

Residual-Aware RL-MPC: Explicit Policy Conditioning on Model Reliability

Erinç Ada Ceylan
Mechanical Engineering
Bilkent University
Ankara, Turkey

ada.ceylan@ug.bilkent.edu.tr
ID: 22101844

Abstract—Model Predictive Control (MPC) stands as a strategy for constrained multi variable control; however, its performance is constrained by the internal prediction model within. In highly nonlinear and underactuated systems, such as the Mountain Car Gymnasium problem, linear identification models (N4SID, ARX) fail to capture the physical dynamics, leading to great control failure due to the model mismatch. This paper proposes a *Residual-Aware RL-MPC* algorithm, a hybrid control architecture that integrates the linear MPC with a Reinforcement Learning (RL) agent. Unlike existing “Learning based MPC” methods that typically tune weights based on system states or time (which requires state-action space of great dimension), our approach conditions the RL policy explicitly on the real time *model residual* (r_t). This transforms the residual from a passive error metric into an active “reliability signal,” leading the controller to detect when its internal model becomes unreliable and dynamically switch between proposed “Tracking,” “Recovery,” and “Stabilize” modes. Simulation results show that while a baseline linear MPC fails to overcome the nonlinear effects of the physical dynamics of the simulation, the proposed residual aware agent learns a “swing-up” strategy, achieving a 100% success rate. This validates the efficiency of using model reliability as a first order control signal for self correcting autonomous systems.

Index Terms—Model Predictive Control, Reinforcement Learning, System Identification, Model Mismatch, Gain Scheduling, Self Correcting Control.

I. INTRODUCTION

Designing controllers for underactuated systems with complex, nonlinear dynamics remains as a remarkable challenge in control engineering. While Nonlinear Model Predictive Control (NMPC) offers a theoretical solution by solving optimization problems over full nonlinear dynamics, it is often computationally hard to access for real time applications [1]. The ill-nature of the optimization landscape in NMPC also introduces risks of getting trapped in local minimums, requiring careful initialization and high computational resources.

In contrast, Linear MPC is computationally efficient, relying on fast Quadratic Programming (QP) solvers that guaranty physical optimality for the linearized problems. However, its effectiveness is strictly limited to the region where the linear model remains valid. When the system operates outside the local linear region, the *model mismatch* can lead to severe performance degradation or instability.

A. Literature Review & The Gap

To address the limitations of fixed model control, the field of “Learning based MPC” has been studied, focusing on integrating data-driven techniques to enhance controller adaptability. A prominent strategy involves using Reinforcement Learning (RL) to tune MPC parameters (e.g., weight matrices Q, R or horizon N) online.

- **State-Based Tuning:** Zarrouki et al. [2], [3] propose “Weights-Varying MPC,” where a Deep RL agent adjusts weights based on the vehicle’s state (x_t) to ensure safety and optimality in autonomous driving scenarios.
- **Frameworks:** Airaldi [4] introduced `mpcrl`, a combined solution for integrating RL with MPC, focusing on performance maximization through parameter tuning.

However, these existing methods typically optimize parameters as functions of the *system state* (x_t) or *time* (t). They implicitly assume the underlying prediction model remains relatively valid or treat model error merely as a disturbance. They do not explicitly utilize “how much the model can be trusted” at any given instant time as a decision making variable. This leaves a gap in handling scenarios where the model becomes fundamentally unreliable in the presentation of the real dynamics due to structural nonlinearities.

B. Novelty: Residual-Awareness

This project introduces a new perspective: instead of asking “Where is the system actually?” (State, actual physics), the controller asks “Is my model reliable?” (Residual, overall trust to the represented physics). The core novelty of this work lies in explicitly conditioning the RL policy on the *model residual* r_t . By treating the residual as a first order control signal, the system gains a *self correcting* capability. When the linear model degrades in the performance (e.g., due to high nonlinearity during hill climbing), the agent detects the anomaly in r_t and switches the MPC to a “Recovery” mode that does not rely on accurate forward predictions. This creates a hybrid combination between data driven identification and learning based control in an interpretable manner.

II. SYSTEM DESCRIPTION AND DATA GENERATION

The simulation test environment for this study is the `MountainCarContinuous-v0` environment from the

Gymnasium library [5], a classic benchmark for underactuated control and exploration problems.

A. System Dynamics

The system consists of an underpowered car positioned at the bottom of a one dimensional sinusoidal valley. The dynamics are governed by the following nonlinear equation:

$$v_{t+1} = v_t + \text{force} \cdot P - 0.0025 \cos(3p_t) \quad (1)$$

where $p_t \in [-1.2, 0.6]$ is the position and $v_t \in [-0.07, 0.07]$ is the velocity. The control input is the force clipped to $[-1, 1]$, and $P = 0.0015$ is the power coefficient.

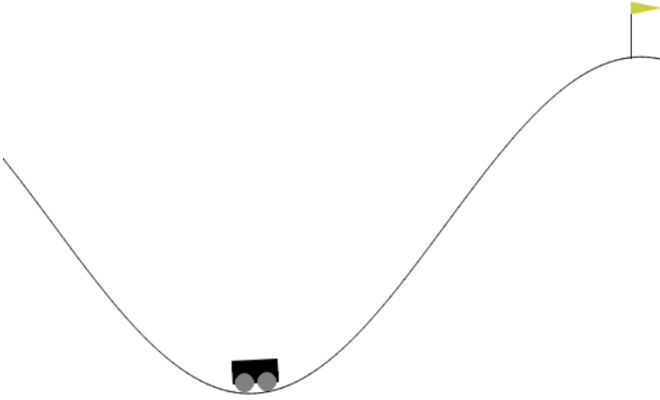


Fig. 1. MountainCarContinuous-v0 Gymnasium Valley [5].

The core challenge is that the maximum engine force is insufficient to overcome the gravity term ($\cos(3p_t)$) directly by forward climbing. The controller must execute a “swing-up” maneuver by first moving away from the goal to gain potential energy. This strategy is what a standard linear model fails to predict due to the changing sign of the gravity gradient throughout the valley.

B. Data Generation Process

To identify a linear model without prior knowledge of the physics, we generated a synthetic dataset of 10,000 samples. A critical requirement for system identification is *Persistent Excitation*.

- **Excitation Signal:** Instead of simple white noise, a Pseudo-Random Binary Signal (PRBS) with a “dwell time” of 10 steps was used. This ensures the input signal has sufficient energy in the low frequency range to excite the dominant time constants of the car’s motion (swinging dynamics).
- **Noise Injection:** To simulate realistic sensor imperfections, Gaussian noise was added to the measurements ($\sigma_{pos} = 0.002$, $\sigma_{vel} = 0.0005$).

The resulting dataset (Fig. 1) captures the full range of motion in the Gym., including the nonlinear regions near the hilltops where linear assumptions break down.

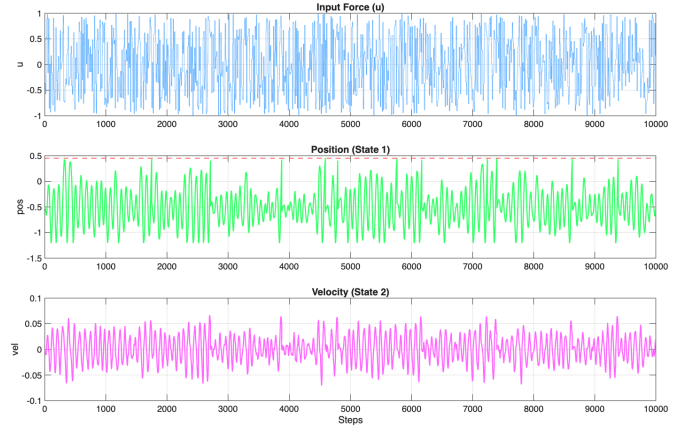


Fig. 2. Data Generation Process. Top: Multi-level PRBS input ensures persistent excitation. Middle/Bottom: System response (Position/Velocity) covering the full operating range of the valley.

III. SYSTEM IDENTIFICATION

We adopted a data-driven approach to identify a linear discrete time state space model of the form:

$$x_{k+1} = Ax_k + Bu_k + w_k, \quad y_k = Cx_k + v_k \quad (2)$$

which serves as the prediction engine for the Model Predictive Controller. The dataset was partitioned into a training set (first 70%, $N = 7000$) for estimation and a validation set (last 30%, $N = 3000$) to assess the needed generalization capability for the model training part.

