

Bridging the Gaps: Learning Verifiable Model-Free Quadratic Programming Controllers Inspired by Model Predictive Control

Yiwen Lu
Zishuo Li
Yihan Zhou
Na Li
Yilin Mo

LUYW20@MAILS.TSINGHUA.EDU.CN
LIZS19@MAILS.TSINGHUA.EDU.CN
ZHOUYH23@MAILS.TSINGHUA.EDU.CN
NALI@SEAS.HARVARD.EDU
YLMO@TSINGHUA.EDU.CN

Abstract

In this paper, we introduce a new class of parameterized controllers, drawing inspiration from Model Predictive Control (MPC). These controllers adopt a Quadratic Programming (QP) structure similar to linear MPC, with problem parameters being learned rather than derived from models. This approach may address the limitations of commonly learned controllers with Multi-Layer Perceptron (MLP) architecture in deep reinforcement learning, in terms of explainability and performance guarantees. The learned controllers not only possess verifiable properties like persistent feasibility and asymptotic stability akin to MPC, but they also empirically match MPC and MLP controllers in control performance. Moreover, they are more computationally efficient in implementation compared to MPC and require significantly fewer learnable policy parameters than MLP controllers. Practical application is demonstrated through a vehicle drift maneuvering task, showcasing the potential of these controllers in real-world scenarios.

1. Introduction

Recent years have witnessed the development of Deep Reinforcement Learning (DRL) in the domain of control (Lillicrap et al., 2015; Duan et al., 2016; Haarnoja et al., 2018), with the locomotion control of agile robots (Xie et al., 2018; Li et al., 2021; Margolis et al., 2022; Rudin et al., 2022) being a notable example. Many successful control applications use Multi-Layer Perceptron (MLP) as the entirety or a part of the policy, which, despite their remarkable empirical performance, face limitations in terms of explainability (Agogino et al., 2019) and performance guarantees (Osinenko et al., 2022). Research efforts have been devoted to the stability verification of MLP controllers (Dai et al., 2021; Zhou et al., 2022), structured controller parameterizations (Srouji et al., 2018; Johanning et al., 2019; Sattar and Oymak, 2020; Ni et al., 2021), or a combination of both (Zinage and Bakolas, 2023), and the learning of explainable and verifiable controllers has remained an active topic.

On the other hand, for ensuring stability and safety, Model Predictive Control (MPC) has been well-studied (Morari and Lee, 1999; Qin and Badgwell, 2003; Schwenzler et al., 2021). Recently, there has been a growing interest on augmenting MPC with learning, a large portion of which are focused on addressing the challenges in designing critical components of MPC like prediction models (Desaraju and Michael, 2016; Soloperto et al., 2018; Hewing et al., 2019), terminal costs and constraints (Brunner et al., 2015; Rosolia and Borrelli, 2017; Abdulfattokhov et al., 2021), and stage costs (Englert et al., 2017; Menner et al., 2019); readers are referred to Hewing et al. (2020) for a comprehensive overview. Most of the methods follow a model-based framework – using data to estimate a model and then perform optimal control at each receding horizon. These methods still

suffer various challenges, such as requiring intensive computation at each time horizon, and making myopic decisions that lead to infeasibility or inefficiency in the long run.

Motivated by the advantages of DRL and MPC, we propose an MPC-inspired but model-free controller. Noting the fact that linear MPC formulates each step as a Quadratic Programming (QP) problem, we consider a parameterized class of controllers with QP structure similar to MPC. However, the key distinction lies in the approach to obtaining the QP problem parameters: instead of deriving them from a model, they are learned. This approach ensures that the resulting controllers not only have theoretical guarantees akin to MPC, but also demonstrate competitive performance and computational ease when compared empirically to MPC and MLP controllers.

Contrasting with works from the learning community, such as [Amos et al. \(2018\)](#); [Ha and Schmidhuber \(2018\)](#); [Hafner et al. \(2019\)](#); [Hansen et al. \(2023\)](#); [LeCun \(2022\)](#), which adapt MPC ideas by incorporating model-based planning into broader learning frameworks, our work retains the classical MPC structure of well-formulated Quadratic Programming (QP). While these learning works aim for versatility and complexity, integrating MPC concepts into complicated learned models, they often lack control-theoretic guarantees central to classical MPC. In contrast, our approach specializes in control tasks with a focus on performance guarantees and computational efficiency. This specialization does not limit the applicability of our method: it empirically generalizes beyond simple linear systems to practical scenarios in real robot systems, as exemplified in an aggressive vehicle control setting.

Our Contribution. In this paper, we propose a new parameterized class of MPC-inspired controllers. Specifically, our controller employs an unrolled QP solver, structured similarly to a deep neural network, with the learnable parameters being those required for determining the underlying QP problem. To train the parameters of the controller, existing DRL methods would be applied, such as PPO ([Schulman et al., 2017](#)). However, in contrast to most DRL-trained controllers, which often lack rigorous theoretical guarantees, our MPC-inspired controller is proven to enjoy verifiable properties like recursive feasibility and asymptotic stability. We also compare the proposed controller on benchmark tasks with other methods such as classical MPC and DRL-trained neural network controllers, showing that our proposed controller enjoys lighter computation and increased robustness. Lastly, though the method is formally introduced and the performance is proved on linear systems, we verify the generalizability of the proposed approach via vehicle drift maneuvering, a challenging nonlinear robotics control task, demonstrating potential in a wider range of real-world applications.

