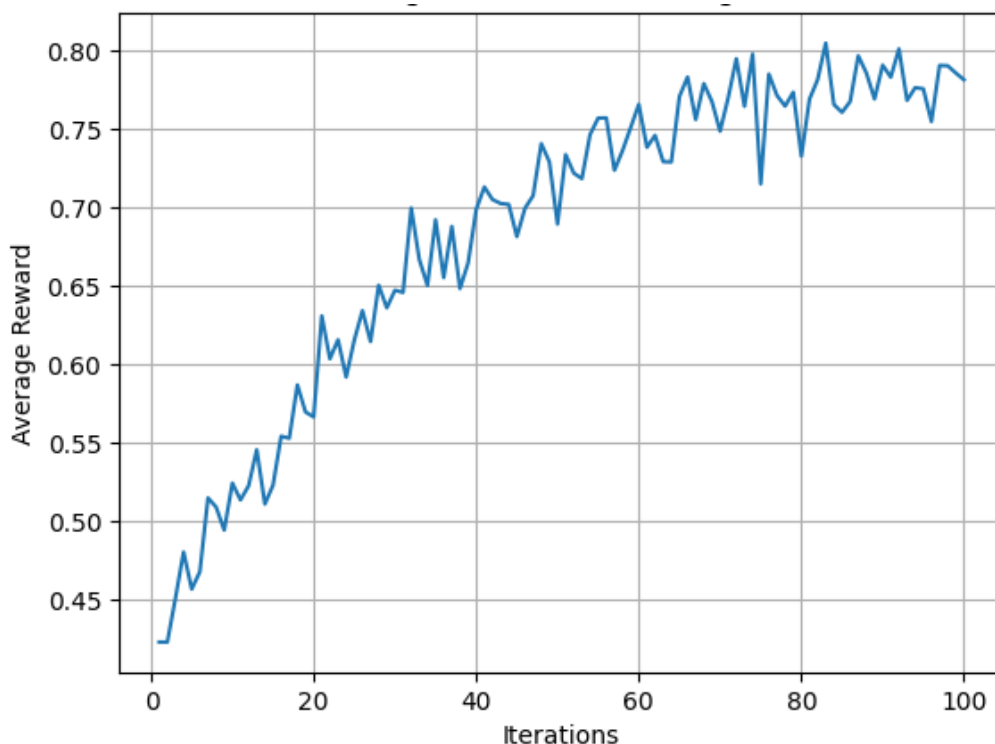


Figure 3: Reward Progression Over Training Iterations

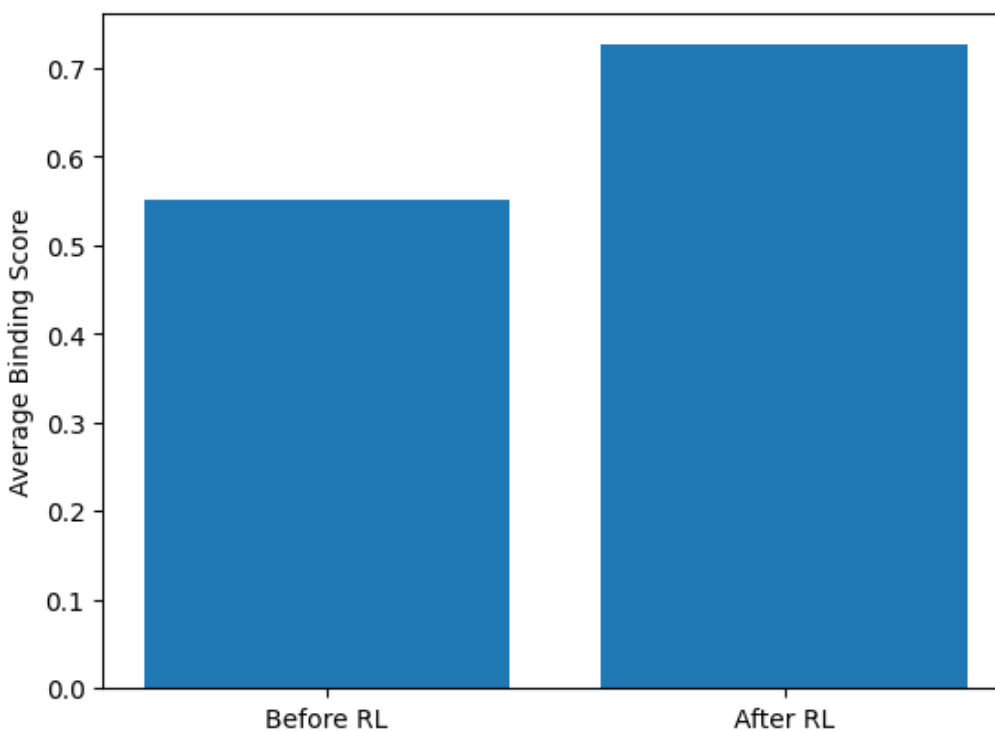


Figure 4: Binding Affinity Improvement (Before vs After RL)