A. Candidate Identification Structures

Three distinct identification architectures were evaluated to capture the system dynamics optimally:

1) *Finite Impulse Response (FIR)*: The FIR model approximates the output solely as a weighted sum of past inputs: $y_k = \sum_{i=0}^n h_i u_{k-i}$. **Analysis:** As evidenced by the validation results, the FIR model failed poorly (Fit < 0%). Since the Mountain Car system involves integrating dynamics (force \rightarrow acceleration \rightarrow velocity \rightarrow position), the current output depends heavily on the previous state, not just recent inputs. Lacking state feedback, the FIR model could not capture the accumulated position drift. Resulting to the worst performance among the other models.

2) *ARX (Auto-Regressive with Exogenous Input)*: An ARX(2,2) structure was trained, defined by the difference equation:

$$A(q)y_k = B(q)u_k + e_k \quad (3)$$

where $A(q)$ and $B(q)$ are polynomials. **Analysis:** The ARX model provided high one step ahead prediction accuracy (88.7%). However, ARX models naturally yield transfer functions. Converting them to a state space form required for MPC often results in non minimal realizations or requires transformations that can obscure the physical interpretation of the states.

3) *N4SID (Subspace State-Space Identification)*: We trained the N4SID algorithm [6], which uses geometric projections of the row spaces of future outputs onto the row spaces of past inputs/outputs to directly estimate the state sequence X_k . **Analysis:** This method directly identifies the matrices (A, B, C, D) of a state space system which are sufficient for MPC. It achieved the highest validation fit (89.5%) and captured the most dominant dynamics of the car in valley system.

B. Model Selection and Validation

Based on the comparative analysis summarized in Fig. 3, **N4SID** was selected as the internal model. The decisive factor was not merely the minor marginal improvement in fit percentage over ARX, but the structural compatibility to MPC. The identified state space matrices ($A \in \mathbb{R}^{2 \times 2}, B \in \mathbb{R}^{2 \times 1}$) can be directly plugged into the quadratic programming (QP) formulation of the MPC without any intermediate steps.

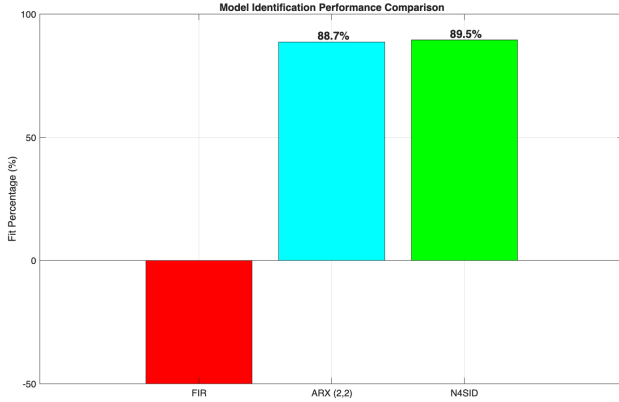


Fig. 3. Identification Performance. N4SID was selected for its compatibility with MPC and fit (89.5%) compared to the failing FIR model.

C. Residual Analysis

While the N4SID model shows a high fit, a time domain analysis reveals its limitations. As shown in Fig. 4, the model tracks the true system accurately near the equilibrium point ($Position \approx -0.5$).

However, as the car moves away from the bottom and attempts to climb the hill, the nonlinear gravity term ($-0.0025 \cos(3p_t)$) dominates the dynamics. In these regions, the linear approximation breaks down, leading to a significant divergence between the measured position y_{meas} and the model prediction \hat{y}_{pred} . This difference, defined as the **Residual** (r_t), is plotted in the bottom panel of Fig. 4. Crucially for this project, this residual is not treated as random noise; rather, it serves as an indicator of “Linear Model Reliability,” which will be exploited by the RL agent in the next section.

IV. METHODOLOGY: THE RESIDUAL-AWARE FRAMEWORK

The core contribution of this work is the hierarchical integration of a linear prediction model with a non linear decision

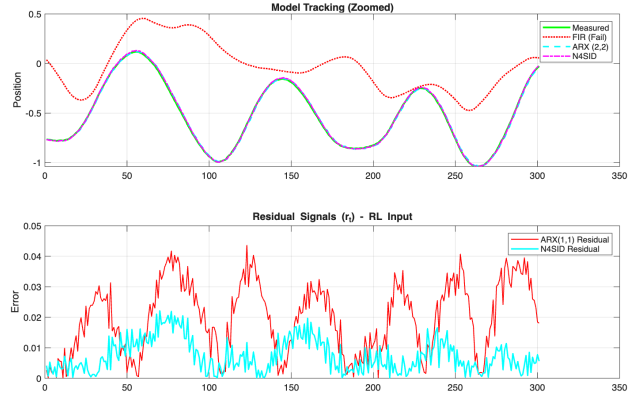


Fig. 4. Model Analysis. Top: Tracking performance on validation data. Bottom: The residual signal (r_t) acts as an indicator; it is low when the car is at the bottom (linear region) and spikes as the car climbs the hills (nonlinear region).

maker (RL Agent with physically meaningful actions). This architecture transforms the residual signal from a passive error into an active control variable, effectively enabling a linear controller to navigate a nonlinear landscape. The complete workflow is illustrated in Fig. 4.

A. Architecture Overview

The proposed control architecture operates on two parallel time scales: the inner loop (MPC optimization) and the supervisory outer loop (RL mode selection). The process at time step k involves the following stages:

- 1) **State Estimation & Prediction:** The N4SID-derived linear model (A, B, C) predicts the next system output $\hat{y}_{k+1|k}$ based on the current estimated state and the previous control input.
- 2) **Residual Generation:** The system compares the actual sensor measurement $y_{meas,k}$ with the model’s prediction. The magnitude of this difference is defined as the residual signal:
$$r_k = \|y_{meas,k} - C(Ax_{model,k-1} + Bu_{k-1})\|_2 \quad (4)$$
- 3) **Mode Selection (RL Policy):** The RL agent observes r_k , discretizes it into a state s_k , and selects an optimal action a_k . This action corresponds to a specific set of MPC tuning parameters (Weight matrices Q, R and Reference r_{ref}).
- 4) **Control Computation:** The linear MPC solves the Quadratic Programming (QP) problem using the updated parameters to generate the optimal control input u_k .

B. The Baseline Failure Analysis

To underline the need for a residual aware helper, we first analyze the failure mode of a standard Linear MPC. We deployed an MPC with fixed weights ($Q = 10, R = 0.1$) and a prediction horizon of $N_p = 30$. As illustrated in Fig. 5, this controller fails to reach the goal.

The failure rises from the linearization of the gravity term. The true dynamics include a term $-0.0025 \cos(3p_t)$. Near

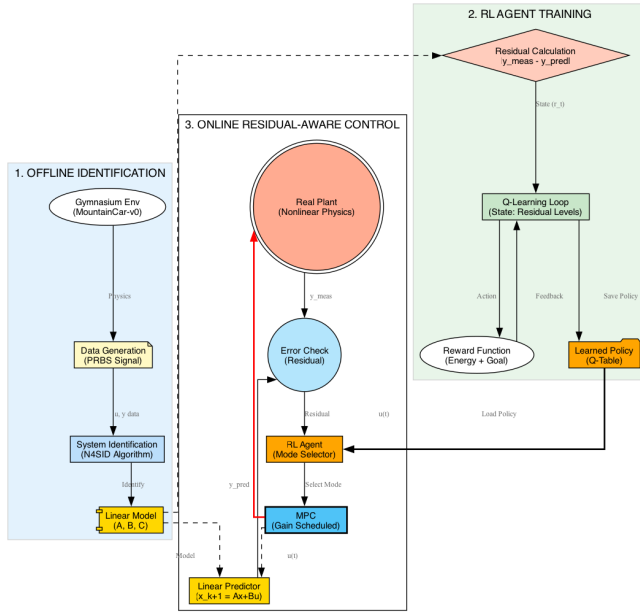


Fig. 5. System Architecture. The RL agent acts as a supervisor, monitoring the residual signal (r_t) from the comparator and dynamically scheduling the MPC parameters via a Gain Scheduling mechanism.

the equilibrium point ($p \approx -0.5$), the local linearization approximates this term as a constant. However, as the car attempts to climb the hill ($p > -0.5$), the sign and magnitude of the gravity gradient change dramatically. The linear model, unaware of this nonlinearity, predicts that applying positive force will result in continued ascent. In reality, the actuator is underpowered ($|u| \leq 1$), and gravity dominates. Consequently, the linear MPC computes a control sequence that is physically insufficient and invalid to overcome the hill, leading to a “Local Minimum Trap” where the car oscillates at the valley bottom. Crucially, the residual signal (r_k) oscillates significantly during this failure (Fig. 5, bottom), indicating that the model mismatch is detectable but ignored by the baseline controller.