2. Problem Formulation and Preliminaries

Notations. We use subscript to denote the time index, e.g., x_k stands for the system state at step k , and we use $x_{0:k}$ as a shorthand notation for the sequence x_0, x_1, \dots, x_k . We use superscript to denote the iteration index in an iterative algorithm, e.g. y^i stands for the variable y at the i -th iteration. We use bracketed subscript to denote slicing operation on a vector, e.g., $v_{[1]}$ denotes the first element of the vector v , and $v_{[1:n]}$ denotes its first n elements. The set of positive definite $n \times n$ matrices is denoted as \mathbb{S}_{++}^n , and the nonnegative orthant of \mathbb{R}^n is denoted as \mathbb{R}_+^n . The Kronecker product of two matrices A and B is denoted as $A \otimes B$. The block diagonal matrix with diagonal blocks A_1, \dots, A_n is denoted as $\text{diag}(A_1, \dots, A_n)$.

2.1. Problem Formulation

In this paper, we consider the discrete-time infinite-horizon constrained linear-quadratic optimal control problem, formulated as follows:

Problem 1 (Infinite-horizon constrained linear-quadratic optimal control)

$$\underset{u_{0:\infty}}{\text{minimize}} \quad \limsup_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} (x_{k+1} - r)^\top Q (x_{k+1} - r) + u_k^\top R u_k, \quad (1)$$

$$\text{subject to} \quad x_{k+1} = A x_k + B u_k, \quad (2)$$

$$u_{\min} \leq u_k \leq u_{\max}, x_{\min} \leq x_{k+1} \leq x_{\max}, \quad (3)$$

where $x_k \in \mathbb{R}^{n_{sys}}$ are state vectors, $u_k \in \mathbb{R}^{m_{sys}}$ are control input vectors, $r \in \mathbb{R}^{n_{sys}}$ is the reference state, $A \in \mathbb{R}^{n_{sys} \times n_{sys}}$ and $B \in \mathbb{R}^{n_{sys} \times m_{sys}}$ are the system and input matrices, $Q \in \mathbb{S}_{++}^{n_{sys}}$ and $R \in \mathbb{S}_{++}^{m_{sys}}$ are the stage cost matrices, and $u_{\min}, u_{\max} \in \mathbb{R}^{m_{sys}}$ and $x_{\min}, x_{\max} \in \mathbb{R}^{n_{sys}}$ are bounds on control input and state respectively. It is assumed without loss of generality that (A, B) is controllable.

2.2. Linear MPC and its QP Representation

Problem 1 is typically computationally intractable due to infinite planning horizon and constraints. A commonly adopted approximation is truncating it to finite horizon N , and solving the problem formulated in Problem 2 at each time step, with x_0 being the current state. The first control input u_0^* from the optimal solution is applied to the system in a receding horizon fashion.

Problem 2 (Naive Linear MPC)

$$\underset{x_{1:N}, u_{0:N-1}}{\text{minimize}} \quad \sum_{k=0}^{N-1} (x_{k+1} - r)^\top Q (x_{k+1} - r) + u_k^\top R u_k, \quad (4)$$

$$\text{subject to} \quad x_{k+1} = A x_k + B u_k, \quad k = 0, \dots, N-1, \quad (5)$$

$$u_{\min} \leq u_k \leq u_{\max}, x_{\min} \leq x_{k+1} \leq x_{\max}, \quad k = 0, \dots, N-1. \quad (6)$$

The above MPC problem can be cast into a Quadratic Programming (QP) problem in the following standard form¹:

Problem 3 (Standard-form QP)

$$\underset{y}{\text{minimize}} \quad \frac{1}{2} y^\top P y + q^\top y, \quad \text{subject to} \quad H y + b \geq 0, \quad (7)$$

where $y \in \mathbb{R}^{n_{qp}}$, $P \in \mathbb{S}_{++}^{n_{qp}}$, $q \in \mathbb{R}^{n_{qp}}$, $H \in \mathbb{R}^{m_{qp} \times n_{qp}}$, and $b \in \mathbb{R}^{m_{qp}}$.