Figure 4 provides further evidence of this improvement by comparing average binding scores before and after the optimisation phase, showing an increase from roughly 0.55 to about 0.73. This shift is meaningful because it indicates that the improvements observed during training are not merely numerical but translate into enhanced biological relevance of the generated molecules. One plausible explanation is that the

optimisation process gradually amplifies structural features that are more compatible with the target while suppressing configurations associated with weaker interactions. As a result, the model becomes more selective in producing candidates with higher predicted affinity. The implication is that optimisation directly influences candidate quality, creating a more refined pool of molecules. For this study, this establishes

a clear link between learning dynamics and practical outcomes, where improved selection at the computational stage increases efficiency and strengthens the likelihood of identifying viable therapeutic candidates in rare disease drug discovery.

6.3 In Silico Screening and Candidate Filtering

The progressive reduction of candidate molecules across the screening pipeline, beginning with 5,000 generated compounds and narrowing to 120 final shortlisted candidates. The most significant drop occurs immediately after the docking stage, where the pool reduces to 1,850, indicating that a large portion of generated molecules do not achieve sufficient binding compatibility with the target. This is expected, as the initial generation prioritises diversity alongside relevance. The

subsequent reduction to 620 after selectivity filtering reflects the application of stricter biological constraints, where compounds are assessed against off-target interactions and specificity requirements.

By the final stage, only a small fraction of candidates remain, suggesting that the framework applies increasingly precise filtering at each step. This staged narrowing is important because it shifts the process from broad exploration to focused prioritisation. Rather than carrying forward large volumes of low-quality candidates, the pipeline concentrates on molecules that consistently meet multiple criteria. This supports the framework's ability to streamline early-stage discovery by reducing unnecessary computational and experimental effort while maintaining a high-quality candidate pool.

