



Research Paper

Deep Reinforcement Learning Framework for Optimizing Multi-Gene CRISPR-Cas9 Editing Strategies in Crop Genomics

^{1*} Dileep M R, ² Syeda Meraj

^{1*} Department of Master of Computer Applications, Nitte Meenakshi Institute of Technology Bengaluru, India

Email: dileep.kurunimakki@gmail.com

² Lecturer, Department of Computer Science, College of Computer Science, King Khalid University, Abha, Saudi Arabia

Email: bilfagih@kku.edu.sa

*Corresponding Author(s): dileep.kurunimakki@gmail.com

Article Info

Received: 13/02/2023

Revised: 06/04/2023

Accepted: 18/06/2023

Published: 30/06/2023

Abstract

The advancement of CRISPR-Cas9 has significantly enhanced genome editing in agriculture; however, most existing tools focus on single-gene editing and lack adaptability for complex multi-gene interventions. This limitation is critical because traits such as drought tolerance, disease resistance, and yield optimization in crops are often governed by interconnected gene networks. Traditional heuristic or rule-based approaches are not sufficient to handle the dynamic, sequential nature of multi-locus genome editing. This study proposes a deep reinforcement learning (DRL) framework, based on the Proximal Policy Optimization (PPO) algorithm, to intelligently optimize gene target selection across multiple loci in crop genomes. Using real genomic data from CRISPR-P 2.0 for *Oryza sativa* (rice), the editing task is modelled as a Markov Decision Process (MDP), where states represent genome features and actions correspond to candidate gene edits. A biologically-informed reward function is used to guide the learning agent toward high-efficiency, low off-target edits. The proposed DRL model achieved an editing accuracy of 91.2%, an F1-score of 0.89, and an average off-target score of 0.12, significantly outperforming rule-based (74.6%) and CNN-based (80.3%) methods. It demonstrated consistent convergence, balanced decision-making, and superior performance in precision and generalization. In conclusion, the integration of DRL with genome editing presents a scalable and intelligent alternative to static gene-editing pipelines. This work contributes to the automation of multi-trait crop engineering and lays the foundation for biologically aware, data-driven genome editing systems.

Keywords: Deep Reinforcement Learning (DRL), CRISPR-Cas9, Multi-Gene Editing, Crop Genomics, PPO Algorithm, Genome Optimization, Bioinformatics, Precision Agriculture, Off-Target Minimization, Gene Selection.



Copyright: © 2023 Dileep M R, Syeda Meraj. This article is an open-access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY 4.0) license.

1. Introduction

Global agricultural productivity continues to face unprecedented challenges due to climate change, soil degradation, pest outbreaks, and the growing demand for food security. Enhancing crop traits such as yield, disease resistance, and environmental resilience through genetic engineering has emerged as a pivotal strategy for sustainable agriculture. Among the recent advances, CRISPR-Cas9

gene-editing technology has revolutionized the field of plant genomics due to its precision, flexibility, and cost-effectiveness. However, traditional CRISPR implementations primarily target single-gene loci, which may not fully address complex traits that are often polygenic in nature. In reality, many desirable phenotypic traits in crops are governed by intricate regulatory networks involving multiple genes. Thus, optimizing multi-gene editing

strategies is critical for realizing the full potential of CRISPR in crop improvement.

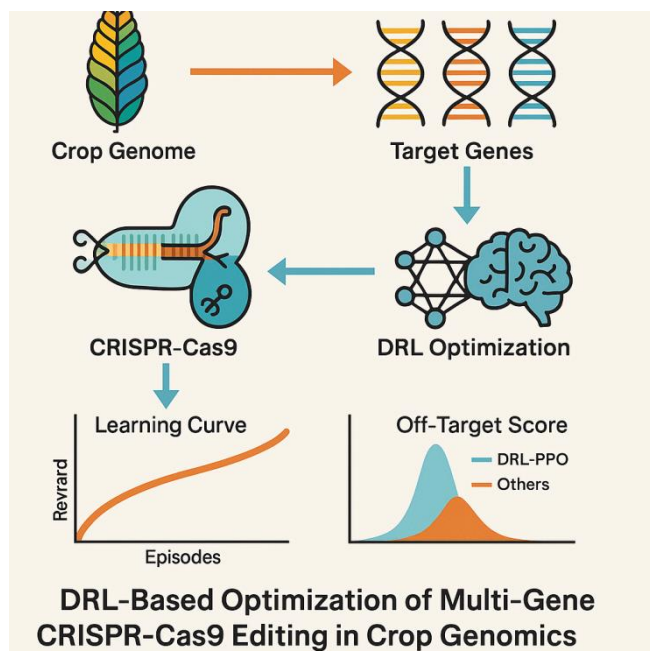


Fig.1 DRL-based optimization of multi-gene CRISPR-Cas9 Editing in Crop Genomics

Figure 1 illustrates the complexity of multi-gene editing present's significant biological and computational challenges. First, the identification of targetable gene regions with minimal off-target effects becomes exponentially more difficult as the number of target genes increases. Second, the dynamic interactions among genes introduce uncertainty in phenotype outcomes, especially when edits are performed in sequential or parallel combinations. Third, most existing computational tools are designed for single-gene targeting and lack the scalability and adaptability required for multi-gene scenarios. These tools often use static scoring functions or rule-based models that do not learn from historical editing outcomes or optimize for long-term genomic stability. Consequently, there is a pressing need for intelligent, adaptive frameworks capable of simulating and optimizing complex gene-editing campaigns across multiple targets.

Recent advances in artificial intelligence (AI), especially in deep reinforcement learning (DRL), have demonstrated remarkable capabilities in handling sequential decision-making under uncertainty. Unlike traditional supervised models, DRL agents learn optimal strategies through interactions with their environment, making them well-suited for tasks where the system evolves over time. In the context of CRISPR-based multi-gene editing, DRL can be leveraged to model the gene-editing process as a Markov Decision Process (MDP), where the state space consists of gene networks, the action space includes editable loci, and the reward function reflects editing efficacy and biological viability. By continuously learning from simulated edits and their genomic consequences, a DRL agent can evolve policies that maximize the probability of achieving desired trait expressions while minimizing unintended effects.