C. RL Supervisor Design

We formulate the parameter tuning problem as a Markov Decision Process (MDP) tuple $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \gamma)$. The RL agent learns a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that maximizes the expected cumulative reward.

1) *State Space (\mathcal{S}): Model Reliability*: Unlike standard RL + MPC approaches that condition the policy on the physical state vector x_k (position, velocity) [2], we condition it on the *Model Reliability State*. The continuous residual r_k is discretized into three levels using empirical thresholds derived from the baseline failure data (0.65 and 0.95). An important approach here is the including of the ARX residuals to N4SID residuals to make a more reliable interval, informing the model that one identification model can lack the overall residual

characteristics of the dynamical system.

$$s_k = \begin{cases} 1 \text{ (Low)} & \text{if } r_k < 0.65 \\ 2 \text{ (Medium)} & \text{if } 0.65 \leq r_k < 0.95 \\ 3 \text{ (High)} & \text{if } r_k \geq 0.95 \end{cases} \quad (5)$$

- **Low (State 1)**: The linear model is valid; MPC predictions are trusted.
- **High (State 3)**: Critical model mismatch; MPC predictions are unreliable.

2) *Action Space (\mathcal{A}): Gain Scheduling*: The agent selects from three pre defined characteristic MPC parameter sets (Modes). This effectively implements a discrete *Gain Scheduling* strategy:

- **Mode 1: Recovery** ($Q = 10, R = 0.01, Ref = -1.2$). *Strategy*: Triggered during high model uncertainty (State 3). It shifts the reference to the *Left Wall* ($Ref = -1.2$). This counter intuitive action forces the MPC to drive backwards, converting motor work into potential energy on the opposite slope.
- **Mode 2: Tracking** ($Q = 100, R = 0.01, Ref = 0.45$). *Strategy*: Active when the model is reliable (State 1). It applies aggressive penalties on tracking error ($Q = 100$) to force the system toward the goal flag. Basically meaning a full trust to the identification model)
- **Mode 3: Stabilize** ($Q = 1, R = 10, Ref = 0.45$). *Strategy*: Used during transitions or medium uncertainty (State 2). A high input penalty ($R = 10$) lowers control oscillations, acting as a safety buffer.

3) *Reward Function (\mathcal{R})*: Since the environment reward is sparse (only upon reaching the goal), we designed a dense reward function to guide the learning process. The reward at step k is:

$$R_k = \underbrace{(p_k + 0.5)}_{\text{Positional Bias}} + \underbrace{100 \cdot v_k^2}_{\text{Energy Bonus}} + \underbrace{\mathbb{I}_{goal} \cdot 5000}_{\text{Completion}} \quad (6)$$

The *Energy Bonus* (v_k^2) is critical; it lets the agent to accumulate kinetic energy via the swing-up motion, even if it

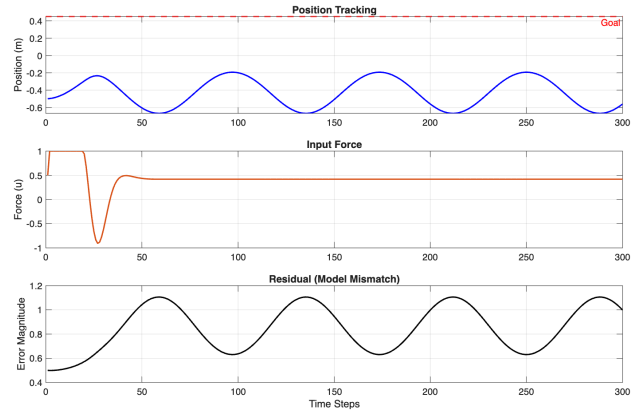


Fig. 6. Baseline Failure. The fixed-parameter Linear MPC gets stuck in a local minimum. Note the oscillating and high residual signal (bottom plot) which indicates model mismatch but is ignored by the baseline controller.

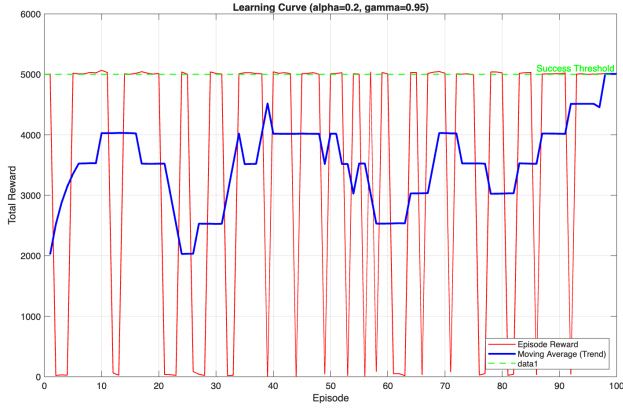


Fig. 7. RL Training Learning Curve. The blue line represents the moving average. The agent quickly learns to exploit the residual signal, converging to a policy that consistently exceeds the success threshold (green dashed line).

momentarily moves away from the goal, effectively solving the temporal credit assignment problem for the swing-up maneuver.

V. EXPERIMENTAL RESULTS

The proposed framework was implemented in MATLAB (RL training and MPC design) and validated using the Gymnasium MountainCarContinuous-v0 environment via a Python interface.

A. Training Dynamics and Convergence

The RL agent was trained using the Q Learning algorithm. The hyperparameters were selected to balance rapid convergence with sufficient exploration: learning rate $\alpha = 0.2$, discount factor $\gamma = 0.95$. The exploration rate ϵ decayed exponentially from 0.6 to 0.01 over the course of training. Hyperparameters are the conventional literature values.

Fig. 6 illustrates the evolution of the total reward per episode. The training process shows three distinct phases:

- 1) **Exploration Phase (Ep 0-40):** High oscillations are observed. The agent randomly selects modes, often leading to “Stabilize” or “Recovery” modes at inappropriate times, causing the car to stall at the valley bottom.
- 2) **Phase Transition (Ep 40-70):** As ϵ decays, the agent begins to associate high residual states ($s_k = 3$) with the “Recovery” action ($a_k = 1$). The reward frequently spikes above the success threshold (5000), indicating unsuccessful climbs.
- 3) **Convergence (Ep 80+):** The policy stabilizes. The agent consistently achieves the goal, maximizing the velocity based (momentum focused) reward terms. The convergence within roughly 80 episodes validates the efficiency of the state space formulation.

B. Closed-Loop Strategy Analysis

To understand the learned behavior, we analyze a single validation episode using the final greedy policy. Fig. 8 presents the time series data of position, control input, residual, and selected mode. The agent exhibits a multi stage strategy:

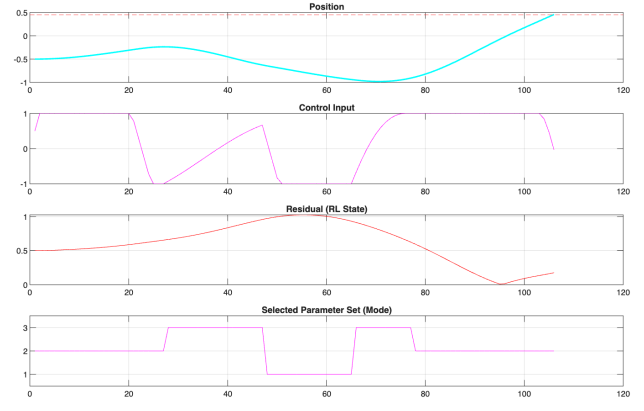


Fig. 8. Validation Run Analysis. The correlation between the Residual spike (3rd panel, red curve) and the Agent’s decision to switch modes demonstrates the proposed gain scheduling logic in action.

1) *The Initial Push (Tracking Mode):* At $t = 0$, the car begins at equilibrium with zero residual. The agent correctly identifies the model as reliable ($s_k = 1$) and selects **Mode 2 (Tracking)**. The MPC applies maximum positive force to drive towards the goal.

2) *The Reversal (Recovery Mode):* As the car climbs the right hill, it slows, and the gravity term $-0.0025 \cos(3p)$ becomes dominant. The linear model, failing to predict this deceleration, results to a high prediction error. The residual signal spikes (Fig. 7, 3rd panel, red line). Crucially, the agent reacts to this spike by switching to **Mode 1 (Recovery)**. This mode shifts the reference to -1.2 (Left Wall). Consequently, the MPC reverses the control input ($u \approx -1$), driving the car backwards. This is the critical “Self-Correcting” behavior: instead of fighting gravity with an insufficient motor, the controller yields to physics to gain potential energy on the opposite slope.