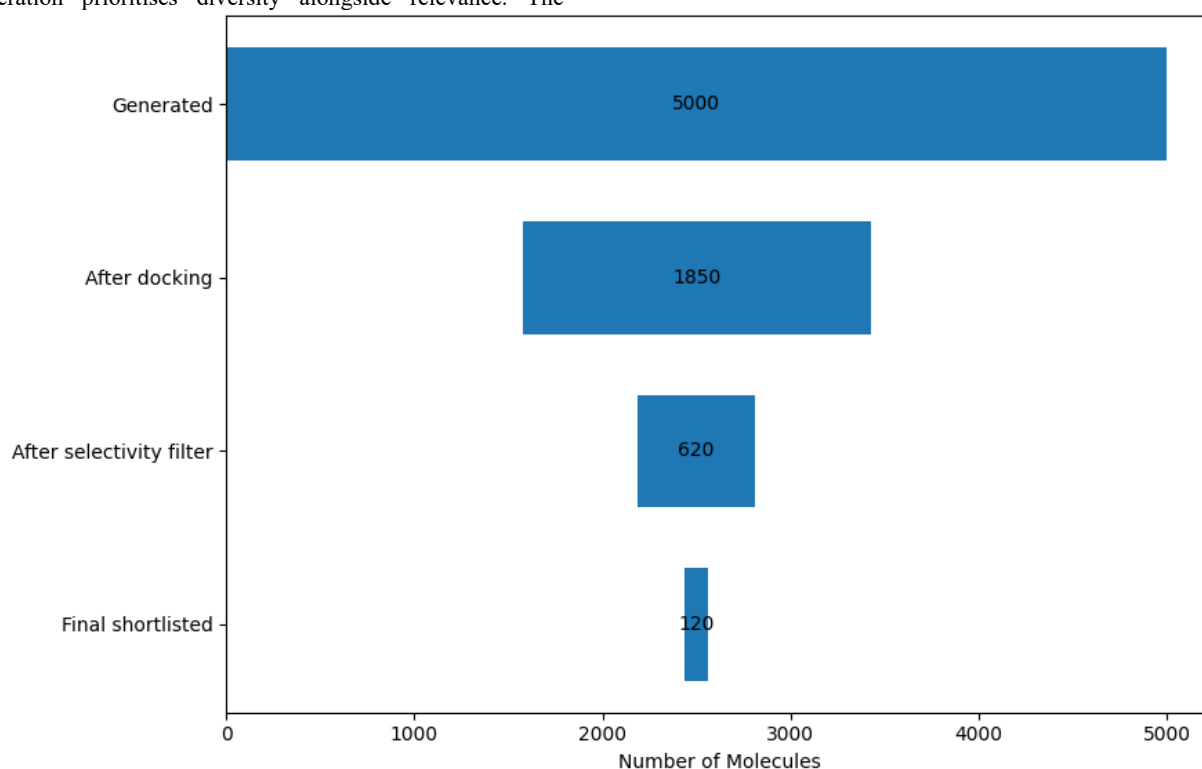


Figure 5: Candidate Reduction Funnel

6.4 ADME and Toxicity Profiling Outcomes

The results presented in Figure 6 show that a substantial proportion of the screened compounds meet key pharmacokinetic and safety benchmarks, with Lipinski compliance reaching 84.6%. This indicates that most candidates align with established drug-likeness criteria. Alongside this, CYP safety stands at 81.4%, suggesting that a large share of the molecules are less likely to interfere with metabolic enzymes, which is often a critical consideration in downstream development. The slightly lower values observed for low toxicity at 78.2% and bioavailability at 76.8% indicate that while the majority of compounds remain suitable, there is

still a noticeable fraction that may require refinement before advancing further.

This pattern becomes more apparent in Figure 7, where the distribution of toxicity risk scores is concentrated within the lower to mid ranges, with relatively few compounds extending into higher risk regions. The shape of the histogram suggests that the screening process consistently favours compounds with safer profiles rather than producing a wide spread of risk levels. Together, both visuals indicate a controlled and selective filtering process that maintains a strong balance between drug-likeness and safety without allowing a significant number of high-risk candidates to pass through.