The translation from Problem 2 to Problem 3 can be performed by using the control sequence $y = [u_0^\top \ \dots \ u_{N-1}^\top]^\top$ as the decision variable, and eliminating the equality constraints (5) by representing the trajectory $x_{1:N}$ using y . The resulting QP problem size and parameters are:

$$n_{qp} = N m_{sys}, \quad m_{qp} = N (m_{sys} + n_{sys}), \quad (8)$$

$$P = B^\top Q B + R, \quad q = 2B^\top Q (A x_0 - r), \quad H = -C B - D, \quad b = e - C A x_0, \quad (9)$$

1. We term the problem as ‘‘Naive Linear MPC’’ since no terminal set or cost is included. Although only the ‘‘naive’’ version is presented for simplicity, one can derive a similar QP formulation for linear MPC with quadratic terminal cost and affine terminal constraint.

where

$$\mathbf{A} = \begin{bmatrix} A \\ A^2 \\ \vdots \\ A^N \end{bmatrix}, \mathbf{B} = \begin{bmatrix} B & & & \\ AB & B & & \\ \vdots & & \ddots & \\ A^{N-1}B & \dots & & B \end{bmatrix}, \begin{aligned} \mathbf{C} &= I_N \otimes \text{diag}(I_{n_{sys}}, -I_{n_{sys}}, \mathbf{0}_{2m_{sys} \times 2m_{sys}}), \\ \mathbf{D} &= I_N \otimes \text{diag}(\mathbf{0}_{2n_{sys} \times 2n_{sys}}, I_{m_{sys}}, -I_{m_{sys}}), \\ \mathbf{e} &= I_N \otimes [x_{\max}^\top \quad -x_{\min}^\top \quad u_{\max}^\top \quad -u_{\min}^\top]^\top, \\ \mathbf{r} &= I_N \otimes r, \mathbf{Q} = I_N \otimes Q, \mathbf{R} = I_N \otimes R. \end{aligned} \quad (10)$$

2.3. Algorithm for Solving QPs

A family of efficient methods for solving QPs is operator splitting algorithms (Ryu and Yin, 2022), which are adopted by existing solvers such as OSQP (Stellato et al., 2020). An iteration of an operator splitting algorithm for solving QPs can generally be represented as a combination of affine transformations and projections on the variable. For example, with the Primal-Dual Hybrid Gradient (PDHG) (Chambolle and Pock, 2011) algorithm, an iteration can be expressed in the succinct form shown as follows :

$$z^{i+1} = \Pi_{\mathbb{R}_+^{m_{qp}}} ((I - 2\alpha G)(z^i + \lambda^i) - 2\alpha\mu), \quad \lambda^{i+1} = G(z^i + \lambda^i) + \mu, \quad (11)$$

where $z = Hy + b \in \mathbb{R}^{m_{qp}}$ is the primal variable of an equivalent form of the original problem (7), $\lambda \in \mathbb{R}^{m_{qp}}$ is a dual variable introduced by the same equivalent form, $\alpha > 0$ is the step size, and the parameters in the iteration are:

$$G = (I + HP^{-1}H^\top)^{-1}, \quad \mu = G(HP^{-1}q - b). \quad (12)$$

Once one obtains an approximate solution z^i , the original variable can be recovered from the equality-constrained QP problem $y^i \in \arg \min \{\frac{1}{2}y^\top Py + q^\top y \mid Hy + b = z^i\}$, whose KKT condition is a linear equation. Hence, by finding the least-square solution of this KKT equation, y^i can be explicitly represented as:

$$y^i = -P^{-1}q + P^{-1}H^\top(HP^{-1}H^\top)^\dagger(z^i - b + HP^{-1}q). \quad (13)$$

Theorem 1 *If $0 < \alpha < 1$ and the problem (7) is feasible, then the iterations (11) yields $y^i \rightarrow y^*$, where y^* is the optimal solution of the original problem. Furthermore, the suboptimality gap satisfies:*

$$p^i - p^* \leq \|\lambda^i\|_2 \|r_{prim}^i\|_2 + \|y^i - y^*\|_2 \|r_{dual}^i\|_2, \quad (14)$$

where p^i, p^* are the primal value at iteration i and the optimal primal value respectively, and r_{prim}^i, r_{dual}^i are the primal and dual residuals defined as follows:

$$r_{prim}^i = Hy^i + b - z^i, r_{dual}^i = Py^i + q + H^\top \lambda^i. \quad (15)$$

The iteration (11) on the primal-dual variable pair (z, λ) can be implemented by interleaving an affine transformation, whose parameters (G, μ) depend on the problem parameters (P, q, H, b) , and a projection of the z -part onto the positive orthant, which is equivalent to ReLU activation in neural networks. Therefore, the sequence of iterations for solving a QP problem resembles a single-layer recurrent neural network (RNN) with weights dependent on the QP parameters (P, q, H, b) , followed by ReLU activation. In the setting where the QP parameters are learnable, as what follows in Section 3, this resemblance facilitates the end-to-end gradient-based tuning of the QP parameters driven by a cost that depends on the solution of the QP.