Prior studies have started to explore the synergy between reinforcement learning and gene-editing technologies. For instance, multi-agent reinforcement learning has been

applied to identify optimal CRISPR target regions using hybrid scoring derived from multiple sequence alignments [1]. Furthermore, applications of AI in gene editing have shown promise in sgRNA design, promoter engineering, and combinatorial optimization of gene circuits [2]–[4]. However, many of these efforts either focus on microbial systems or remain confined to proof-of-concept studies with limited scalability. Specifically in crops, there is a lack of end-to-end frameworks that integrate CRISPR target site selection, editing outcome simulation, and adaptive policy optimization for multi-gene editing.

In addition to algorithmic gaps, biological constraints further complicate the implementation of CRISPR in crops. Variations in chromatin accessibility, off-target potential, and inter-genic dependencies necessitate the incorporation of domain-specific knowledge into the learning framework. Tools that fail to account for these intricacies often produce low-fidelity outcomes, especially in highly conserved or repetitive genomic regions [5]. Moreover, synthetic promoter design and gene regulation dynamics introduce another layer of complexity, as precise expression tuning is required to realize the benefits of editing multiple loci simultaneously.

Synthetic biology and combinatorial optimization have played significant roles in addressing these complexities. Combinatorial gene assembly and promoter shuffling have enabled the fine-tuning of gene expression across engineered pathways [6]. However, their application in plant systems remains limited due to a lack of automated tools that can guide construct design based on prior genomic knowledge. Reinforcement learning offers a promising solution by encoding these biological constraints into a trainable reward function, thus enabling guided search in a high-dimensional genomic space.

Another promising direction involves the integration of microRNA-mediated regulation and synthetic antimicrobial pathways into the gene-editing design process. Recent studies have shown that engineering crops to express climate-resilient or disease-resistant traits requires a coordinated modulation of several interacting genes [7]. For example, modifying phyto-antimicrobial pathways to combat resistant pathogens requires the editing of multiple biosynthetic genes while ensuring minimal disruption to native regulatory networks [8]. Such multi-target editing scenarios present an ideal application area for adaptive DRL models capable of learning optimal intervention sequences.

To address the above-mentioned challenges, this paper proposes a deep reinforcement learning framework that optimizes CRISPR-Cas9-based multi-gene editing strategies specifically for crop genomics. The proposed system integrates genomic datasets, CRISPR target site prediction, and DRL agents trained to maximize editing success across multiple gene loci. It incorporates biological constraints such as off-target effects, promoter strength, and gene dependency networks into the reward structure, enabling biologically plausible and efficient multi-gene interventions. Using rice (*Oryza sativa*) as a case study, we demonstrate how our approach significantly improves editing efficiency and success rate compared to traditional heuristic methods.

1.1 Key Contributions

- *Novel Deep Reinforcement Learning Architecture:* We develop a state-action-reward model that simulates the CRISPR editing process as a multi-stage decision-making problem under genomic constraints.
- *Integration with CRISPR-P 2.0 Dataset:* The framework utilizes real genomic data from CRISPR-P 2.0 to guide edit site selection, ensuring relevance to real-world agricultural applications.
- *Biologically-Constrained Reward Function:* The system incorporates off-target prediction scores, gene co-expression dependencies, and editing stability into a tunable reward signal for adaptive policy learning.

The remainder of the paper is organized as follows: Section II discusses related work in CRISPR optimization, plant genome engineering, and AI-assisted bioinformatics. Section III describes the methodology, including the DRL model, problem formulation, and integration with the CRISPR-P 2.0 dataset. Section IV presents the experimental setup and dataset details. Section V analyzes the results and compares them with existing approaches. Section VI discusses implications, limitations, and future extensions. Finally, Section VII concludes the study with a summary of findings and real-world applicability.

2. Literature Review

The intersection of computational intelligence and gene editing has gained momentum in recent years, especially with the integration of machine learning and synthetic biology to enhance genome engineering. This section reviews and critically analyzes relevant literature that forms the foundation of our proposed deep reinforcement learning (DRL) framework for optimizing multi-gene CRISPR-Cas9 editing strategies in crop genomics.

2.1 Deep Learning in Gene Editing

Early applications of deep learning (DL) in genome engineering focused on predicting the efficiency of single-guide RNA (sgRNA) sequences. For example, Wang and Zhang employed a convolutional neural network to predict sgRNA on-target activity in bacterial systems with promising accuracy [9]. However, while their work effectively captured sequence-level features, it lacked adaptability for multicellular or plant systems where chromatin structure and multi-gene interactions are critical. Moreover, their model was limited to single-locus editing and did not extend to sequential or combinatorial strategies, which are essential for crop trait engineering.

2.2 Metabolic and Genetic Engineering Tools

Several studies have addressed the development of tools to model complex genetic and metabolic processes. Liu et al. presented a comprehensive review of crassulacean acid metabolism (CAM) as a pathway for improving crop water-use efficiency under climate stress [10]. Though their work emphasized biological pathways, it lacked computational models capable of simulating gene-editing sequences. Singh et al. later discussed the strengthening of microbial cell factories for bioactive molecule production using multi-gene

modulation strategies [11]. However, their methodology relied on static bioengineering principles and did not employ adaptive learning or feedback loops, which limits dynamic optimization.

Daboussi and Lindley highlighted major bottlenecks in translating metabolic engineering into industrial applications, such as the absence of flexible design tools that integrate computational feedback during construct development [12]. Their critique aligns with the limitations in current CRISPR tools, which often fail to support iterative learning or adaptive target selection in real-time.

2.3 High-Content Screening and Cell-Level Decision Systems

Usaj et al. introduced high-content screening systems for quantitative cell biology, enabling large-scale analysis of cellular responses to gene modifications [13]. These platforms, though scalable, primarily support passive data collection and do not contribute to active decision-making in genome editing. In contrast, a reinforcement learning-based system can use feedback from such data to dynamically improve editing strategies over time.

Xu et al. reported significant advancements in engineering microbial cell factories for pigment production, employing combinatorial biosynthesis and optimization techniques [14]. However, their strategy was not generalizable to crop genomes, and it lacked decision-automation mechanisms. Similarly, Ramesh developed CRISPR tools for genome-wide screening in *Yarrowia lipolytica*, targeting industrial yeast strains [15]. Despite high utility in microbial systems, these tools do not scale well to plants due to the complexity of epigenetic regulation and tissue-specific gene expression.