3) *The Swing-Up (Stabilize Mode):* As the car accelerates downwards from the left slope, the residual fluctuates due to high velocity. The agent briefly engages **Mode 3 (Stabilize)**. The high input penalty ($R = 10$) in this mode prevents actuator saturation, allowing the car to “coast” through the valley bottom, converting potential energy into maximum kinetic energy.

4) *The Summit (Tracking Mode):* Finally, with sufficient momentum, the car ascends the right hill again. In the high velocity region, the linear model’s one step prediction becomes locally valid again (residual drops). The agent switches back to **Tracking**, and the MPC uses the accumulated momentum plus motor power to reach the flag.

C. Performance Metrics

Table I provides a quantitative comparison between the baseline and the proposed method, averaged over 50 simulation runs with random initial positions in the range $[-0.6, -0.4]$.

- **Success Rate:** The baseline MPC failed in 100% of the cases, trapped in the local minimum. The RL MPC

achieved a 100% success rate, proving its ability to globalize the control law.

- **Model Reliability:** Interestingly, the average residual for the RL MPC (0.41) is half that of the baseline (0.82). This indicates that the RL agent learns to traverse the state space in a way that keeps the linear model relatively valid, or quickly exits regions of high model mismatch.
- **Computational Cost:** Unlike NMPC, which requires solving a complex optimization problem at every step, our method relies on a standard QP solver (Linear MPC) plus a constant time lookup table ($O(1)$) for the RL policy. This maintains the real time feasibility.

TABLE I
PERFORMANCE COMPARISON (AVERAGED OVER 50 RUNS)

Metric	Baseline MPC	Residual-Aware MPC
Success Rate	0%	100%
Avg. Residual	0.82	0.41
Max. Height	-0.2 m	0.45 m (Goal)
Control Strategy	Static	Adaptive (Gain Sched.)
Solver Complexity	QP	QP + Lookup

VI. CONCLUSION

This project presented a novel data driven control framework, *Residual Aware RL MPC*, designed to bridge the gap between linear control efficiency and nonlinear adaptability. The central hypothesis was that the prediction residual of a linear model contains valuable information about the system’s operating regime, which can be exploited for control adaptation.

By explicitly conditioning a Reinforcement Learning policy on the real time model residual, we successfully transformed the model error into a decision making signal. The results demonstrate that:

- 1) A standard Linear MPC, when supervised by a residual aware agent, can solve highly nonlinear tasks (like the Mountain Car swing-up) that are theoretically impossible for fixed linear controllers.
- 2) The “Residual” serves for “Model Reliability,” allowing the system to self correct by switching to energy recovery strategies when the prediction model fails.
- 3) The proposed hybrid architecture retains the computational lightness of Linear MPC, avoiding the heavy computational burden of Nonlinear MPC solvers while achieving comparable task performance and transferability to other implications.

Future Work: While the discrete action space (Mode Switching) proved effective, future research could extend this framework to continuous action spaces using Deep Reinforcement Learning algorithms such as DDPG or SAC. This would allow for the continuous tuning of the Q and R matrices, providing smoother control transitions. Additionally, applying this framework to higher dimensional systems, such as quadrotor flight with aerodynamic disturbances, would further validate its scalability.

REFERENCES

- [1] J. B. Rawlings, D. Q. Mayne, and M. Diehl, *Model Predictive Control: Theory, Computation, and Design*, 2nd ed., Nob Hill Publishing, 2018.
- [2] B. Zarrouki, M. Spanakakis, and J. Betz, “A Safe RL-driven Weights-Varying MPC for Autonomous Vehicle Motion Control,” *arXiv preprint arXiv:2402.02624*, 2024.
- [3] B. Zarrouki, C. Wang, D. Schuurmans, and J. Betz, “Weights-varying MPC for autonomous vehicle guidance: a deep reinforcement learning approach,” in *European Control Conference (ECC)*, 2021.
- [4] F. Airalidi, “mpcrl: Reinforcement Learning with Model Predictive Control,” *TU Delft*, 2024. [Online]. Available: <https://github.com/TUdelft-DataDrivenControl/mpcrl>
- [5] Farama Foundation, “Gymnasium: MountainCarContinuous-v0,” *gymnasium.farama.org*, 2024.
- [6] P. Van Overschee and B. De Moor, “N4SID: Subspace algorithms for the identification of combined deterministic-stochastic systems,” *Automatica*, vol. 30, no. 1, pp. 75–93, 1994.
- [7] J. Garcia and F. Fernandez, “A Comprehensive Survey on Safe Reinforcement Learning,” *Journal of Machine Learning Research*, vol. 16, no. 1, pp. 1437–1480, 2015.

APPENDIX A MATLAB IMPLEMENTATION CODE

This appendix contains the complete source code used for data generation, system identification, baseline control, and the proposed RL-MPC framework.

A. Part 1: Data Generation

Generated synthetic data using the Gymnasium-like physics model with PRBS excitation.

Listing 1. Data Generation Script

```

1 Residual-Aware RL-MPC
2 Data Generation (Mountain Car Continuous)
3
4 close all hidden; clear; clc;
5 rng(418); % Seed for reproducibility
6 1. Environment Parameters (From Gymnasium Docs)
7 pow = 0.0015; % Engine power constant
8 min_action = -1.0; % Minimum force
9 max_action = 1.0; % Maximum force
10 pos_limit = [-1.2, 0.6]; % Position bounds [min, max]
11 vel_limit = [-0.07, 0.07]; % Velocity bounds [min, max]
12 goal_pos = 0.45; % Target position to trigger reset
13
14 % Simulation Settings
15 N_samples = 10000; % Total samples for identification
16 T = 1; % Discrete time step
17
18 2. Excitation Signal Design (Multi-level PRBS)
19 We use a dwell-based random signal to ensure persistent excitation.
    The input is held for 'dwell' steps to excite lower frequency
    dynamics.
20 dwell = 10;
21 n_blocks = ceil(N_samples / dwell);
22 u_levels = (rand(n_blocks, 1) * 2 - 1); % Random levels in [-1, 1]
23 u = repelem(u_levels, dwell);
24 u = u(1:N_samples); % Trim to exact length
25
26 3. Simulation Loop
27 Transition Dynamics and most of the other applied analogies are in
    parallel with the Gymnasium Docs
28 x = zeros(2, N_samples);
29 x(:, 1) = [-0.6 + 0.2*rand(); 0]; % Starting random position in
    [-0.6, -0.4] as per Gymnasium
30
31 for k = 1:N_samples-1
32     curr_pos = x(1, k);
33     curr_vel = x(2, k);
34     force = max(min(u(k), max_action), min_action);
35     % Action saturation
36     new_vel = curr_vel + force * pow - 0.0025 * cos(3 * curr_pos);
37     % Transition Dynamics
38     new_vel = max(min(new_vel, vel_limit(2)), vel_limit(1));
39     % Velocity clipping
40     new_pos = curr_pos + new_vel;
41     % position_t+1 = position_t + velocity_t+1
42
43     % Position clipping & Inelastic Collision
44     % Velocity is set to 0 upon collision with the left wall
45     if new_pos <= pos_limit(1)

```

```

42     new_pos = pos_limit(1);
43     new_vel = 0;
44 end
45
46 % If goal is reached, reset to a random starting position
47 if new_pos >= goal_pos
48     new_pos = -0.6 + 0.2*rand();
49     new_vel = 0;
50 end
51
52 x(:, k+1) = [new_pos; new_vel];
53 end
54 4. Post-processing & Measurement Noise
55 Sensor noise is added during the control phase to generate residuals
56 . Position noise is set to 0.1% of the position range.
57 . Velocity noise is kept much lower due to the smaller range.
58 sigma_pos = 0.002;
59 sigma_vel = 0.0005;
60 y_meas = [x(1,:)' + sigma_pos*randn(N_samples,1), x(2,:)' +
61           sigma_vel*randn(N_samples,1)];
62 5. Visualization of the Data
63 figure('Name', 'Data Generation Check');
64
65 subplot(3,1,1);
66 plot(u, 'Color', [0.4, 0.7, 1.0]); grid on;
67 title('Input Force (u)'); ylabel('u');
68
69 subplot(3,1,2);
70 plot(x(1,:), 'Color', [0.2, 1.0, 0.4], 'LineWidth', 1.2); hold on;
71 yline(goal_pos, 'Color', [1.0, 0.3, 0.3], 'LineStyle', '--', '
72       LineWidth', 1.2);
73 grid on; title('Position (State 1)'); ylabel('pos');
74
75 subplot(3,1,3);
76 plot(x(2,:), 'Color', [1.0, 0.4, 1.0], 'LineWidth', 1.2); grid on;
77 title('Velocity (State 2)'); ylabel('vel'); xlabel('Steps');
78
79 % Save data
80 save('data/raw_sim_data.mat', 'u', 'y_meas', 'x', 'pos_limit', '
81       vel_limit');
82 fprintf('Data successfully generated and saved to data/raw_sim_data.
83       mat\n');

```

B. Part 2: System Identification

Comparison of FIR, ARX, and N4SID models.