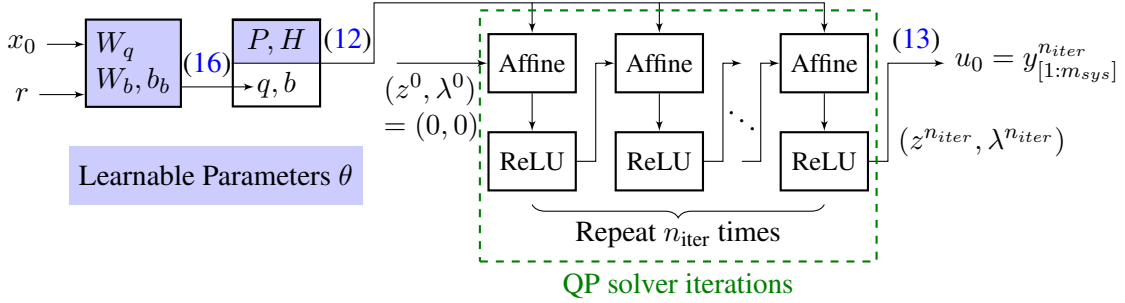


Figure 1: Proposed control policy architecture. The controller solves a QP problem in form (7), whose parameters P, H are shared across all initial state and reference (x_0, r) , while q, b depend affinely on (x_0, r) with weights W_q, W_b and biases b_q, b_b (see (16)). An approximate solution to the QP problem, $y^{n_{iter}}$, is obtained by running a n_{iter} QP solver iterations (11) followed by a transform (13), whose first m_{sys} dimensions are used as the current control input u_0 .

3. Learning Model-Free QP Controllers

This section introduces a framework to learn control laws based on the QP problem (7) through a reinforcement learning approach. MPC, which derives the parameters (P, q, H, b) from a model-based prediction over a finite horizon, may face limitations in control performance, due to short-sightedness (Erez et al., 2012) and model inaccuracies (Forbes et al., 2015). Furthermore, from a computational perspective, the derived QP problem may require a large number of iterations to solve. Our approach seeks to address these challenges by enabling the direct learning of these parameters, thereby bypassing the restrictions of model-based prediction. The policy architecture facilitating this learning process is shown in Figure 1.

The high level idea of our method is to parameterize our control policy—the mapping from x_0 and r to the control action u_0 , using the QP iterative algorithm as discussed in the previous section. This parameterization is represented by an RNN, which could be learned by applying deep reinforcement learning methods. Specifically, here we highlight several critical design components of the policy architecture. The first two are regarding the parameterization of the policy, i.e., *what* to learn:

State-independent matrices P, H : Note from (9) that for the MPC controller, the matrices P, H holds the same across different initial and reference state (x_0, r) 's. Motivated by this fact, the matrices P, H are also state-independent in the proposed policy architecture, i.e., only one matrix P and one matrix H needs to be learned for a specific system. Additionally, to ensure the positive definiteness of P , we use the factor L_P in the Cholesky decomposition $P = L_P L_P^\top$ instead of the matrix P itself as the learnable parameter, and force its diagonal elements to be positive via a softplus activation (Zheng et al., 2015), a commonly applied trick for learning positive definite matrices (Haarnoja et al., 2016; Lutter et al., 2019).

Affine transformations yielding vectors q, b : Drawing inspirations from MPC (9), we restrict the vectors q to depend linearly on the current state x_0 and the reference state r , and the vector b to depend affinely on x_0 :

$$q(x_0, r; W_q) = W_q [x_0^\top \ r^\top]^\top, \quad b(x_0, r; W_b, b_b) = W_b x_0 + b_b, \quad (16)$$

where W_q, W_b (resp. b_b) are learnable matrices (resp. vector) of proper dimensions.

The above described parameterization strategy ensures that when the chosen problem dimensions n_{qp}, m_{qp} match the dimensions of the QP translated from MPC, then the MPC policy is within

the family of parameterized policies defined by the proposed architecture. In other words, the proposed controller can be viewed as a generalization of MPC. Meanwhile, the state-independence / state-affineness restrictions on the problem significantly narrow down the class of policies compared to MLP policies or MPC-akin policies with generic function approximator components (Amos et al., 2018), a crucial restriction that diminishes the number of learnable parameters as well as facilitates the theoretical guarantees.

Another two design components determine *how* the parameters are learned:

Unrolling with a fixed number of iterations: To solve the QP problem and differentiate the solution with respect to the problem parameters, we deploy a fixed number n_{iter} of QP solver iterations described in Section 2.3, and differentiate through the computational path of these iterations, a practice known as unrolling (Monga et al., 2021). Unlike implicit differentiation methods (Amos and Kolter, 2017; Amos et al., 2018; Agrawal et al., 2019), which differentiate through the optimality condition and hence requires the forward pass of the solver to reach the stationary point, our method directly differentiates the solution after n_{iter} iterations, and can obtain a correct gradient even if the stationary point is not reached within these iterations. According to our empirical results, a small number of iterations would suffice for good control performance (e.g., $n_{iter} = 10$), which mitigates the computational burden of the unrolling process. Intuitively, the sufficiency of a small n_{iter} can be accredited to the model-free nature of the proposed method, which, by discarding the restrictions imposed by model-based prediction (see (10)), gains the flexibility to learn a QP problem that not only optimizes the controller performance, but also is easy to solve.