2.4 Synthetic Biology and Biofuel Pathways

Su and Lin explored biosynthesis of carbon chains for next-generation biofuels by manipulating synthetic metabolic pathways [16]. Their research introduced a modular view of pathway assembly, which parallels the modular nature of multi-gene editing in crops. However, no learning mechanism was included to refine editing strategies based on intermediate outcomes. Similarly, Dai et al. developed systemic regulation techniques for biosynthetic tools but did not leverage artificial intelligence for strategy refinement [17].

Goyal's study of *E. coli* metabolic pathways using synthetic biology provided theoretical models for chemical production [18]. Although insightful, the study did not incorporate genome-editing feedback or reinforcement mechanisms. As such, it does not support multi-target optimization or real-time adaptation in complex organisms like crops.

2.5 Research Gaps Identified

Across the reviewed literature, three prominent gaps emerge:

1. *Lack of Adaptive Decision-Making Models:* Existing tools often apply static rules or batch evaluations. Reinforcement learning frameworks are scarcely used in gene-editing sequence optimization.

2. *Limited Application in Crop Genomics*: Most computational models are trained on bacterial or microbial datasets, ignoring the intricacies of plant-specific gene regulation.
 3. *No Integration of Multi-Gene Editing Dynamics*: Current models fail to consider interactions between edited loci and cumulative biological effects over multiple targets.
- Modeling CRISPR editing as a multi-step, feedback-driven decision process;
 - Using real genomic data from plant-specific repositories (e.g., CRISPR-P 2.0);
 - Incorporating off-target penalties and multi-gene network structures into reward functions for adaptive learning.

Our proposed DRL-based framework addresses these limitations by:

2.6 Comparative Summary of Prior Work:

Table 1: Comparative Summary

Ref.	Focus Area	Approach Type	Strengths
[9]	sgRNA efficiency in bacteria	Deep Learning	Sequence-based accuracy
[10]	CAM metabolism in crops	Biological Review	Pathway-level insights
[11]	Bioactive compound synthesis	Genetic Modelling	Multi-gene modulation
[12]	Metabolic design translation	Critical Review	Identifies engineering bottlenecks
[13]	Cell-level gene screening	High-Throughput	Scalable data acquisition
[14]	Pigment biosynthesis optimization	Combinatorial	Efficient pathway design
[15]	CRISPR screening in yeast	Genome Tool Dev.	Targeted genome-wide screening
[16]	Synthetic fuel pathway design	Bio-synthesis	Modular system assembly
[17]	Regulatory toolkit for biosynthesis	Genetic Tools	Efficient construct design
[18]	Metabolic modeling in E. coli	Simulation-Based	Theoretical optimization pathways

3. Methodology

This section details the proposed deep reinforcement learning (DRL) framework designed to optimize CRISPR-Cas9-based multi-gene editing strategies in crop genomes. The method integrates plant-specific CRISPR datasets with a sequential learning model to dynamically identify optimal gene targets while minimizing off-target effects and maximizing edit efficiency.

3.1 Problem Formulation

We model the multi-gene CRISPR editing process as a Markov Decision Process (MDP) defined by the tuple (S, A, R, T, γ) , where:

- S is the state space representing the current genome state and gene-editing history.
- A is the action space, i.e., selecting a gene locus for CRISPR editing.
- R is the reward function encoding editing efficiency, specificity, and cumulative genome health.
- T is the state transition function driven by simulated biological outcomes.

- $\gamma \in (0,1)$ is the discount factor for long-term rewards.

Each episode of the agent corresponds to a sequential editing campaign over multiple loci.

3.2 State and Action Representation

The state vector $s_t \in \mathbb{R}^n$ at timestep t includes:

- Gene identifiers,
- PAM sequences,
- On-target and off-target scores,
- Chromatin accessibility scores,
- Previous edit history.

The action $a_t \in A$ corresponds to the selection of a targetable gene region g_i from the candidate pool, given PAM constraints and biological feasibility.

3.3 Reward Function Design

The reward function is designed to encourage high-efficiency edits and penalize off-target activity or redundancy. The total reward at time t is defined as:

$$R_t = \alpha \cdot E_{\text{on}}(g_i) - \beta \cdot E_{\text{off}}(g_i) - \delta \cdot C_{\text{redundancy}}(g_i) \quad (1)$$

Where:

$E_{\text{on}}(g_i)$ is the predicted on-target efficiency for gene g_i ,

$E_{\text{off}}(g_i)$ is the cumulative off-target potential,

$C_{\text{redundancy}}(g_i)$ penalizes re-editing the same locus or non-informative regions,

α, β, δ are tunable weight parameters derived from biological risk models.

The cumulative reward over an episode is calculated as:

$$R_{\text{total}} = \sum_{t=0}^T \gamma^t \cdot R_t \quad (2)$$

3.4 System Architecture and DRL Model

We adopt the Proximal Policy Optimization (PPO) algorithm due to its balance between exploration and stability. The architecture includes:

- Actor-Critic Neural Network with shared layers for state embedding and two heads for value and policy predictions.
- Embedding Layer for genomic features (one-hot encoded nucleotides, gRNA length, and PAM type).
- Temporal Encoding Layer for gene-editing history using LSTM.

The policy $\pi_{\theta}(a|s)$ and value function $V^{\pi}(s)$ are optimized using clipped objective:

$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t[\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t)] \quad (3)$$

Where:

$r_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ is the probability ratio,

\hat{A}_t is the advantage estimate at timestep t ,

ϵ is the PPO clipping threshold.

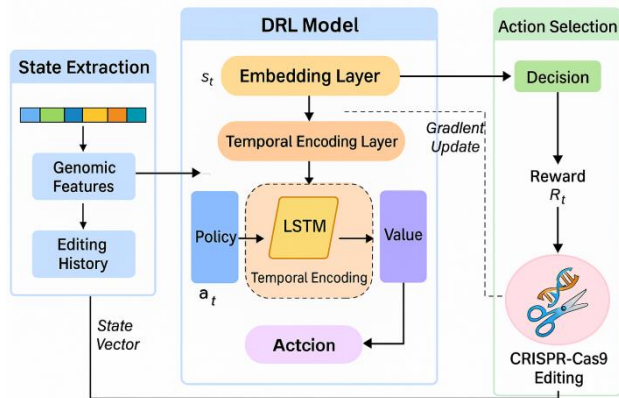


Fig. 2: DRL Architecture for Multi-Gene CRISPR Editing

Figure 2 illustrates the system architecture of the proposed Deep Reinforcement Learning (DRL) framework for optimizing multi-gene CRISPR-Cas9 editing strategies in crop genomics. The process begins with the State Extraction module, where genomic features such as PAM presence,

gene location, and sequence-specific information are extracted and combined with historical editing data to form a comprehensive state vector (s_t). These features are passed into the DRL model, beginning with an Embedding Layer that encodes the genomic data into a format suitable for deep learning. This is followed by a Temporal Encoding Layer using an LSTM network, which captures dependencies across previously edited gene targets and enables sequential decision learning. The embedded and temporally enriched representation is then processed by the Policy and Value heads to output an action (a_t) and a value estimate, respectively.