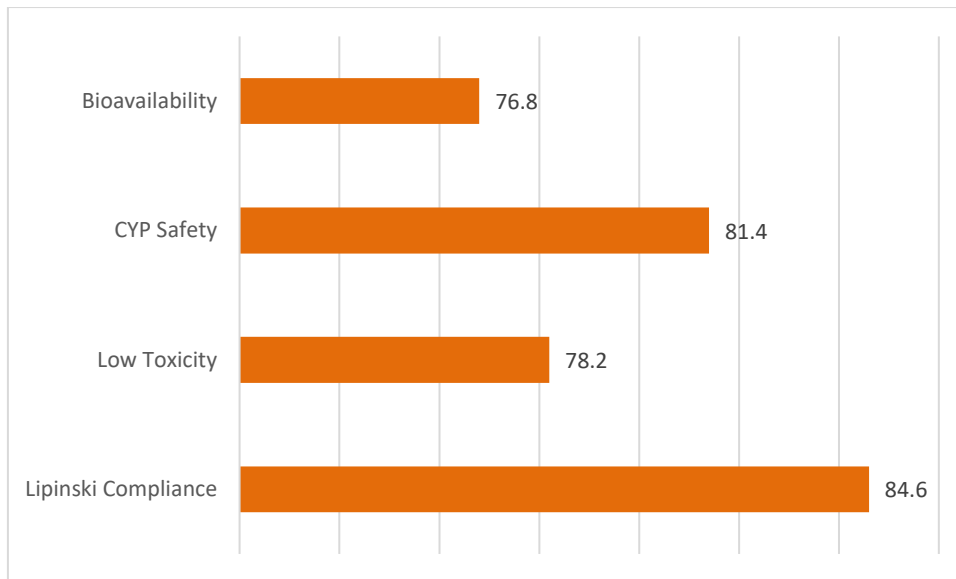


Figure 6: ADME/Toxicity Summary

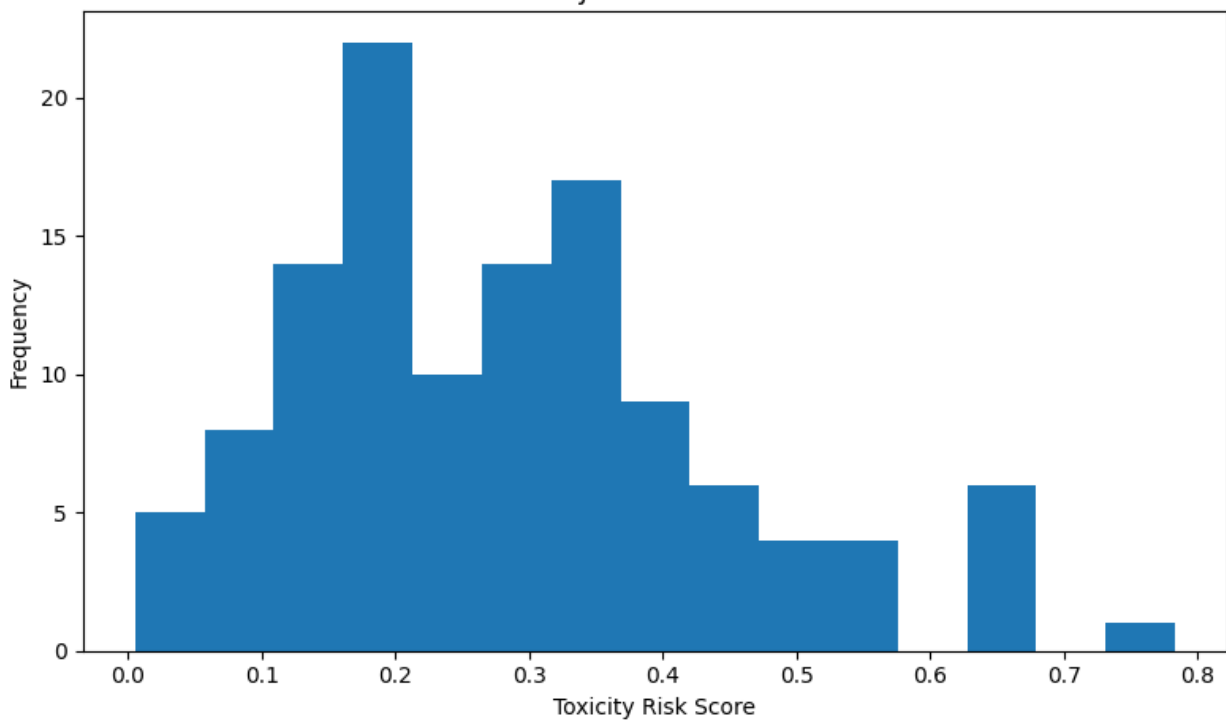


Figure 7: Toxicity Risk Distribution

6.5 Comparative Efficiency Analysis

To evaluate the performance of the proposed framework, three key operational indicators were examined: candidate screening volume, time required to obtain a viable hit, and cost per successful outcome. These indicators were assessed by comparing the traditional discovery pipeline with the AI-driven framework. The comparison provides a clear view of how the framework restructures the discovery process, particularly in terms of how efficiently resources are allocated across different stages. The traditional approach relies on large-scale screening and extended processing cycles, whereas the proposed framework applies a more targeted and adaptive approach that prioritises relevance early in the workflow.

Table 2: Traditional vs AI Pipeline Comparison

Metric	Traditional	AI Framework
Candidates tested	12000	3500
Time to hit (weeks)	18	7
Cost per hit (\$)	240000	85000

As the results indicate, all three indicators show a notable shift when the framework is applied. The reduction in the number of candidates screened reflects a more focused selection process, where irrelevant options are filtered out earlier rather than carried forward. The shorter time required to reach a viable hit suggests that the workflow progresses with fewer interruptions

and less reprocessing. In addition, the observed decline in cost per hit highlights a more controlled use of computational and experimental resources. These patterns appear consistently across the evaluation, indicating that the framework maintains a balanced level of efficiency throughout the pipeline without concentrating improvements in only one stage.