Reinforcement learning with residual minimization: The control policy described above, parameterized by $\theta = (L_P, H, W_q, b_q, W_b, b_b)$, can serve as a drop-in replacement for standard policy networks, and be optimized using various off-the-shelf policy-based or actor-critic RL algorithms, such as PPO (Schulman et al., 2017), SAC (Haarnoja et al., 2018) and DDPG (Lillicrap et al., 2015). However, apart from the standard RL loss, we also include a regularization term for minimizing the residuals given by the QP solver embedded in the policy. Given a dataset \mathcal{D} of transition samples, it is defined as follows:

$$\ell_{res}(\theta; \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{k=1}^{|\mathcal{D}|} \|Hy_k^{n_{iter}} + b_k - z_k^{n_{iter}}\|_2^2 + \|Py_k^{n_{iter}} + q_k + H^\top \lambda_k^{n_{iter}}\|_2^2, \quad (17)$$

which, motivated by the result stated in Theorem 1 that small residuals are indicative of near-optimality, encourages the learned QP problems to be easy to solve. From above, the procedure of policy learning using an RL algorithm is shown in Algorithm 1.

4. Performance Guarantees of Learned QP Controller

In this section, we propose a method for establishing performance guarantees of a learned QP controller with the architecture described in Section 3. We provide sufficient conditions for persistent feasibility and asymptotic stability of the closed-loop system under a QP controller, which parallel the theoretical guarantees for linear MPC (Borrelli et al., 2017). For simplicity, we consider the stabilization around the origin, i.e., $r = 0$, but the method of analysis can be extended to the general case. Additionally, we assume throughout the section that the optimal solution of the learned QP problem is attained, which can be ensured by allowing the QP solver to run sufficient iterations until convergence when deploying.

Algorithm 1: Framework of Learning of QP Controllers

Input: Simulation environment Env with nominal dynamics (2), RL algorithm RL , policy architecture π_θ shown in Fig. 1, regularization coefficient ρ_{res}

Output: Optimized policy parameters $\theta = (L_P, H, W_q, b_q, W_b, b_b)$

```

1 for  $epoch = 1, 2, \dots$  do
2   | Interact with  $Env$  using current policy  $\pi_\theta$  to collect a dataset  $\mathcal{D}$ 
3   | Compute RL loss, denoted by  $\ell_{RL}(\theta; \mathcal{D})$ 
4   | Compute residual loss  $\ell_{res}(\theta; \mathcal{D})$  using (17)
5   | Update  $\theta$  according to the loss  $\ell_{RL}(\theta; \mathcal{D}) + \rho_{res}\ell_{res}(\theta; \mathcal{D})$ 
6 end
    
```

Denote the property under consideration as \mathcal{P} . Suppose that a certificate to \mathcal{P} , given the initial state x_0 is in a polytopic set \mathcal{X}_0 , can be written in the following form:

$$\min_{x_0 \in \mathcal{X}_0, u_0, \nu} \{f(x_0, u_0, \nu) | g(x_0, u_0, \nu) \leq 0, u_0 = \pi_\theta(x_0)\} \geq 0 \Rightarrow \mathcal{P} \text{ holds when } x_0 \in \mathcal{X}_0, \quad (18)$$

where π_θ denotes the θ -parameterized control policy described in Section 3, ν is an auxiliary variable, and f, g are quadratic (possibly nonconvex) functions. The optimization problem in the LHS of (18) can be expressed as a bilevel problem by explicitly expanding the control policy π_θ as:

$$\pi_\theta(x_0) = y_{[1:m_{sys}]}^*, y^* \in \arg \min \left\{ (1/2)y^\top P y + q^\top y \mid H y + b \geq 0 \right\},$$

where $q = W_q x_0, b = W_b x_0 + b_b$. (19)

Replacing the inner-level problem in (19) by its KKT condition, the verification problem in (18) can be cast into a nonconvex Quadratically Constrained Quadratic Program (QCQP) with variables x_0, ν, y, μ . Various computationally tractable methods for lower bounding the optimal value of a QCQP are available, such as Lagrangian relaxation (d'Aspremont and Boyd, 2003) and the method of moments (Lasserre, 2001), and once a nonnegative lower bound is obtained, the property \mathcal{P} is verified.

Verification of persistent feasibility and asymptotic stability both fall into the framework described above. The conclusions are stated as follows:

Theorem 2 (Certificate for Persistent Feasibility) *The control policy (19) if persistently feasible (i.e., gives a valid control input that keeps the next state inside the bounds at every step) for all initial states $x_0 \in \mathcal{X}_0 = \{x | Gx \leq c\}$, if the optimal value of the following nonconvex QCQP is nonnegative:*

$$\begin{aligned} & \underset{x_0, \nu, y, \mu}{\text{minimize}} && -\nu^\top (G(Ax_0 + By_{[1:m_{sys}]}) - c), \\ & \text{subject to} && Gx_0 \leq c, \nu \geq 0, \mathbf{1}^\top \nu = 1, \\ & && Py + W_q x_0 - H^\top \mu = 0, Hy + W_b x_0 + b_b \geq 0, \mu \geq 0, \mu^\top (Hy + W_b x_0 + b_b) = 0. \end{aligned}$$