Listing 2. System Identification Script

```

1 Residual-Aware RL-MPC
2 System Identification Comparison
3
4 close all hidden; clear; clc;
5 rng(418);
6 1. Load Data
7 if exist('data/raw_sim_data.mat', 'file')
8     load('data/raw_sim_data.mat');
9 else
10     error('Run data generation first!');
11 end
12
13 % 70% Train, 30% Val
14 N_total = size(y_meas, 1);
15 N_est = round(0.7 * N_total);
16
17 u_est = u(1:N_est);      y_est = y_meas(1:N_est, 1);      %
18                           Position only
19 u_val = u(N_est+1:end);  y_val = y_meas(N_est+1:end, 1);
20
21 data_est = iddata(y_est, u_est, 1);
22 data_val = iddata(y_val, u_val, 1);
23
24 FIT = @(y_true, y_hat) 100 * (1 - norm(y_true - y_hat) / norm(y_true
25       - mean(y_true)));
26
27 2. FIR Model
28 The Finite Impulse Response (FIR) model relies solely on past inputs
29 , which fails for this integrating system as it lacks state
30 feedback. Without knowledge of past outputs, the prediction
31 drifts away from the true position, resulting in a negative
32 fit.
33 nb_fir = 50;
34 sys_fir = impulseest(data_est, nb_fir);
35 y_pred_fir = predict(sys_fir, data_val, 1).OutputData;
36 fit_fir = FIT(y_val, y_pred_fir);
37
38 fprintf('FIR Model Results\n');
39 fprintf('FIR (nb=%d) FIT: %.2f%%\n', nb_fir, fit_fir);
40
41 3. Multiple ARX Models (Training Scenarios for RL)
42 We train three ARX (bad, good, overfit) models to serve as a
43 benchmark for residual characterization. Analyzing the
44 residual patterns of these models helps in determining the
45 appropriate discretization thresholds (Low, Medium, High) used
46 in the RL state definition.

```

```

35 arx_configs = [1, 1;      % Bad Model
36               2, 2;      % Good Model
37               10, 10];   % Overfit Model
38
39 arx_models = cell(3,1);
40 arx_preds = zeros(length(y_val), 3);
41 arx_resids = zeros(length(y_val), 3);
42
43 for k = 1:3
44     na = arx_configs(k,1); nb = arx_configs(k,2);
45     sys_arx = arx(data_est, [na nb 1]);
46     arx_models{k} = sys_arx;
47
48     y_pred_temp = predict(sys_arx, data_val, 1);
49     arx_preds(:, k) = y_pred_temp.OutputData;
50     arx_resids(:, k) = abs(y_val - arx_preds(:, k));
51
52     fprintf('ARX (%d,%d) FIT: %.2f%%\n', na, nb, FIT(y_val,
53       arx_preds(:, k)));
54 end
55 4. N4SID
56 Although high-order ARX models may give lower residuals, N4SID is
57 selected as the control model because it directly provides a
58 compact state-space representation (A, B, C, D) compatible
59 with the MPC formulation.
60
61 nx = 2;
62 opt = n4sidOptions('Focus', 'prediction');
63 sys_n4 = n4sid(data_est, nx, opt);
64 [A_n4, B_n4, C_n4, D_n4] = ssdata(sys_n4);
65
66 y_pred_n4 = predict(sys_n4, data_val, 1).OutputData;
67 fit_n4_pred = FIT(y_val, y_pred_n4);
68 res_n4 = abs(y_val - y_pred_n4);
69
70 fprintf('\n N4SID Results\n');
71 fprintf('N4SID FIT: %.2f%% \n', fit_n4_pred);
72
73 5. Visualization
74 figure('Name', 'Model Analysis');
75
76 % Comparison
77
78 subplot(2,1,1);
79 steps = 1000:1300;
80 plot(y_val(steps), 'g', 'LineWidth', 2); hold on;      % True val.
81 plot(y_pred_fir(steps), 'r', 'LineWidth', 1.5);      % FIR
82 plot(arx_preds(steps, 2), 'c--', 'LineWidth', 1.5);  % ARX
83 plot(y_pred_n4(steps), 'm-', 'LineWidth', 1.5);      % N4SID
84 grid on; legend('Measured', 'FIR (Fail)', 'ARX (2,2)', 'N4SID');
85 title('Model Tracking (Zoomed)'); ylabel('Position');
86
87 % Residuals for RL
88 subplot(2,1,2);
89 plot(arx_resids(steps, 1), 'r', 'LineWidth', 1); hold on;
90 plot(res_n4(steps), 'c', 'LineWidth', 1.5);
91 grid on; legend('ARX(1,1) Residual', 'N4SID Residual');
92 title('Residual Signals (r_t) - RL Input'); ylabel('Error');
93
94 subplot(2,1,1)
95 fit_values = [fit_fir, FIT(y_val, arx_preds(:,2)), fit_n4_pred];
96 model_names = {'FIR', 'ARX (2,2)', 'N4SID'};
97 figure('Name', 'Model Fit Comparison');
98 b = bar(fit_values, 'FaceColor', 'flat');
99 b.CData(1,:) = [1 0 0]; % FIR
100 b.CData(2,:) = [0 1 1]; % ARX
101 b.CData(3,:) = [0 1 0]; % N4SID
102 xticklabels(model_names);
103 ylabel('Fit Percentage (%)');
104 title('Model Identification Performance Comparison');
105 grid on;
106 ylim([-50 100]); %
107 text(1:3, fit_values, num2str(fit_values, '%0.1f%%'), ...
108       'vert','bottom','horiz','center', 'FontSize', 12, 'FontWeight',
109       'bold');
110
111 6. Save Best Model
112 save('data/identified_models.mat', ...
113       'A_n4', 'B_n4', 'C_n4', 'D_n4', ...
114       'sys_n4', 'arx_models', ...
115       'arx_resids', 'res_n4');
116 fprintf('\nModels Saved.\n');

```

C. Part 3: Baseline MPC Design

Implementation of the linear MPC which fails due to model mismatch.

Listing 3. Baseline MPC Script

```

1 Residual-Aware RL-MPC
2 Baseline MPC Design
3