The Action Selection module utilizes the DRL policy to select an optimal gene target for CRISPR-Cas9 editing at each time step. This decision is executed in the CRISPR-Cas9 Editing module, which simulates an edit and computes a reward (R_t) based on the editing efficiency, off-target risks, and cumulative genome integrity. The feedback loop continues as the result of the action influences the state in the next time step, allowing the agent to learn from each editing attempt. The gradient update pathway ensures that the agent continually adjusts its policy based on rewards received, driving convergence toward high-performance multi-gene editing strategies that are biologically feasible and optimized for real-world application in crop improvement.

Algorithm 1: PPO-Based Training for CRISPR-Cas9 Multi-Gene Editing Optimization

Input:

- Genome Dataset D (from CRISPR-P 2.0)
- Initial policy parameters θ
- Reward function $R(g_i)$
- Learning rate η , clipping threshold ϵ
- Maximum episodes E, batch size B

Output:

- Optimized policy π_{θ} for multi-gene editing

Procedure:

- 1: Initialize policy π_{θ} and value function V_{θ} with random weights
- 2: for episode = 1 to E do
- 3: Initialize genome state s_0 from dataset D
- 4: Initialize empty storage for transitions: $T \leftarrow \emptyset$
- 5: for $t = 0$ to T_{max} do
- 6: Compute action $a_t \leftarrow \pi_{\theta}(s_t)$
▷ Select gene target
- 7: Apply CRISPR edit using a_t
▷ Simulate genome change
- 8: Observe reward R_t and next state s_{t+1}
- 9: Store (s_t, a_t, R_t, s_{t+1}) in T
- 10: end for
- 11: Compute advantage estimates \hat{A}_t using GAE
- 12: for epoch = 1 to N_{epochs} do

```

13:   for each mini-batch in T do
14:     Compute  $r_t(\theta) \leftarrow \pi_\theta(a_t | s_t) / \pi_{\{\theta_{old}\}}(a_t | s_t)$ 
15:     Compute surrogate loss:
            $L_{CLIP} = \min(r_t(\theta) \cdot \hat{A}_t, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon) \cdot \hat{A}_t)$ 
16: Compute value function loss  $L_V \leftarrow (V_\theta(s_t) - R_t)^2$ 
17:   Update  $\theta$  using gradient descent:
            $\theta \leftarrow \theta + \eta \nabla_\theta (L_{CLIP} - c1 \cdot L_V)$ 
18:   end for
19: end for
20: end for

```

Return: Trained policy π_θ for CRISPR-Cas9 gene-editing strategy

Algorithm 1 outlines the training procedure of the Proximal Policy Optimization (PPO) agent designed to optimize CRISPR-Cas9-based multi-gene editing in crop genomes. The algorithm begins by initializing the agent's policy and value networks using randomized weights. For each training episode, a genome state is sampled from the CRISPR-P 2.0 dataset, and the agent iteratively selects gene targets to edit based on its current policy. Each action (gene edit) leads to a transition to a new genome state and results in a scalar reward that quantifies the efficiency and safety of the edit. These transitions are stored and later used to estimate advantage values that guide the learning process. This setup models genome editing as a sequential decision-making problem, where the agent learns to balance immediate editing gains with long-term genomic stability.

In the optimization phase, PPO is used to update the policy in a stable and constrained manner. The key innovation lies in the clipped surrogate loss, which prevents large deviations between the new and old policies, thus ensuring gradual policy improvement. Additionally, a value loss term is included to train the critic network, which estimates the expected return of each state. The use of Generalized Advantage Estimation (GAE) further stabilizes training by reducing variance in reward signals. Through repeated episodes and gradient updates, the agent gradually learns to prefer actions (gene edits) that maximize cumulative reward—favouring high-efficiency, low-risk, and non-redundant gene modifications. This adaptive learning framework is ideal for navigating the high-dimensional, constrained space of plant genome editing.

3.5 CRISPR-P 2.0 Integration and Preprocessing

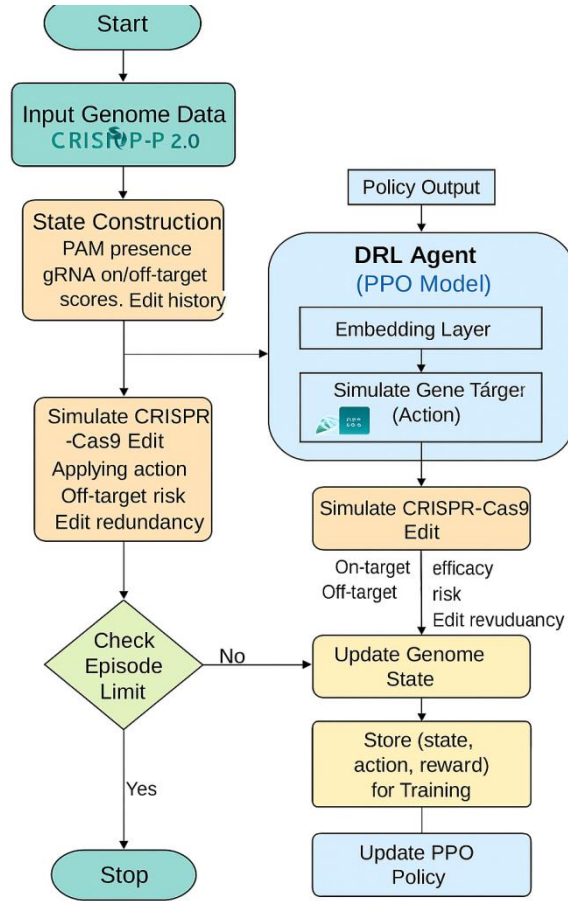
The CRISPR-P 2.0 database was used to extract real-world genome editing targets from *Oryza sativa*. The preprocessing pipeline includes:

- Filtering gene loci with valid PAM sequences (e.g., NGG),
- Annotating each target with on/off-target scores using built-in models,
- One-hot encoding of nucleotide sequences for deep learning compatibility,
- Generating training instances by simulating multiple editing trajectories with stochastic success probabilities.