To certify asymptotic stability, we consider the Lyapunov function of a stabilizing baseline MPC, and attempt to show that the Lyapunov function decreases along all trajectories even if the

learned QP controller is deployed instead of the baseline MPC. A similar technique has been applied to the stability analysis of approximate MPC (Schwan et al., 2023). To formalize this idea, we define the following notations: $l(x, u) = x^\top Qx + u^\top Ru$ is the stage cost; the baseline MPC policy has horizon N , terminal constraint $x_N \in \mathcal{X}_f$ and terminal cost $V_N(x_N)$; the function $J(x_0, u_{0:N-1}) = \sum_{k=0}^{N-1} l(x_k, u_k) + V_N(x_N)$, where $x_{k+1} = Ax_k + Bu_k$, is the objective function of the baseline MPC. To ensure that the baseline MPC is stabilizing as long as it is feasible, one can choose \mathcal{X}_f to be an invariant set under a stabilizing linear feedback controller $u = Kx$, and $V_N(x)$ to be the cost-to-go under $u = Kx$. Based on these notations, a certificate for asymptotic stability can be stated as follows:

Theorem 3 (Certificate for Asymptotic Stability) *The closed-loop system under the control policy (19) is asymptotically stable for all initial states $x_0 \in \mathcal{X}_0 = \{x | Gx \leq c\}$, if there exists $\epsilon > 0$ and $N \in \mathbb{N}^*$, such that the optimal value of the following nonconvex QCQP is nonnegative:*

$$\begin{aligned} & \underset{x, \bar{u}_{0:N-1}, y, \mu}{\text{minimize}} && J(x_0, \bar{u}_{0:N-1}) + l(x_0, \bar{u}_0) - J(x_0, (y_{[1:m_{sys}]}, \bar{u}_{1:N-1})) - \epsilon \|x_0\|^2, \\ & \text{subject to} && Gx_0 \leq c, x_N(x_0, \bar{u}_{0:N-1}) \in \mathcal{X}_f, \\ & && Py + W_q x_0 - H^\top \mu = 0, Hy + W_b x_0 + b_b \geq 0, \mu \geq 0, \mu^\top (Hy + W_b x_0 + b_b) = 0. \end{aligned}$$

5. Benchmarking Results

In our empirical evaluations, we aim to answer the following questions:

- How does the learned QP controller compare with common baselines (MPC, RL-trained MLP) on nominal linear systems?
- Can the learned QP controller handle modeling inaccuracies and disturbances?
- Does the method generalize to real-world robot systems with modeling inaccuracy and non-linearity?

We briefly describe the experimental setup and typical results in this section, with complete details on systems, setup, hyperparameters, baseline definitions, and additional results in the supplementary materials. Code is available at <https://github.com/yiwenlu66/learning-qp>.

5.1. Results on Nominal Systems

We compare the Learned QP (LQP) controller with MPC and MLP baselines on benchmark systems like the quadruple tank (Johansson, 2000) and cartpole (Geva and Sitte, 1993), generating random initial states and references across 10^4 trials. For MPC, we evaluate variants with and without manually tuned terminal costs over short (2 steps) and long (16 steps) horizons, all implemented using OSQP (Stellato et al., 2020), a solver known for its efficiency in MPC applications (Forgione et al., 2020), with default solver configurations. Both LQP and MLP are trained using PPO, maintaining consistent reward definitions and RL hyperparameters. We incrementally increase the MLP size until further increases yield negligible performance improvements, selecting this size for comparison. The LQP is assessed in both small ($n_{qp} = 4, m_{qp} = 24$) and large ($n_{qp} = 16, m_{qp} = 96$) configurations, approximately aligning with the QP problem sizes from short- and long-horizon MPC. All training (including simulation and policy update) are performed on a single NVIDIA RTX 4090 GPU, with the small and large configurations taking 1.2 hours and 2.7 hours respectively.

The results of the benchmarking experiments are summarized in Table 1. In terms of control performance, LQP demonstrates comparable effectiveness to both MPC and MLP baselines. A

Table 1: Performance comparison on benchmark systems. Methods: MPC(N) = naive MPC (Problem 2) with horizon N ; MPC-T(N) = MPC with horizon N and manually tuned terminal cost; RL-MLP = reinforcement learning controller with MLP policy; LQP(n_{qp}, m_{qp}) = proposed learned QP controller with problem dimensions (n_{qp}, m_{qp}). Metrics: Fail% = percentage of early-terminated trials due to constraint violation; Cost = average LQ cost until termination; P-Cost = average cost with penalty for constraint violation; FLOPs = floating point operations per control step (reported as median_{+(max−median)} for variable data); #Params = number of learnable policy parameters. Best is highlighted in **bold**, and second best is underlined.