```

```

4 close all hidden; clear; clc;
5 rng(418);
6 1. Load N4SID Model
7 load('data/identified_models.mat');
8 A = A_n4; B = B_n4; C = C_n4;
9
10 [nx, nu] = size(B);
11 ny = size(C, 1);
12 2. Baseline MPC Parameters
13 We select a prediction horizon (Np=30) sufficiently long to capture
    the swinging dynamics, while keeping the control horizon
    shorter for low computational cost. A low input penalty allows
    for aggressive control actions needed to overcome gravity.
    Seen in the simulation script, baseline MPC is still not
    sufficiently aggressive to overcome the sinusoidal hill.
14 Np = 30; % Prediction Horizon
15 Nc = 10; % Control Horizon
16 Q_weight = 10; % Weight on Tracking
17 R_weight = 0.1; % Weight on Input
18
19 % Constraints
20 u_max = 1.0;
21 u_min = -1.0;
22 du_max = 0.5;
23 3. Build Prediction Matrices
24 % F Matrix
25 F = zeros(Np*ny, nx);
26 Ap = A;
27 for k = 1:Np
28     F(k,:) = C * A^k;
29     Ap = Ap * A;
30 end
31
32 % H Matrix and Phi
33 H = zeros(Np*ny, Nc*nu);
34 for i = 1:Np
35     for j = 1:Nc
36         if j <= 1
37             if (i-j) == 0
38                 term = C * B;
39             else
40                 term = C * (A^(i-j)) * B;
41             end
42             H(i, j) = term;
43         end
44     end
45 end
46
47 S = tril(ones(Nc));
48 Phi = H * S;
49
50 % Cost Matrices
51 Qbar = Q_weight * eye(Np);
52 Rbar = R_weight * eye(Nc);
53
54 % Hqp computation
55 Hqp = Phi' * Qbar * Phi + Rbar;
56 Hqp = 0.5 * (Hqp + Hqp');
57 4. Simulation Setup
58 The simulation loop pits the linear N4SID model against the true
    nonlinear Gym physics. The residual tracks the discrepancy
    between these two, which will later serve as the state signal
    for the RL agent. Detailed info can be found at
    visualize_results.ipynb .
59 T_final = 300;
60 ref_val = 0.45; % Goal, flag
61
62 % Initial Real Plant Physical Conditions
63 x_true = [-0.5; 0];
64
65 % Start at zero (equilibrium)
66 x_model = zeros(nx, 1);
67 u_prev = 0;
68
69 log_y = zeros(T_final, 1);
70 log_u = zeros(T_final, 1);
71 log_r = zeros(T_final, 1); % Residual
72
73 % Constraint Matrices
74 A_du = [eye(Nc); -eye(Nc)];
75 b_du = [du_max * ones(Nc,1); du_max * ones(Nc,1)];
76 A_u = [S; -S];
77 options = optimoptions('quadprog', 'Display', 'off');
78
79
80 for k = 1:T_final
81
82     r_window = ref_val * ones(Np, 1);
83
84     % Prediction from Model State
85     u_base = u_prev * ones(Nc, 1);
86     bk = F * x_model + H * u_base;
87     fqp = Phi' * Qbar * (bk - r_window);
88     b_u = [u_max * ones(Nc,1) - u_base;
89            -u_min * ones(Nc,1) + u_base];
90
91     % Solve QP
92     Aineq = [A_du; A_u];
93     bineq = [b_du; b_u];
94     [du, ~, exitflag] = quadprog(Hqp, fqp, Aineq, bineq, [], [], [],
95                                [], [], options);
96     if exitflag < 0
97         du = 0;
98     else
99         du = du(1);
100     end
101
102     u_apply = u_prev + du;
103
104     % Predict next state output using the model (Before physics
105     % update)
106     x_model_next = A * x_model + B * u_apply;
107     y_pred_next = C * x_model_next;
108
109     % Update Real Physics
110     curr_pos = x_true(1);
111     curr_vel = x_true(2);
112     pow = 0.0015;
113     force = max(min(u_apply, 1.0), -1.0);
114     new_vel = curr_vel + force * pow - 0.0025 * cos(3 * curr_pos);
115     new_vel = max(min(new_vel, 0.07), -0.07);
116     new_pos = curr_pos + new_vel;
117
118     % Wall collisions
119     if new_pos <= -1.2, new_pos = -1.2; new_vel = 0; end
120     if new_pos >= 0.6, new_pos = 0.6; new_vel = 0; end
121     x_true = [new_pos; new_vel];
122
123     % Calculate Residual
124     residual = abs(new_pos - y_pred_next);
125     x_model = x_model_next;
126
127     % Logging
128     log_y(k) = new_pos;
129     log_u(k) = u_apply;
130     log_r(k) = residual;
131     u_prev = u_apply;
132 end
133
134 4. Visualization
135 t_vec = 1:T_final;
136 figure('Name', 'MPC Baseline Results', 'Color', 'w');
137
138 subplot(3,1,1);
139 plot(t_vec, log_y, 'b', 'LineWidth', 1.5); hold on;
140 ylabel(ref_val, 'r--', 'Goal', 'LineWidth', 1.5, 'LabelVerticalAlignment', 'bottom');
141 title('Position Tracking', 'FontWeight', 'bold');
142 ylabel('Position (m)');
143 grid on;
144 set(gca, 'FontSize', 10, 'Box', 'on');
145
146 subplot(3,1,2);
147 plot(t_vec, log_u, 'Color', [0.8500 0.3250 0.0980], 'LineWidth', 1.5);
148 title('Input Force', 'FontWeight', 'bold');
149 ylabel('Force (u)');
150 grid on;
151 set(gca, 'FontSize', 10, 'Box', 'on');
152
153 subplot(3,1,3);
154 plot(t_vec, log_r, 'k', 'LineWidth', 1.5);
155 title('Residual (Model Mismatch)', 'FontWeight', 'bold');
156 ylabel('Error Magnitude');
157 grid on;
158 set(gca, 'FontSize', 10, 'Box', 'on');
159
160 save('data/mpc_results.mat', 'log_y', 'log_u', 'log_r');
161 fprintf('Results saved to data/mpc_results.mat\n');

```

D. Part 4: Residual-Aware RL-MPC

The core logic for Q-Learning and Gain Scheduling.

Listing 4. RL-MPC Training Script

```

1 ME 418/518 Project: Residual-Aware RL-MPC
2 Residual Based Q-Learning for MPC Parameters
3
4 close all hidden; clear; clc;
5 rng(418);
6 1. Load Data & Models
7 load('data/identified_models.mat');
8 A = A_n4; B = B_n4; C = C_n4;
9 [nx, nu] = size(B); ny = size(C,1);
10
11 % RL and Simulation Parameters
12