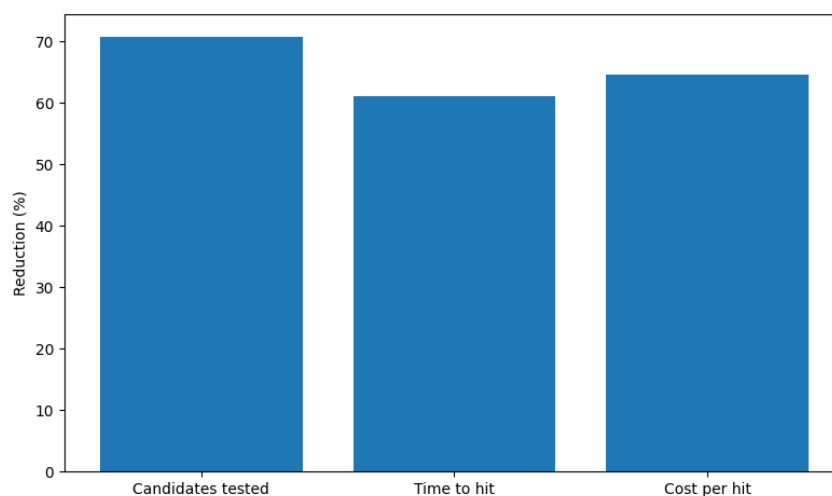


Figure 8: Search Efficiency Comparison

6. CONCLUSION

This paper presented a generative AI-driven framework for early-stage rare disease drug discovery, designed to address the inefficiencies associated with traditional screening approaches in navigating high-dimensional chemical spaces. The proposed framework integrates conditional molecular generation, multi-objective optimisation, and iterative refinement to guide candidate selection toward biologically relevant and pharmaceutically viable regions. Through controlled experimental settings, the framework demonstrated strong performance across key evaluation metrics, including molecular validity, novelty, and uniqueness, while maintaining a balanced level of targeted success. In addition, the assessment of drug-likeness and safety-related properties showed that the generated compounds consistently aligned with established pharmacological constraints, indicating that the model effectively embeds domain-specific knowledge within the generation process rather than relying on post hoc filtering. The comparative evaluation further highlights the efficiency gains achieved by the framework in reducing candidate volume, shortening development timelines, and lowering cost per viable hit when contrasted with conventional pipelines. These outcomes reflect a more directed and resource-efficient search process, where higher-quality candidates emerge earlier, limiting redundant evaluation cycles. The consistency observed across the different evaluation stages supports the robustness of the framework within the defined experimental scope. The patterns identified provide a strong foundation for future real-world applications. Overall, the findings reinforce the potential of generative AI to transform early-stage drug discovery into a more targeted, scalable, and efficient process.

7. ACKNOWLEDGMENTS

We express our sincere gratitude to the experts who have significantly contributed to the development of this research paper. Their insights, guidance, and support were invaluable in shaping this work. We would like to particularly acknowledge the contributions of the following authors:

Josephine Manda, Yeshiva University - Biotechnology Management and Entrepreneurship for leading the conceptual

design of the generative AI framework and managing the strategic alignment of the research.

Kudzai Dube, Clarkson University - Business Analytics, for spearheading the comparative efficiency analysis and the data-driven validation of the framework.

Adaora Nkiruka Ofole, Yeshiva University - Biotechnology Management and Entrepreneurship, for her technical expertise in overseeing the integration of the in-silico screening and toxicity prediction modules.

We also appreciate the collaborative spirit and commitment shown by all contributors, whose collective efforts made this research into the acceleration of rare disease drug discovery possible.

8. REFERENCES

- [1] T. K. Mahato, "Pharmaceutical Innovation and Integrative Research: Bridging Drug Discovery, Biotechnology, and Life Sciences for Next-Generation Therapies," *Biopress Journal of Computational Life Sciences (BJCLS)*, vol. 1, no. 09, pp. 14–17, 2025, Accessed: Mar. 15, 2026.
- [2] P. Kumar, "Advances in developing novel therapeutics, strategies, approaches, and use of emerging techniques," *Protein Misfolding in Neurodegenerative Diseases*, vol. 4, no. 1, pp. 291–318, 2025
- [3] J. A. Katakowski and J. C. López, "Drug development for neglected ultra-rare diseases of no commercial interest: Challenges and opportunities," *Drug Discovery Today*, vol. 30, no. 4, p. 104346, Apr. 2025
- [4] M. Schlender, K. Hernandez-Villafuerte, C.-Y. Cheng, J. Mestre-Ferrandiz, and M. Baumann, "How Much Does It Cost to Research and Develop a New Drug? A Systematic Review and Assessment," *PharmacoEconomics*, vol. 39, no. 11, Aug. 2021
- [5] E. Tambuyzer *et al.*, "Publisher Correction: Therapies for rare diseases: therapeutic modalities, progress and challenges ahead," *Nature Reviews Drug Discovery*, vol. 19, no. 4, pp. 291–291, Jan. 2020