Method	Quadruple Tank					Cartpole Balancing				
	Fail%	Cost	P-Cost	FLOPs	#Params	Fail%	Cost	P-Cost	FLOPs	#Params
MPC(2)	16.59	236.1	275.7	95K _{+1.2M}	-	100.0	1.36	129	67K _{+814K}	-
MPC(16)	4.27	228.3	237.2	22M _{+52M}	-	46.86	<u>0.34</u>	8.39	3.9M _{+47M}	-
MPC-T(2)	4.23	239.6	248.5	470K _{+779K}	-	100.0	1.41	122	89K _{+792K}	-
MPC-T(16)	3.22	224.8	231.5	26M _{+49M}	-	4.74	0.30	0.79	51M _{+50M}	-
RL-MLP	0.03	266.7	266.7	<u>23K</u>	11K	3.23	0.57	0.91	<u>87K</u>	43K
LQP(4, 24)	0.18	272.5	272.8	14K	0.3K	<u>3.49</u>	0.76	1.12	14K	0.2K
LQP(16, 96)	<u>0.13</u>	<u>227.3</u>	227.6	208K	<u>2.6K</u>	4.11	0.44	<u>0.87</u>	208K	<u>2.2K</u>

benefit of LQP is its independence from manual tuning of the terminal cost, a necessity in MPC methods.

Regarding computational efficiency, LQP stands out for its minimal demand for achieving similar control performance. This efficiency stems from LQP’s fixed number of unrolled QP solver iterations. While MPC’s computation cost varies based on implementation, the light computation of LQP is still noteworthy, especially in scenarios with tight computational limits. For example, the LQP(4, 24) configuration, despite having lowest FLOPs among all methods, still manages acceptable control performance.

Finally, in terms of the number of learnable policy parameters, LQP requires substantially fewer than the RL-MLP. This not only hints at LQP’s suitability for memory-constrained embedded systems, but also opens up possibilities for online few-shot learning, which can be left further exploration.

5.2. Validation of Robustness

This subsection is concerned with the robustness of the learned QP controller against modeling inaccuracies and disturbances. Instead of the nominal dynamics (2), we now consider the following perturbed dynamics:

$$x_{k+1} = (A + \Delta A)x_k + (B + \Delta B)u_k + w_k,$$

where $\Delta A, \Delta B$ are parametric uncertainties, and w_k is a disturbance. LQP and MLP are trained using domain randomization (Tobin et al., 2017; Mehta et al., 2020), where the simulator randomly sample these uncertain components during training. Robust MPC baselines, including tube MPC (Mayne et al., 2005) and

Table 2: Performance comparison on quadruple tank system with process noise and parametric uncertainties. Notations are similar to those in the caption of Table 1. Computation time instead of FLOPs per control step is used as the metric for computational efficiency since it is difficult to obtain the exact FLOPs from the robust MPC baselines.

Method	Fail%	Cost	P-Cost	Time(s)	#Params
MPC-T(16)	82.6	216.8	713.4	0.25 _{+0.56}	-
Tube	81.9	<u>233.3</u>	597.9	2.22 ₊₄₄	-
Scenario	<u>16.4</u>	236.9	273.2	5.21 ₊₁₈	-
RL-MLP	40.9	272.1	387.5	3×10^{-5}	11K
LQP(4, 24)	1.4	253.6	<u>256.4</u>	2×10^{-5}	0.3K
LQP(16, 96)	1.4	240.7	243.4	3×10^{-4}	<u>2.6K</u>

scenario MPC (Bernardini and Bemporad, 2009) implemented by the do-mpc toolbox (Fiedler et al., 2023), are included for comparison.

The results in Table 2 highlight LQP’s enhanced robustness, as it achieves the highest success rate and lowest constraint-violation-penalized cost among all methods tested. Also, it requires significantly less online computation compared to robust MPC methods, benefitting from domain randomization known for its effectiveness in empirical RL and robotics (Loquercio et al., 2019; Margolis et al., 2022). Additionally, LQP shows better generalizability than RL-MLP, potentially due to its structural design, although further empirical analysis is necessary to fully understand these benefits.

5.3. Application Example on a Real-World System: Vehicle Drift Maneuvering

LQP is also evaluated on a challenging real-world control task, namely, the drift maneuvering of a 1/10 scale RC car, similar to the problem studied in Yang et al. (2022); Domberg et al. (2022); Lu et al. (2023). The objective is to track the yaw rate, side slip angle, and velocity references, such that the car enters and maintains a drifting state. Despite the high nonlinearity of the system, the proposed controller formally introduced on linear systems successfully generalizes to this task. As shown in Fig. 2, the learned QP controller can successfully track the references and maintain the drifting state, achieving similar performance to previously reported RL-trained MLP methods on the same task (Domberg et al., 2022).

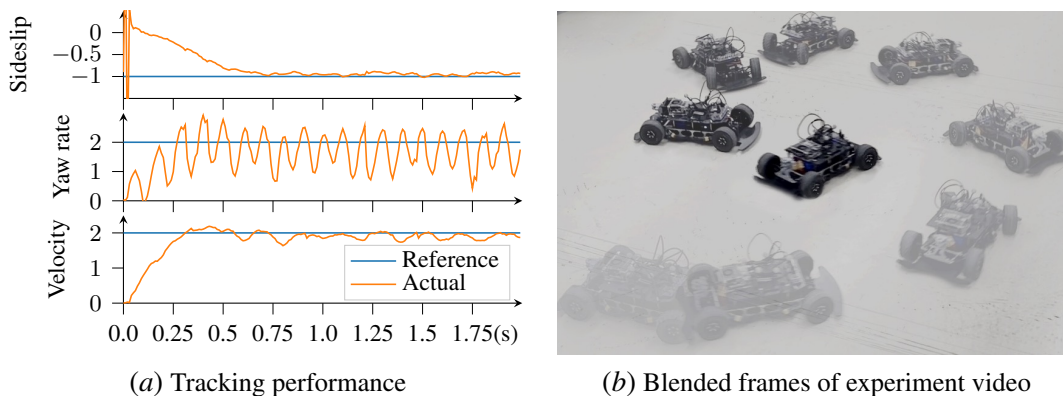


Figure 2: Result of deploying learned QP controller to the vehicle drift maneuvering task.