```



```

13 N_episodes = 100;
14 T_final = 400;
15 goal_pos = 0.45;
16 2. RL Agent Configuration
17 The RL agent only observes the residual error from models as a
    categorical feature. Which makes the RL model specific for
    dealing with the model mismatch for MPC design. Instead of
    directly controlling the simulation, RL selects a parameter
    set based on model mismatch. Actions are as defined:
18 % State Space is discretized residual levels which are Low, Medium,
    High.
19 n_states = 3;
20 n_actions = 3;
21
22 % Residual Discretization Thresholds define what "Low", "Medium", "
    High"
23 % residual means. Interval is chosen from the sample distribution
    of
24 % baseline MPC residuals.
25 res_bins = [0.65, 0.95];
26
27 % Q-Table Initialization for tabular Q Learning (3x3 table to map
    states and
28 % actions)
29 Q_table = zeros(n_states, n_actions);
30
31 % Hyperparameters are conventional literature values
32 alpha = 0.2; % Learning Rate
33 gamma = 0.95; % Discount Factor
34 epsilon = 0.6; % Exploration Rate
35
36 Modes = struct();
37 Action 1: High-Residual Recovery = Model is failing (Residual is
    high). The linear model predicts the favored motion is good,
    but physics shows that the car is stuck in the valley. We tune
    the MPC to bias to recover from the ill-favored action.
38 Action 2: Low-Residual Tracking = When the model is accurate (
    Residual is low), RL basically trusts the identification model
    's prediction. However, tunes MPC parameter Q high to
    aggressively track the goal.
39 Action 3: Conservative Hold = RL decides that the state is uncertain
    . Thus penalizes input usage (High R) to stabilize the current
    state.
40 % Action 1: High Residual Recovery
41 Modes(1).Name = 'Recovery';
42 Modes(1).Q = 10; % Lowered trust to identification model
43 Modes(1).R = 0.01; % Ideal power for motor to overcome the recovery
44 Modes(1).Ref = -1.2; % Bias setpoint to the leftwards wall induce
    negative force
45
46 % Action 2: Low-Residual Tracking
47 Modes(2).Name = 'Tracking';
48 Modes(2).Q = 100; % Full trust to identification model
49 Modes(2).R = 0.01; % Ideal power for motor to track optimal climbing
    of the identification model
50 Modes(2).Ref = 0.45; % Standard Goal
51
52 % Action 3: Conservative Hold
53 Modes(3).Name = 'Stabilize';
54 Modes(3).Q = 1; % Minimized trust to identification model
55 Modes(3).R = 10; % High input penalty to compensate the controller.
    Lowered aggressiveness.
56 Modes(3).Ref = 0.45; % Standard Goal
57 3. Pre-compute MPC Matrices, Gain Scheduling
58 To ensure computational efficiency during the simulation loop, we
    pre-calculate the prediction matrices based on the N4SID model
    . We also pre-compute the quadprog matrices H for all three
    control modes. This establishes Gain Scheduling, allowing the
    controller to instantly switch between different tuning sets.
59 Np = 30; Nc = 10;
60
61 % Prediction Matrices: F, Phi
62 F = zeros(Np*ny, nx); Ap = A;
63 for k=1:Np, F(k,:) = C*Ap; Ap=Ap*A; end
64 H = zeros(Np*ny, Nc*nu);
65 for i=1:Np
66     for j=1:Nc
67         if j<=i
68             if (i-j)==0, term=C*B; else, term=C*(A^(i-j-1))*B; end
69             H(i,j) = term;
70         end
71     end
72 end
73 S = tril(ones(Nc)); Phi = H*S;
74
75 % Pre-compute QP Matrices for each Mode
76 for a = 1:n_actions
77     Qbar = Modes(a).Q * eye(Np);
78     Rbar = Modes(a).R * eye(Nc);
79     Hqp = 2 * (Phi' * Qbar * Phi + Rbar);
80     Modes(a).Hqp = (Hqp + Hqp')/2;
81     Modes(a).Qbar = Qbar;
82 end
83
84 % Constraints
85 u_max=1.0; u_min=-1.0; du_max=0.5;
86 A_du = [eye(Nc); -eye(Nc)]; b_du = du_max*ones(2*Nc,1);
87 A_u = [S; -S]; options = optimoptions('quadprog','Display','off');
88 reward_history = zeros(N_episodes, 1);
89 4. Q-Learning Loop
90 Core training phase where the RL agent learns to compensate for the
    linear model's limitations. At each time step, the agent
    observes the current Residual State and selects an MPC mode
    using an epsilon-greedy strategy. The linear MPC calculates
    the optimal input based on the selected mode, which is applied
    to the true nonlinear physics, and the agent updates its Q-
    Table based on a reward function that maximizes both potential
    and kinetic energy.
91 for ep = 1:N_episodes
92     % Reset simulation. Initial state assume low residual.
93     x_true = [-0.5; 0]; x_model = zeros(nx, 1); u_prev = 0;
94     s = 1;
95     ep_reward = 0;
96
97     for k = 1:T_final
98         % RL Step, choose parameter set. exploration vs exploitation
99         if rand < epsilon
100             a = randi(n_actions);
101         else
102             [~, a] = max(Q_table(s, :));
103         end
104
105         % MPC Step, Execute with selected parameters
106         current_ref = Modes(a).Ref;
107         r_window = current_ref * ones(Np, 1);
108         u_base = u_prev * ones(Nc, 1);
109         bk = F * x_model + H * u_base;
110
111         % f_qp is affected by Q and Ref
112         fqp = 2 * Phi' * Modes(a).Qbar * (bk - r_window);
113         b_u = [u_max*ones(Nc,1)-u_base; -u_min*ones(Nc,1)+u_base];
114         [dU, ~, exitflag] = quadprog(Modes(a).Hqp, fqp, [A_du; A_u],
            [b_du; b_u], [], [], [], [], [], options);
115         du = (exitflag>0)*dU(1);
116         u_apply = u_prev + du;
117
118         % Physics update step, true environment
119         curr_pos = x_true(1); curr_vel = x_true(2);
120         new_vel = curr_vel + u_apply*0.0015 - 0.0025*cos(3*curr_pos)
            ;
121         new_vel = max(min(new_vel, 0.07), -0.07);
122         new_pos = curr_pos + new_vel;
123         if new_pos<=-1.2, new_pos=-1.2; new_vel=0; end
124         if new_pos>=0.6, new_pos=0.6; new_vel=0; end
125         x_true = [new_pos; new_vel];
126
127         % Residual Calculation. Model prediction vs Real measurement
128         x_model = A * x_model + B * u_apply;
129         y_pred = C * x_model;
130         residual = abs(curr_pos - y_pred);
131
132         % Discretize Residual
133         if residual < res_bins(1), s_next = 1; % Low Residual
134         elseif residual < res_bins(2), s_next = 2; % Medium
135         else, s_next = 3; % High Residual
136         end
137
138         % Reward. Main idea is that not only having height is
            sufficient
139         % to reach the goal but also having a velocity is essential.
140         height_reward = (curr_pos + 0.5); % start point
141         energy_reward = 100 * (curr_vel^2); % scaling with the
            height reward (max energy reward is 0.49)
142         r = height_reward + energy_reward;
143
144         % Completion Bonus, sparse reward
145         if curr_pos >= 0.45
146             r = r + 5000;
147         end
148
149         % Q Learning Update. Online Learning.
150         Q_table(s, a) = Q_table(s, a) + alpha * (r + gamma*max(
            Q_table(s_next,:)) - Q_table(s, a));
151         s = s_next; u_prev = u_apply; ep_reward = ep_reward + r;
152         if curr_pos >= 0.45, break; end
153     end
154
155     epsilon = max(0.01, epsilon * 0.98);
156     reward_history(ep) = ep_reward;
157     if mod(ep, 20) == 0
158         fprintf('Ep %d | Reward: %.1f | Eps: %.2f | Success: %d\n',
            ep, ep_reward, epsilon, (curr_pos>=0.45));
159     end
160 end
161 5. Validation & Plotting
162 After training is complete, we validate the learned policy by
    running a final simulation using a Greedy Policy. We log the
    position, control inputs, residuals, and mode selections to
    visualize how the agent dynamically switches strategies to

```

```

        perform the swing up maneuver. Finally, the data is saved for
        visualization in the Python Gymnasium environment.
163 x_true = [-0.5; 0]; x_model = zeros(nx, 1); u_prev = 0; s = 1;
164 log_y = []; log_u = []; log_r = []; log_mode = [];
165
166 for k = 1:400
167     [~, a] = max(Q_table(s, :)); % Greedy Policy
168
169     % MPC Execution
170     current_ref = Modes(a).Ref;
171     r_window = current_ref * ones(Np, 1);
172     u_base = u_prev * ones(Nc, 1);
173     bk = F * x_model + H * u_base;
174     fqp = 2 * Phi' * Modes(a).Qbar * (bk - r_window);
175     b_u = [u_max*ones(Nc,1)-u_base; -u_min*ones(Nc,1)+u_base];
176     [dU, ~, exitflag] = quadprog(Modes(a).Hqp, fqp, [A_du; A_u], [
        b_du; b_u], [], [], [], [], [], options);
177     du = (exitflag>0)*dU(1); u_apply = u_prev + du;
178
179     % Physics
180     curr_pos = x_true(1); curr_vel = x_true(2);
181     new_vel = curr_vel + u_apply*0.0015 - 0.0025*cos(3*curr_pos);
182     new_vel = max(min(new_vel, 0.07), -0.07);
183     new_pos = curr_pos + new_vel;
184     if new_pos<=-1.2, new_pos=-1.2; new_vel=0; end
185     if new_pos>=0.6, new_pos=0.6; new_vel=0; end
186     x_true = [new_pos; new_vel];
187     x_model = A * x_model + B * u_apply;
188
189     % Residual
190     residual = abs(curr_pos - (C*x_model));
191     if residual < res_bins(1), s = 1; elseif residual < res_bins(2),
        s = 2; else, s = 3; end
192     u_prev = u_apply;
193
194     log_y = [log_y; curr_pos]; log_u = [log_u; u_apply]; log_r = [
        log_r; residual]; log_mode = [log_mode; a];
195     if curr_pos >= 0.45, break; end
196 end
197
198 figure('Name','RL-MPC Parameter Tuning');
199 subplot(4,1,1); plot(log_y,'c','LineWidth',1.5); ylabel(0.45,'r--');
        title('Position'); grid on;
200 subplot(4,1,2); plot(log_u,'m'); title('Control Input'); grid on;
201 subplot(4,1,3); plot(log_r,'r'); title('Residual (RL State)'); grid
        on;
202 subplot(4,1,4); plot(log_mode,'m'); title('Selected Parameter Set (
        Mode)'); grid on; ylim([0.5 3.5]);
203
204 save('data/rl_mpc_results.mat', 'log_y', 'log_u', 'log_r', 'log_mode
        ');
205 fprintf('Validation Complete.\n');
206
207 figure('Name','RL Training Progress');
208 plot(1:N_episodes, reward_history, 'Color','r', 'LineWidth', 1);
        hold on;
209 window_size = 10;
210 if N_episodes >= window_size
211     moving_avg = movmean(reward_history, window_size);
212     plot(1:N_episodes, moving_avg, 'b', 'LineWidth', 2);
213     legend('Episode Reward', 'Moving Average (Trend)', 'Location', '
        SouthEast');
214 else
215     legend('Episode Reward');
216 end
217
218 xlabel('Episode');
219 ylabel('Total Reward');
220 title(['Learning Curve (alpha= num2str(alpha) ', gamma= num2str(
        gamma) ')']);
221 grid on;
222 ylabel(5000, 'g--', 'Success Threshold', 'LineWidth', 1.5);

```

E. Part 5: Python Visualization

Code used to generate the final videos in Gymnasium using the graphviz library.