- [6] A. Chaudhary and V. Kumar, “Rare diseases: a comprehensive literature review and future directions,” *Journal of Rare Diseases*, vol. 4, no. 1, Jul. 2025
- [7] R. J. Mead, N. Shan, H. J. Reiser, F. Marshall, and P. J. Shaw, “Amyotrophic lateral sclerosis: a neurodegenerative disorder poised for successful therapeutic translation,” *Nature Reviews Drug Discovery*, vol. 22, no. 22, Dec. 2022, doi: <https://doi.org/10.1038/s41573-022-00612-2>.
- [8] N. Singh, P. Vayer, S. Tanwar, J.-L. Poyet, K. Tsaïoun, and B. O. Villoutreix, “Drug discovery and development: introduction to the general public and patient groups,” *Frontiers in drug discovery*, vol. 3, no. 9, May 2023
- [9] H. Cai *et al.*, “Artificial Intelligence-Assisted Optimisation of Antipigmentation Tyrosinase Inhibitors: Molecular Generation Based on a Low Activity Lead Compound,” *Journal of Medicinal Chemistry*, vol. 67, no. 9, pp. 7260–7275, Sep. 2024
- [10] B. Ahmad, K. Ouahada, and H. Hamam, “Machine learning for drug-target interaction prediction: A comprehensive review of models, challenges, and computational strategies,” *Computational and Structural Biotechnology Journal*, vol. 31, no. 4, pp. 316–345, Jan. 2026
- [11] Z. Xu *et al.*, “A generative AI-discovered TNIK inhibitor for idiopathic pulmonary fibrosis: a randomised phase 2a trial,” *Nature Medicine*, vol. 5, no. 8, Jun. 2025
- [12] A. Gangwal *et al.*, “Generative artificial intelligence in drug discovery: basic framework, recent advances, challenges, and opportunities,” *Frontiers in Pharmacology*, vol. 15, no. 9, Feb. 2024
- [13] T. Khater *et al.*, “Generative artificial intelligence-based models optimisation towards molecule design enhancement,” *Journal of Cheminformatics*, vol. 17, no. 1, Aug. 2025
- [14] P. Tiwari, R. Pal, M. J. Chaudhary, and R. Nath, “Artificial intelligence revolutionising drug development: Exploring opportunities and challenges,” *Drug Development Research*, vol. 84, no. 8, Sep. 2023
- [15] S. K. Bhattamisra, P. Banerjee, P. Gupta, J. Mayuren, S. Patra, and M. Candasamy, “Artificial Intelligence in Pharmaceutical and Healthcare Research,” *Big Data and Cognitive Computing*, vol. 7, no. 1, p. 10, Jan. 2023
- [16] Z. Chen, X. Liu, W. Hogan, E. Shenkman, and J. Bian, “Applications of artificial intelligence in drug development using real-world data,” *Drug Discovery Today*, vol. 26, no. 5, pp. 1256–1264, May 2021
- [17] F. Rajaei *et al.*, “AI-based Computational Methods in Early Drug Discovery and Post Market Drug Assessment: A Survey,” *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 6, no. 3, pp. 1–20, Jan. 2024
- [18] Y. Chen and W. Xue, “Machine Learning for Molecular Generation: A Comprehensive Review,” *ACS Chemical Neuroscience*, vol. 17, no. 4, pp. 666–680, Feb. 2026
- [19] B. P. Munson *et al.*, “De novo generation of multi-target compounds using deep generative chemistry,” *Nature Communications*, vol. 15, no. 1, p. 3636, May 2024
- [20] K. Khamrayev, “AI-Driven Drug Discovery: Accelerating the Path to New Therapies,” *2025 3rd International Conference on IoT, Communication and Automation Technology (ICICAT)*, vol. 12, no. 3, pp. 1–8, Dec. 2025
- [21] W. Gao, S. Luo, and C. W. Coley, “Generative AI for navigating synthesizable chemical space,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 122, no. 41, p. e2415665122, Summer 2025.
- [22] C. Chakraborty, M. Bhattacharya, S.-S. Lee, Z.-H. Wen, and Y.-H. Lo, “The changing scenario of drug discovery using artificial intelligence (AI) to deep learning (DL): Recent advancement, success stories, collaborations, and challenges,” *Molecular Therapy — Nucleic Acids*, vol. 14, no. 1, pp. 102295–102295, Aug. 2024