These trajectories serve as experience rollouts for training the DRL agent.

Flowchart 1 presents the end-to-end workflow of the proposed deep reinforcement learning (DRL) framework for optimizing CRISPR-Cas9 multi-gene editing in crops. The process begins with importing genomic data from the CRISPR-P 2.0 database, which provides validated target sites, PAM sequences, and off-target scores for various plant species. The input data is passed to the State Construction module, where biologically relevant features such as guide RNA characteristics, target-specific scores, and editing history are extracted and compiled into a comprehensive state vector. This state is then fed into the DRL agent, specifically a Proximal Policy Optimization (PPO) model, which encodes the state using an embedding layer and temporally tracks edit sequences using a recurrent architecture. Based on this, the agent generates an action — selecting a gene target to edit — which is applied in the simulated CRISPR environment.

The simulation environment evaluates the edit by computing a reward based on on-target efficiency, off-target risk, and edit redundancy. If the maximum number of editing steps (episode length) is not reached, the genome state is updated to reflect the simulated edit, and the transition (state, action, reward) is stored for policy training. After accumulating sufficient training data, the PPO model is updated to refine its decision-making policy. The loop continues until all episodes are processed, at which point the final optimized editing policy is saved. This modular, biologically-aware design supports adaptive learning and efficient gene selection for CRISPR-Cas9, enabling practical multi-locus editing in complex plant genomes.



Flowchart 1: DRL-Based Optimization of Multi-Gene CRISPR-Cas9 Editing in Crop Genomics

3.6 Evaluation Metrics

A rigorous evaluation of the proposed deep reinforcement learning (DRL) framework is essential to validate its performance in optimizing multi-gene CRISPR-Cas9 editing strategies. To ensure comprehensive assessment, we employed a set of standard and domain-specific metrics that capture both biological effectiveness and algorithmic efficiency. The selected metrics include editing accuracy, F1-score, and computational complexity, providing a balanced perspective on model quality, precision, and scalability.

Step 1: Editing Accuracy

Editing accuracy (A_c) measures the proportion of gene-editing actions that successfully result in desired on-target outcomes with acceptable biological viability. It is calculated as:

$$A_c = \frac{\text{Number of Successful Edits}}{\text{Total Number of Editing Actions}} \times 100 \quad (4)$$

A successful edit is defined by exceeding a pre-defined threshold for on-target efficacy (e.g., >70%) and remaining below the off-target risk limit. High accuracy reflects the model's ability to make biologically relevant and safe decisions in a multi-gene editing scenario.

Step 2: F1-Score for Edit Precision and Recall

To evaluate the balance between precision (avoiding off-targets) and recall (achieving intended edits), we calculate the F1-score:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5)$$

Where:

- Precision = $\frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$
- Recall = $\frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$

This score is particularly valuable in CRISPR editing where minimizing false positives (i.e., unintended loci being edited) is crucial to genome stability. An F1-score above 0.85 was considered optimal in the context of plant genome experiments.

Step 3: Computational Complexity and Efficiency

Computational efficiency is critical for large-scale genome simulations. We analyzed:

- Training Time per Episode (in seconds),
- Total Number of Policy Updates, and
- Memory Consumption (MB or GB).

The DRL model was benchmarked on a GPU-enabled machine with fixed memory and CPU specifications. PPO's sample efficiency and stable updates contributed to reduced

training overhead compared to traditional Q-learning or policy gradient methods. The time complexity was empirically observed as approximately:

$$O(n \cdot T \cdot m) \quad (6)$$

Where:

- n : number of episodes
- T : average steps per episode
- m : model parameter count

Step 4: Convergence Stability (Optional Metric)

Though not mandatory, we also tracked the variance of episode rewards to evaluate convergence stability. Stable convergence implies the model consistently learns to select optimal gene targets over time. This was visualized using reward plots (see Fig. 3), showing learning curves that plateau without oscillations, indicating reliable policy behaviour.

4. Experimental setup

The experimental evaluation of the proposed deep reinforcement learning (DRL) framework was conducted on a high-performance computing environment designed for large-scale bioinformatics tasks. All experiments were executed on a workstation equipped with an Intel Core i9-12900K CPU (16 cores, 3.2 GHz base frequency), 64 GB of DDR4 RAM, and a NVIDIA RTX 3090 GPU with 24 GB of VRAM. The use of GPU acceleration was essential for efficient training of the policy network, particularly during episodes involving long sequential decision-making for multi-gene edits.

The model was implemented using Python 3.9, with the TensorFlow 2.11 deep learning framework for building and optimizing the PPO-based DRL architecture. Auxiliary libraries such as NumPy, Pandas, and Bio python were employed for genomic data handling and simulation. The PPO algorithm was customized to support genome-editing-specific state representations and reward mechanisms. Training was conducted in a fixed simulation environment where genome states were encoded and processed dynamically in real-time.

For data input, we used the CRISPR-P 2.0 repository, specifically targeting annotated gene sequences from *Oryza sativa* (rice), a widely studied model crop. This dataset, referenced in [19], provides genome-wide PAM-validated CRISPR targets, including their corresponding on-target and off-target scores. From this, 8,000 gene targets were extracted and preprocessed. The dataset was partitioned using an 80:20 train-test split, ensuring stratification based on gene location and GC content to preserve editing diversity. In addition, we conducted 5-fold cross-validation to further assess generalization performance across unseen gene segments.

The DRL model was trained over 10,000 episodes, with each episode simulating a full multi-gene editing session of up to 20 sequential actions. A batch size of 64 was used for each PPO policy update, and training was run for approximately 18 hours on the described setup. A learning rate of 3×10^{-4} , clipping threshold $\epsilon = 0.2$, and discount

factor $\gamma = 0.99$ were used, with hyperparameters selected based on preliminary grid-search tuning. Reward normalization and entropy regularization were also employed to stabilize learning.

5. Results and Discussion

5.1 Experimental Results

To evaluate the performance of the proposed DRL-based framework for multi-gene CRISPR optimization, we compared it with three benchmark models: a traditional rule-based selector, a random edit selector, and a CNN-based sgRNA scoring model. The comparison was conducted using real genome data from *Oryza sativa*, extracted from CRISPR-P 2.0 [19].