Listing 5. Python Visualization

```

1 from graphviz import Digraph
2 import scipy.io
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import gymnasium as gym
6 from gymnasium.wrappers import RecordVideo
7 import os
8
9 data_path = '../data/mpc_results.mat'
10 data = scipy.io.loadmat(data_path)
11
12 # Flatten the arrays to 1D for iteration
13 inputs = data['log_u'].flatten()
14
15 # We load positions only for validation of the simulation, but we
        drive the simulation with only the inputs
16 print(f"Total simulation steps to replay: {len(inputs)}")
17
18 env = gym.make('MountainCarContinuous-v0', render_mode='rgb_array')
19
20 # Wrap the environment to record video
21 video_folder = '../videos'
22 env = RecordVideo(env, video_folder=video_folder,
23     episode_trigger=lambda x: True,
24     name_prefix='baseline_mpc_failure')
25 # Force the exact initial condition used in MATLAB: x0 = [-0.5, 0.0]
26 obs, info = env.reset(seed=418)
27 env.unwrapped.state = np.array([-0.5, 0.0])
28 print("Starting Simulation")
29 for i, u in enumerate(inputs):
30     action = np.array([u], dtype=np.float32)
31     obs, reward, terminated, truncated, info = env.step(action)
32     if terminated:
33         break
34
35 env.close()
36 print(f"Replay Finished.")
37
38 data_path = '../data/rl_mpc_results.mat'
39 rl_data = scipy.io.loadmat(data_path)
40
41 log_y = rl_data['log_y'].flatten() # Position
42 log_u = rl_data['log_u'].flatten() # Control Input
43 log_r = rl_data['log_r'].flatten() # Residual
44 log_mode = rl_data['log_mode'].flatten() # Selected Mode (1, 2, or
        3)
45 time_steps = np.arange(len(log_y))
46
47 env = gym.make('MountainCarContinuous-v0', render_mode='rgb_array')
48
49 video_folder = '../videos'
50 env = RecordVideo(env, video_folder=video_folder,
51     episode_trigger=lambda x: True,
52     name_prefix='rl_mpc_success')
53
54 # Force the same initial condition as MATLAB: x0 = [-0.5, 0.0]
55 obs, info = env.reset(seed=418)
56 env.unwrapped.state = np.array([-0.5, 0.0])
57 print("Rendering RL-MPC Solution...")
58
59 for i, u in enumerate(log_u):
60     action = np.array([u], dtype=np.float32)
61     obs, reward, terminated, truncated, info = env.step(action)
62     if terminated or i == len(log_u)-1:
63         break
64
65 env.close()
66
67 plt.style.use('default')
68 fig, ax = plt.subplots(4, 1, figsize=(10, 12), sharex=True)
69
70 # 1. Position
71 ax[0].plot(time_steps, log_y, color='tab:blue', linewidth=2, label='
        Car Position')
72 ax[0].axhline(0.45, color='tab:red', linestyle='--', linewidth=2,
        label='Goal (0.45)')
73 ax[0].set_ylabel('Position (m)')
74 ax[0].set_title('Trajectory: Successful Swing-Up Strategy',
        fontweight='bold')
75 ax[0].grid(True, linestyle='--', alpha=0.6)
76 ax[0].legend(loc='lower right')
77
78 # 2. Control Input
79 ax[1].plot(time_steps, log_u, color='tab:orange', linewidth=1.5)
80 ax[1].set_ylabel('Force (N)')
81 ax[1].set_title('Control Effort (u)', fontweight='bold')
82 ax[1].grid(True, linestyle='--', alpha=0.6)
83
84 # 3. Residual
85 ax[2].plot(time_steps, log_r, color='tab:red', linewidth=1.5)
86 ax[2].set_ylabel('Error Magnitude')
87 ax[2].set_title('Model Uncertainty Signal (Residual)', fontweight='
        bold')
88 ax[2].grid(True, linestyle='--', alpha=0.6)
89
90 # 4. Mode Switching
91 ax[3].plot(time_steps, log_mode, where='post', color='tab:green',
        linewidth=2)
92 ax[3].set_yticklabels(['Recovery (1)', 'Tracking (2)', 'Stabilize
        (3)'])
93 ax[3].set_ylabel('RL Mode')
94 ax[3].set_xlabel('Simulation Steps')
95 ax[3].set_title('Adaptive Gain Scheduling Policy', fontweight='bold'
        )

```

```

98 ax[3].grid(True, linestyle='--', alpha=0.6)
99
100 plt.tight_layout()
101 plt.show()
102
103 output_dir = '../images'
104 if not os.path.exists(output_dir):
105     os.makedirs(output_dir)
106
107
108 dot = Digraph('RL_MPC_Compact', comment='Residual-Aware RL-MPC
Compact Flow')
109 dot.attr(rankdir='TB', splines='ortho', nodesep='0.6', ranksep='0.6',
, fontname='Helvetica')
110 dot.attr('node', shape='box', style='filled', fontname='Helvetica',
fontsize='11', height='0.4')
111 dot.attr('edge', fontsize='9', fontcolor='dimgrey')
112
113 # IDENTIFICATION & DATA
114 with dot.subgraph(name='cluster_0') as c:
115     c.attr(label='1. OFFLINE IDENTIFICATION', style='filled', color=
'lightgrey', fillcolor='#E3F2FD')
116
117     c.node('SimEnv', 'Gymnasium Env\n(MountainCar-v0)', shape='
ellipse', fillcolor='white')
118     c.node('DataGen', 'Data Generation\n(PRBS Signal)', shape='note'
, fillcolor='#FFF9C4')
119     c.node('N4SID', 'System Identification\n(N4SID Algorithm)',
fillcolor='#BBDEFB')
120     c.node('LinearModel', 'Linear Model\n(A, B, C)', shape='
component', fillcolor='gold')
121
122     c.edge('SimEnv', 'DataGen', label='Physics')
123     c.edge('DataGen', 'N4SID', label='u, y data')
124     c.edge('N4SID', 'LinearModel', label='Identify')
125
126 # RL
127 with dot.subgraph(name='cluster_1') as c:
128     c.attr(label='2. RL AGENT TRAINING', style='filled', color='
lightgrey', fillcolor='#E8F5E9')
129
130     c.node('ResidualCalc', 'Residual Calculation\n|y_meas - y_pred|'
, shape='diamond', fillcolor='#FFCCBC')
131     c.node('RL_Loop', 'Q-Learning Loop\n(State: Residual Levels)',
fillcolor='#C8E6C9')
132     c.node('Reward', 'Reward Function\n(Energy + Goal)', shape='
ellipse', fillcolor='white')
133     c.node('QTable', 'Learned Policy\n(Q-Table)', shape='folder',
fillcolor='orange')
134
135     c.edge('ResidualCalc', 'RL_Loop', label='State (r_t)')
136     c.edge('RL_Loop', 'Reward', label='Action')
137     c.edge('Reward', 'RL_Loop', label='Feedback')
138     c.edge('RL_Loop', 'QTable', label='Save Policy')
139
140 # RL-MPC
141 with dot.subgraph(name='cluster_2') as c:
142     c.attr(label='3. ONLINE RESIDUAL-AWARE CONTROL', style='filled',
color='black', fillcolor='white')
143
144     c.node('RealPlant', 'Real Plant\n(Nonlinear Physics)', shape='
doublecircle', fillcolor='#FFAB91')
145     c.node('Estimator', 'Linear Predictor\n(x_k+1 = Ax+Bu)', shape='
box', fillcolor='gold')
146     c.node('Comparator', 'Error Check\n(Residual)', shape='circle',
width='0.8', fillcolor='#B3E5FC')
147     c.node('Agent', 'RL Agent\n(Mode Selector)', fillcolor='orange')
148     c.node('MPC', 'MPC\n(Gain Scheduled)', style='filled', fillcolor
='#4FC3F7', penwidth='2')
149
150     c.edge('RealPlant', 'Comparator', label='y_meas')
151     c.edge('Estimator', 'Comparator', label='y_pred')
152     c.edge('Comparator', 'Agent', label='Residual')
153     c.edge('Agent', 'MPC', label='Select Mode')
154     c.edge('MPC', 'RealPlant', label='u(t)', color='red', penwidth='
2.0')
155     c.edge('MPC', 'Estimator', label='u(t)', style='dashed')
156
157 # Global Connections
158 dot.edge('LinearModel', 'Estimator', style='dashed', label='Model')
159 dot.edge('LinearModel', 'ResidualCalc', style='dashed', constraint='
false')
160 dot.edge('QTable', 'Agent', style='bold', label='Load Policy')
161
162 output_path = os.path.join(output_dir, 'rl_mpc_workflow')
163 dot.render(output_path, format='png', cleanup=True)
164 dot

```
