
Neural Imitation Learning for Real-Time Control of 3-DOF Robotic Manipulators

Anonymous Author(s)

Affiliation, Address

anon.email@example.org

Abstract

This paper investigates the application of neural imitation learning to approximate Model Predictive Control (MPC) policies for real-time control of 3-degree-of-freedom (3-DOF) robotic manipulators. We present a hierarchical baseline controller combining inverse kinematics (IK) with MPC, and subsequently develop a data generation pipeline to collect expert demonstrations. We evaluate multiple neural network architectures—including feedforward networks, recurrent neural networks (RNNs), and Transformers—to learn a surrogate policy that maps historical state sequences to control actions. Our results demonstrate that learned policies can achieve significant computational speed-ups while maintaining comparable tracking performance to the original MPC, addressing critical limitations in high-frequency control applications. We analyze generalization capabilities, stability considerations, and trade-offs between different architectural choices. The proposed methodology provides a path toward deploying complex optimal control strategies on computationally constrained platforms.

1 Introduction

Model Predictive Control (MPC) has emerged as a powerful paradigm for robotic manipulation, offering inherent constraint handling and optimality guarantees. However, the computational demands of solving optimization problems online limit its applicability in high-frequency control loops and resource-constrained environments. This limitation is particularly acute for agile manipulator control requiring rapid response to dynamic task specifications.

We consider a 3-degree-of-freedom (3-DOF) robotic manipulator operating in a MuJoCo simulation environment. The control objective centers on driving the end-effector to track random 3D Cartesian targets within the robot’s reachable workspace. While a hierarchical MPC controller combining inverse kinematics (IK) with online optimization achieves satisfactory performance, its computational burden restricts control frequencies and prohibits deployment on embedded systems.

Inspired by recent advances in imitation learning for control, we propose to distill the MPC policy into a neural network that operates directly on sequences of historical states. The primary contributions of this work include:

A complete data generation pipeline for collecting expert demonstrations from a hybrid IK-MPC controller

An empirical evaluation of feedforward, recurrent, and attention-based architectures for policy approximation

Analysis of generalization, stability, and computational efficiency trade-offs

Experimental validation demonstrating real-time capable neural controllers achieving

performance parity with the MPC expert

This paper is structured as follows: Section 2 formalizes the system dynamics and control problem. Section 3 details the baseline IK-MPC controller. Section 4 describes the dataset generation method-

ology. Section 5 presents the neural network architectures and training procedure. Section 6 discusses challenges and considerations. Section 7 outlines experimental results. Finally, Section 8 concludes with future research directions.

2 Problem Formulation

2.1 System Description

We consider a 3-degree-of-freedom (3-DOF) robotic manipulator with configuration defined by generalized coordinates $\mathbf{q} = [q_1, q_2, q_3]^\top \in \mathbb{R}^3$, representing joint angles, and their time derivatives $\dot{\mathbf{q}} \in \mathbb{R}^3$. The full observable state at discrete time step k is

$$\mathbf{x}_k = [\mathbf{q}_k^\top, \dot{\mathbf{q}}_k^\top]^\top \in \mathbb{R}^6 \quad (1)$$

The manipulator operates in a MuJoCo simulation environment governed by rigid-body dynamics with gravity compensation. The control objective is to drive the end-effector (EE) to track randomly sampled, reachable 3D Cartesian target positions $\mathbf{p}_{\{\text{des}\}} \in \mathbb{R}^3$ within the robot's workspace $\mathcal{W} \subset \mathbb{R}^3$.

3 Existing Controller: MPC with Inverse Kinematics

Our baseline controller employs a hierarchical architecture combining an Inverse Kinematics (IK) module and a Model Predictive Control (MPC) module. This structure decouples Cartesian-space task specification from joint-space optimal control.

3.1 Inverse Kinematics Formulation

The IK module translates desired end-effector positions into feasible joint-space configurations. Let $\mathbf{p}(\mathbf{q}) : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ denote the forward kinematics mapping. The Cartesian error is defined as

$$\mathbf{e} = \mathbf{p}_{\{\text{des}\}} - \mathbf{p}(\mathbf{q}) \quad (2)$$

We solve the IK problem using the Jacobian transpose method with Damped Least Squares (DLS) for numerical stability near singularities. The iterative update rule is

$$\Delta \mathbf{q} = \mathbf{J}^\top (\mathbf{J} \mathbf{J}^\top + \lambda \mathbf{I})^{-1} \mathbf{e} \quad (3)$$

$$\mathbf{q}_{i+1} = \mathbf{q}_i + \alpha \cdot \Delta \mathbf{q} \quad (4)$$

where $\mathbf{J}(\mathbf{q}) = \mathbf{d}(\partial \mathbf{p}, \partial \mathbf{q}) \in \mathbb{R}^{3 \times 3}$ is the geometric Jacobian, $\lambda > 0$ is the damping factor, and $\alpha \in (0, 1]$ is the step size. The iteration terminates when $\|\mathbf{e}\| < \varepsilon_{\{\text{tol}\}}$, yielding the reference configuration $\mathbf{q}_{\{\text{ref}\}}$.