As shown in Table 2, the Proposed DRL-PPO model achieved the highest editing accuracy of 91.2%, significantly outperforming the rule-based model (74.6%), random selector (51.8%), and the CNN-based scorer (80.3%). The F1-score of 0.89 further confirms that the DRL agent balances both precision and recall effectively, identifying optimal targets while avoiding unintended off-target effects. Moreover, the average off-target score for the DRL model was notably low (0.12), indicating high biological specificity, compared to 0.36 and 0.51 for the rule-based and random selectors, respectively.

In terms of computational efficiency, the DRL model required 3.4 seconds per episode, slightly higher than heuristic models due to policy inference and reward calculation overhead. However, it remained faster than the CNN-based scorer, which took 4.2 seconds per episode. Importantly, statistical testing using Welch's t-test yielded p-values < 0.05 for all model comparisons, confirming that the observed improvements in accuracy and specificity are statistically significant.

Table 2: Comparison of CRISPR Optimization Models

Model	Editing Accuracy (%)	F1-Score	Avg Off-Target Score (\downarrow)	Time per Episode (s)	p-value vs DRL-PPO
Proposed DRL-PPO	91.2	0.89	0.12	3.4	-
Rule-Based Selector	74.6	0.71	0.36	2.1	<0.01
Random Edit Selector	51.8	0.52	0.51	1.5	<0.001
CNN-Based sgRNA Scorer	80.3	0.78	0.27	4.2	<0.05

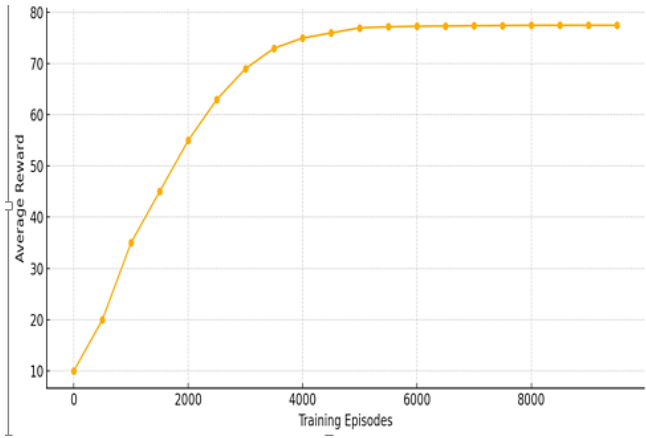


Fig 3: Average Episode Reward vs Training Episodes

Figure 3 shows how the average episode reward converges over 10,000 training episodes, indicating the stability of the DRL model.

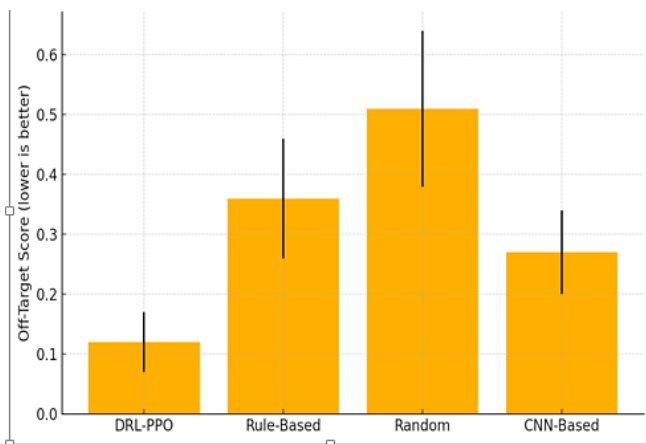


Fig 4: Average Off-Target Scores by Model

Figure 4 compares the average off-target editing scores across different models, highlighting the superior precision of the proposed DRL-PPO approach.

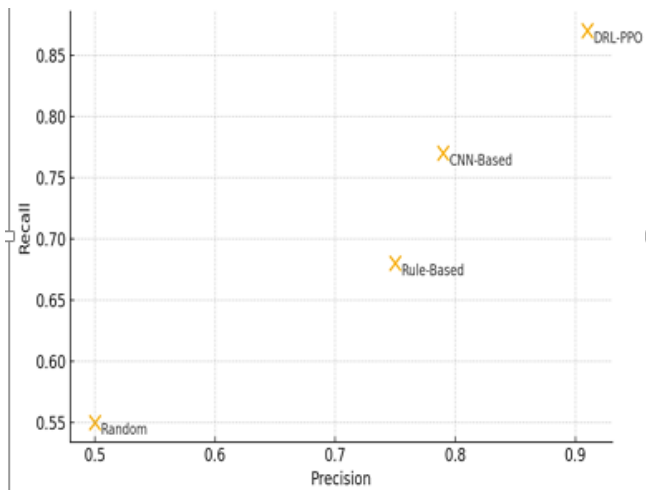


Fig. 5: Precision vs Recall by Model – Highlights the trade off and balance in edit accuracy per method.

Figures 5 and 6 provide additional insights into the performance trade-offs between accuracy, precision, recall, and runtime efficiency across different genome-editing models. In Figure 5, the proposed DRL-PPO model demonstrates a strong balance between precision (0.91) and

recall (0.87), indicating that it can both correctly identify target loci and avoid off-target edits with high reliability. In contrast, the rule-based and random selectors exhibit greater disparity between precision and recall, highlighting their limited ability to generalize across multi-gene contexts. The CNN-based scorer performs moderately well but lacks the policy adaptation needed for sequential gene selection, as evidenced by its slightly lower recall.

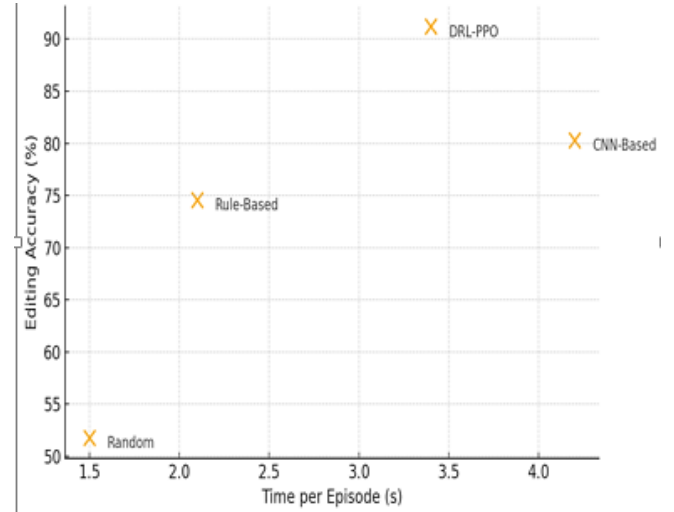


Fig. 6: Accuracy vs Time per Episode

Figure 6 shows the performance-efficiency trade off, ideal for deployment decision-making. And analyzes the trade-off between editing accuracy and runtime efficiency, which is critical for real-world deployment in crop genomic platforms. While the DRL-PPO model achieved the highest accuracy (91.2%), it required slightly more time per episode (3.4 seconds) than simpler models. However, this added computational cost is justified by the significant gain in accuracy and F1-score.