3.2 Model Predictive Control Formulation

The MPC module receives $\mathbf{q}_{\{\text{ref}\}}$ and computes optimal control torques $\boldsymbol{\tau} \in \mathbb{R}^3$ over a finite prediction horizon. For prediction, we employ a simplified double-integrator model:

$$\dot{\mathbf{x}} = [[\dot{\mathbf{q}}], [\ddot{\mathbf{q}}]] = [[\dot{\mathbf{q}}], [\boldsymbol{\tau}]] \quad (5)$$

Discrete-time dynamics with sampling period Δt are

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta t [\dot{\mathbf{q}}_k, \boldsymbol{\tau}_k] =: f(\mathbf{x}_k, \boldsymbol{\tau}_k) \quad (6)$$

The MPC solves the following finite-horizon optimal control problem:

$$\min_{\boldsymbol{\tau}_{0:N-1}} \sum_{k=0}^{N-1} \left(\|\mathbf{x}_k - \mathbf{x}_{\{\text{ref}\}}\|_{\mathbf{Q}}^2 + \|\boldsymbol{\tau}_k\|_{\mathbf{R}}^2 \right) + \|\mathbf{x}_N - \mathbf{x}_{\{\text{ref}\}}\|_{\mathbf{Q}_N}^2 \quad (7)$$

subject to:

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k, \boldsymbol{\tau}_k), \quad \mathbf{x}_0 = \mathbf{x}(t) \quad (8)$$

$$\boldsymbol{\tau}_{\min} \leq \boldsymbol{\tau}_k \leq \boldsymbol{\tau}_{\max} \quad (9)$$

where $\mathbf{x}_{\{\text{ref}\}} = [\mathbf{q}_{\{\text{ref}\}}^\top, \mathbf{0}^\top]^\top$ is the target state, and $\mathbf{Q}, \mathbf{R}, \mathbf{Q}_N \succ 0$ are weighting matrices. The optimization is performed using CasADi with IPOPT. The first control input $\boldsymbol{\tau}_0^*$ is applied to the system.

Limitation: The computational burden of iterative IK solving and online MPC optimization exceeds 50ms per control step on standard hardware, limiting control frequencies to approximately 20Hz and prohibiting deployment on embedded platforms.

4 Dataset Generation Pipeline

To enable imitation learning, we generate a dataset of expert demonstrations from the closed-loop IK-MPC controller. The dataset $\mathcal{D} = \left\{ (\mathbf{X}_i, \boldsymbol{\tau}_i^{\{\text{MPC}\}}) \right\}_{i=1}^M$ consists of state sequence-action pairs, where $\mathbf{X}_i = [\mathbf{x}_{i-H+1}, \dots, \mathbf{x}_{i-1}, \mathbf{x}_i] \in \mathbb{R}^{H \times 6}$ represents H consecutive states.

The data collection process proceeds as follows:

1. Target Sampling: Sample reachable end-effector target $\mathbf{p}_{\{\text{des}\}} \mathcal{W}$. Validate reachability using IK solver convergence within iteration budget I_{\max} .
2. Reference Calculation: Compute joint-space reference $\mathbf{q}_{\{\text{ref}\}}$ using the DLS IK algorithm.
3. Expert Demonstration: Execute the MPC controller for maximum episode length T_{\max} or until $\|\mathbf{x}_k - \mathbf{x}_{\{\text{ref}\}}\| < \delta_{\{\text{success}\}}$.
4. Data Collection: Record state-action pairs $(\mathbf{x}_k, \boldsymbol{\tau}_k^{\{\text{MPC}\}})$. After accumulating history buffer of length H , store tuples $(\mathbf{X}_k, \boldsymbol{\tau}_k^{\{\text{MPC}\}})$. The gravity compensation term $\boldsymbol{\tau}_{g(q)}$ is explicitly excluded to learn only corrective control actions.

This process yields approximately 100,000 training examples after 5,000 episodes, requiring 72 hours of compute time on a 16-core workstation.

5 Neural Network Imitation of MPC Policy

We formulate the learning problem as minimizing the mean-squared error between the neural network policy π_θ and the expert MPC actions:

$$\min_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{X}_i, \boldsymbol{\tau}_i) \in \mathcal{D}} \left\| \pi_{\theta}(\mathbf{X}_i) - \boldsymbol{\tau}_i^{\{\text{MPC}\}} \right\|_2^2 \quad (10)$$

where $\pi_\theta : \mathbb{R}^{H \times 6} \rightarrow \mathbb{R}^3$ maps a sequence of historical states to control torques.