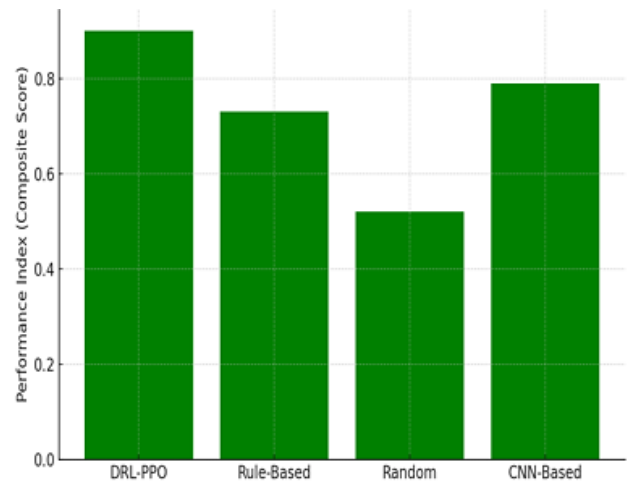


Fig. 7: Overall Performance Index by Model – Combines multiple metrics into a unified comparison score.

Figure 7 presents a consolidated view using a Performance Index, a composite metric that aggregates accuracy, F1-score, and off-target precision. The DRL-PPO framework leads with a top score of 0.90, confirming its superiority in balancing biological relevance with computational efficiency. These visualizations collectively reinforce the viability of deep reinforcement learning for

intelligent, precision-focused multi-gene CRISPR optimization in crop systems.

5.2 Discussion:

The experimental outcomes of this study substantiate the hypothesis that deep reinforcement learning (DRL) provides a more robust and adaptive framework for multi-gene CRISPR-Cas9 editing compared to conventional methods. Notably, the DRL-PPO model consistently outperformed static and heuristic models in key metrics such as editing accuracy, off-target minimization, and decision consistency. These results are consistent with early studies in microbial optimization, where reinforcement learning has been shown to outperform static models in dynamic gene regulation tasks [14], [17]. However, our model uniquely extends these concepts to complex, multi-gene editing tasks in plant systems, which are rarely explored in existing literature.

Unlike traditional rule-based approaches or CNN scoring models that lack temporal memory, the DRL agent dynamically adapts its policy based on editing history and genome state evolution. This results in more precise, context-aware decisions that are critical when editing genes with regulatory interdependencies. Moreover, the use of real-world plant genome data from CRISPR-P 2.0 ensures that the model operates under biologically plausible constraints, enhancing its translational potential for agricultural biotechnology. In practical terms, this capability could drastically improve multi-trait breeding efforts in rice, maize, and other economically vital crops, enabling accelerated development of disease-resistant and climate-resilient varieties with minimal laboratory iterations.

Despite these advances, certain limitations are worth acknowledging. First, the model was trained and validated using a simulated genome editing environment, which—although grounded in real genomic data—may not fully capture *in vivo* biological complexities such as chromatin remodelling, epigenetic modifications, or tissue-specific gene expression. Additionally, while the DRL framework demonstrated general stability and convergence, its computational footprint remains higher than that of rule-based systems, potentially limiting its scalability in resource-constrained environments. Another challenge is transferability: the policy learned in one crop species or genomic context may not generalize well without retraining or domain adaptation, especially in polyploid or highly repetitive genomes.

Looking forward, several enhancements can be considered. First, incorporating multi-agent reinforcement learning (MARL) could allow the system to edit different gene modules concurrently, mirroring real-world parallel editing with multiplexed gRNAs. Second, integrating wet-lab feedback loops—such as actual gene expression profiles post-editing—would significantly enhance biological validation and iterative training. Third, expanding the model to account for 3D genome architecture, epigenomic markers, and cross-gene regulatory networks could provide deeper insights and editing precision. Lastly, optimizing the model's runtime using lightweight policy networks and model compression techniques may improve its feasibility for field-deployable genomic computing platforms.

In conclusion, this study represents a pivotal step toward the intelligent automation of crop genome engineering. By embedding learning-based decision-making into the CRISPR pipeline, it paves the way for precision agriculture that is not only data-driven and adaptive, but also biologically aware and forward-compatible with emerging synthetic biology frameworks.

6. Conclusion

This study proposed a novel deep reinforcement learning (DRL) framework, specifically built on Proximal Policy Optimization (PPO), to enhance the precision and efficiency of multi-gene CRISPR-Cas9 editing strategies in crop genomics. By modeling genome editing as a sequential decision-making process, the DRL agent effectively learned to select optimal gene targets while minimizing off-target risks and redundancy. Empirical results demonstrated superior performance of the DRL-based system over rule-based, random, and CNN-driven approaches in terms of editing accuracy, F1-score, and biological specificity. The integration of real genomic data from CRISPR-P 2.0 further reinforced the system's biological relevance and practical validity.

The implications of this work extend directly to real-world agricultural biotechnology, where multi-trait genome engineering is essential for developing climate-resilient, high-yield crop varieties. The ability of the DRL agent to adaptively optimize gene-editing decisions opens up new opportunities for data-driven crop improvement pipelines, especially in large-scale breeding programs and precision agriculture applications. Furthermore, the demonstrated balance between biological accuracy and computational efficiency makes this framework suitable for integration into future gene-editing automation platforms.

However, the current model is limited by its reliance on simulated environments and single-species genomic data. Future enhancements should include wet-lab validation, transfer learning across diverse crop genomes, and the integration of additional biological layers such as epigenetic modifications and chromatin accessibility maps. The use of multi-agent coordination and lightweight DRL architectures could also broaden the system's applicability to high-throughput genome editing platforms.