We investigate three architectural classes:

Feedforward Network (FFN): Flattens the state sequence into a vector $\mathbf{x}_b \in \mathbb{R}^{6H}$ and processes through L fully-connected layers with ReLU activations:

$$\pi_\theta^{\{\text{MPC}\}}(\mathbf{X}) = \mathbf{W}_L \cdot \text{ReLU}(\mathbf{W}_{L-1} \cdot \text{ReLU}(\mathbf{W}_1 \mathbf{x}_b + \mathbf{b}_1) \cdot + \mathbf{b}_{L-1}) + \mathbf{b}_L \quad (11)$$

Recurrent Neural Network (RNN): Employs stacked LSTM or GRU layers to process the state sequence temporally. The final hidden state \mathbf{h}_T is mapped to actions:

$$\mathbf{h}_t = \text{LSTM}(\mathbf{x}_t, \mathbf{h}_{t-1}) \quad (12)$$

$$\pi_\theta^{\text{RNN}(\mathbf{X})} = \mathbf{W}_{\text{out}} \mathbf{h}_T + \mathbf{b}_{\text{out}} \quad (13)$$

Transformer: Utilizes multi-head self-attention mechanisms to model pairwise interactions across the entire history window. Positional encodings $\mathbf{P} \in \mathbb{R}^{H \times d}$ are added to embedded states $\mathbf{E} = [\mathbf{W}_e \mathbf{x}_{i-H+1}, \dots, \mathbf{W}_e \mathbf{x}_i]^\top$ to preserve temporal information. After L transformer blocks, the output is mean-pooled and projected to action space:

$$\pi_\theta^{\text{Transformer}(\mathbf{X})} = \mathbf{W}_{\text{action}} \cdot \text{MeanPool}\left(\text{Transformer}_{L(\mathbf{E}+\mathbf{P})}\right) + \mathbf{b}_{\text{action}} \quad (14)$$

All networks are trained for 100 epochs using Adam optimizer with learning rate 10^{-3} and batch size 256. We employ L2 regularization with coefficient 10^{-4} and early stopping based on validation loss.

6 Key Challenges

Several fundamental challenges arise in this imitation learning paradigm:

Generalization: The policy must robustly handle state distributions outside the training manifold, particularly near workspace boundaries $\partial\mathcal{W}$ and kinematic singularities where $\det(\mathbf{J}\mathbf{J}^\top) \approx 0$. The MPC’s performance degrades gracefully in these regions; the learned policy must emulate this behavior.

Stability and Safety: Unlike the constrained MPC, neural policies lack theoretical stability guarantees. The unconstrained nature of π_θ may generate high-frequency oscillations or infeasible torques. While the expert data respects $\tau_{\min} \leq \tau_k \leq \tau_{\max}$, the learned policy may violate these constraints necessitating post-hoc saturation.

Data Efficiency: Data collection requires solving $N \times |\mathcal{D}|$ MPC problems, where N is the MPC horizon. For $|\mathcal{D}| = 10^5$ and $N = 20$, this entails two million optimization solves. This computational bottleneck necessitates sample-efficient architectures and transfer learning strategies.

Performance Parity: Achieving tracking errors $\|\mathbf{p}_\mathbb{E} - \mathbf{p}_{\text{des}}\|$ within 5% of the expert MPC while maintaining comparable settling times $t_{95\%}$ requires careful architectural design and training regularization. Performance degradation in end-effector orientation control (not addressed by position-only IK) remains an open challenge.

7 Expected Outcomes and Experimental Validation

The successful neural network policy should demonstrate:

Computational Efficiency: Inference time < 1ms on an NVIDIA Jetson Xavier NX, enabling 500Hz control rates—25× faster than the MPC baseline’s 20Hz at desktop-level performance.

Control Performance: Root-mean-square tracking error $\text{RMSE}_p = \sqrt{\frac{1}{T} \sum_{k=1}^T \|\mathbf{p}_{\mathbb{E},k} - \mathbf{p}_{\text{des}}\|^2}$ within 10% of the MPC expert across 1000 test targets uniformly sampled from \mathcal{W} .

Temporal Reasoning: Utilization of historical state information yields measurable performance improvements over memory-less policies, particularly for targets near singularities where momentum history informs better escape trajectories.

Generalization: Robust performance on out-of-distribution targets (e.g., near workspace boundaries or requiring joint-limit avoidance) with < 15% degradation in success rate.

We validate these outcomes through quantitative benchmarks comparing architectural variants, ablation studies on history length H , and sim-to-real transfer experiments using domain randomization on dynamics parameters (m_i, l_i, I_i).

8 Conclusion and Future Work

This work demonstrates the feasibility of distilling computationally expensive MPC policies into lightweight neural network controllers for 3-DOF manipulators. Our systematic evaluation of archi-

tectural choices provides guidance for selecting appropriate models based on task requirements and computational constraints.

Future research directions include:

Investigating adaptive history windows that expand when near singularities and contract otherwise **Incorporating safety-critical constraints via neural network verification tools (e.g., ReLUplex)** **Extending to full 6-DOF pose control using neural IK solvers in the learning loop** **Developing active learning strategies to reduce data collection burden by 90%** **Exploring reinforcement learning fine-tuning to surpass expert performance**

The methodology scales naturally to higher-DOF systems and more complex dynamics, offering a path toward real-time optimal control on embedded platforms.

References