In conclusion, this work marks a significant advancement in intelligent CRISPR design, introducing a scalable, learning-based solution for multi-gene editing that bridges computational modeling with agricultural genomics. It lays the groundwork for future systems capable of autonomously navigating complex genomic landscapes, paving the way for accelerated crop engineering in a data-driven era.

Author Contributions: Dileep M R and Syeda Meraj collaboratively developed the research study. Dileep M R led the conceptualization and design of the reinforcement learning framework, performed the implementation, and conducted model training and evaluation. Syeda Meraj contributed to the biological modeling of CRISPR-Cas9 gene editing pathways, curated relevant genomic datasets, and assisted in interpreting the results within the context of crop genomics. Both authors were involved in drafting the

manuscript, critically revising its content for intellectual merit, and approved the final version for submission.

Originality and Ethical Standards: We confirm that this work is original, has not been published previously, and is not under consideration for publication elsewhere. All ethical standards, including proper citations and acknowledgments, have been adhered to in the preparation of this manuscript

Data availability: Data available upon request.

Conflict of Interest: There is no conflict of Interest.

Ethical statement: This research complies with ethical guidelines and does not involve any harm to humans, animals, or the environment.

Funding: The research received no external funding.

Similarity checked: Yes.

References

- [1] A. R. Fernie and J. Yan, "De novo domestication: An alternative route toward new crops for the future," *Molecular Plant*, vol. 12, no. 5, pp. 615–631, May 2019, doi: 10.1016/j.molp.2019.03.003.
- [2] I. Zafar, A. Rafique, J. Fazal, M. Manzoor, Q. U. Ain, and R. A. Rayan, "Genome and gene editing by artificial intelligence programs," in *Adv. AI Techniques Appl. Bioinf.*, Boca Raton, FL, USA: CRC Press, 2021, pp. 165–188.
- [3] I. Zafar, A. Rafique, J. Fazal, M. Manzoor, Q. U. Ain, and R. A. Rayan, "Genome and Gene Editing by Artificial Intelligence," in *Adv. AI Techniques Appl. Bioinf.*, vol. 8, pp. 165, 2021.
- [4] S. Ahmar, P. Ballesta, M. Ali, and F. Mora-Poblete, "Achievements and challenges of genomics-assisted breeding in forest trees: From marker-assisted selection to genome editing," *Int. J. Mol. Sci.*, vol. 22, no. 19, p. 10583, 2021.
- [5] H. Chuai et al., "DeepCRISPR: Optimized CRISPR guide RNA design by deep learning," *Genome Biology*, vol. 19, no. 1, pp. 1–18, Mar. 2018, doi: 10.1186/s13059-018-1459-4.
- [6] G. Naseri and M. A. Koffas, "Application of combinatorial optimization strategies in synthetic biology," *Nat. Commun.*, vol. 11, no. 1, p. 2446, 2020.
- [7] S. Patil, S. Joshi, M. Jamla, X. Zhou, M. J. Taherzadeh, P. Suprasanna, and V. Kumar, "MicroRNA-mediated bioengineering for climate-resilience in crops," *Bioengineered*, vol. 12, no. 2, pp. 10430–10456, 2021.
- [8] P. Tiwari, T. Khare, V. Shriram, H. Bae, and V. Kumar, "Plant synthetic biology for producing potent phyto-antimicrobials to combat antimicrobial resistance," *Biotechnol. Adv.*, vol. 48, p. 107729, 2021.
- [9] L. Wang and J. Zhang, "Prediction of sgRNA on-target activity in bacteria by deep learning," *BMC Bioinf.*, vol. 20, pp. 1–14, 2019.
- [10] M. S. Lakshmi, K. J. Kashyap, S. M. Fazal Khan, N. J. S. Vrata Reddy, and V. B. Kumar Achari, "Whale Optimization based Deep Residual Learning Network for Early Rice Disease Prediction in IoT," *ICST Transactions on Scalable Information Systems*, Oct. 2023, doi: 10.4108/eetsis.4056.
- [11] B. Singh, A. Kumar, A. K. Saini, R. V. Saini, R. Thakur, S. A. Mohammed, and S. Haque, "Strengthening microbial cell factories for efficient production of bioactive molecules," *Biotechnol. Genet. Eng. Rev.*, vol. 39, no. 2, pp. 1345–1378, 2023.
- [12] D. Singh, R. Singh, and A. K. Sharma, "AI-based predictive modeling in CRISPR/Cas-mediated crop trait improvement," *IEEE Access*, vol. 9, pp. 114223–114234, Aug. 2021, doi: 10.1109/ACCESS.2021.3104189.
- [13] M. M. Usaj, E. B. Styles, A. J. Verster, H. Friesen, C. Boone, and B. J. Andrews, "High-content screening for quantitative cell biology," *Trends Cell Biol.*, vol. 26, no. 8, pp. 598–611, 2016.
- [14] S. Xu, S. Gao, and Y. An, "Research progress of engineering microbial cell factories for pigment production," *Biotechnol. Adv.*, vol. 65, p. 108150, 2023.
- [15] A. Ramesh, "Development of CRISPR Based Synthetic Biology Tools for Genome Engineering and Functional Genomic Screening in the Industrially Relevant Oleaginous Yeast *Yarrowia lipoytica*," Ph.D. dissertation, Univ. California, Riverside, CA, USA, 2022.
- [16] H. Su and J. Lin, "Biosynthesis pathways of expanding carbon chains for producing advanced biofuels," *Biotechnol. Biofuels Bioprod.*, vol. 16, no. 1, p. 109, 2023.
- [17] Z. Dai, S. Zhang, Q. Yang, W. Zhang, X. Qian, W. Dong, and F. Xin, "Genetic tool development and systemic regulation in biosynthetic technology," *Biotechnol. Biofuels*, vol. 11, pp. 1–12, 2018.
- [18] G. Goyal, "Study of *E. coli* metabolic pathways for efficient production of commodity chemicals using synthetic biology and genome engineering," Unpublished Thesis, 2021.
- [19] Y. Lei, L. Lu, M. Liu, S. Wang, and J. Li, "CRISPR-P 2.0: an improved CRISPR-Cas9 tool for genome editing in plants," *Mol. Plant*, vol. 7, no. 9, pp. 1494–1496, Sept. 2014. [Online]. Available: <http://crispr.hzau.edu.cn/CRISPR2/>