6. Conclusion

This work presents a novel class of Quadratic Programming (QP) controllers inspired by Model Predictive Control (MPC). The proposed controllers not only retain the theoretical guarantees akin to MPC, but also exhibit desirable empirical performance and computational efficiency. Benchmarks including applications in real-world scenarios like vehicle drift maneuvering, further validate the effectiveness and robustness of our approach.

Notably, the unrolled QP solver is structured similarly to a deep neural network, indicating the suitability of the proposed policy architecture as a drop-in replacement for standard policy networks in RL. This opens up possibilities for combining the architecture with various RL methods, such as meta-learning (Finn et al., 2017) and safety-constrained RL (Achiam et al., 2017; Yu et al., 2022), which can be left for future investigation.

References

- Shokhjakon Abdufattokhov, Mario Zanon, and Alberto Bemporad. Learning convex terminal costs for complexity reduction in mpc. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pages 2163–2168. IEEE, 2021.
- Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In *International conference on machine learning*, pages 22–31. PMLR, 2017.
- Adrian Agogino, Ritchie Lee, and Dimitra Giannakopoulou. Challenges of explaining control. In *2nd ICAPS Workshop on Explainable Planning (XAIP’19)*, 2019.
- Akshay Agrawal, Brandon Amos, Shane Barratt, Stephen Boyd, Steven Diamond, and J Zico Kolter. Differentiable convex optimization layers. *Advances in neural information processing systems*, 32, 2019.
- Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks. In *International Conference on Machine Learning*, pages 136–145. PMLR, 2017.
- Brandon Amos, Ivan Jimenez, Jacob Sacks, Byron Boots, and J Zico Kolter. Differentiable mpc for end-to-end planning and control. *Advances in neural information processing systems*, 31, 2018.
- Daniele Bernardini and Alberto Bemporad. Scenario-based model predictive control of stochastic constrained linear systems. In *Proceedings of the 48th IEEE Conference on Decision and Control (CDC) held jointly with 2009 28th Chinese Control Conference*, pages 6333–6338. IEEE, 2009.
- Francesco Borrelli, Alberto Bemporad, and Manfred Morari. *Predictive control for linear and hybrid systems*. Cambridge University Press, 2017.
- Florian D Brunner, Mircea Lazar, and Frank Allgöwer. Stabilizing model predictive control: On the enlargement of the terminal set. *International Journal of Robust and Nonlinear Control*, 25(15): 2646–2670, 2015.
- Antonin Chambolle and Thomas Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of mathematical imaging and vision*, 40:120–145, 2011.
- Hongkai Dai, Benoit Landry, Lujie Yang, Marco Pavone, and Russ Tedrake. Lyapunov-stable neural-network control. *arXiv preprint arXiv:2109.14152*, 2021.
- Alexandre d’Aspremont and Stephen Boyd. Relaxations and randomized methods for nonconvex qcqps. *EE392o Class Notes, Stanford University*, 1:1–16, 2003.
- Vishnu R Desraj and Nathan Michael. Experience-driven predictive control. *Robot Learning and Planning (RLP 2016)*, page 29, 2016.
- Fabian Domberg, Carlos Castelar Wemmers, Hiren Patel, and Georg Schildbach. Deep drifting: Autonomous drifting of arbitrary trajectories using deep reinforcement learning. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 7753–7759. IEEE, 2022.

- Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. In *International conference on machine learning*, pages 1329–1338. PMLR, 2016.
- Peter Englert, Ngo Anh Vien, and Marc Toussaint. Inverse kkt: Learning cost functions of manipulation tasks from demonstrations. *The International Journal of Robotics Research*, 36(13-14): 1474–1488, 2017.
- Tom Erez, Yuval Tassa, and Emanuel Todorov. Infinite-horizon model predictive control for periodic tasks with contacts. *Robotics: Science and systems VII*, page 73, 2012.
- Felix Fiedler, Benjamin Karg, Lukas Lüken, Dean Brandner, Moritz Heinlein, Felix Brabender, and Sergio Lucia. do-mpc: Towards fair nonlinear and robust model predictive control. *Control Engineering Practice*, 140:105676, 2023.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- Michael G Forbes, Rohit S Patwardhan, Hamza Hamadah, and R Bhushan Gopaluni. Model predictive control in industry: Challenges and opportunities. *IFAC-PapersOnLine*, 48(8):531–538, 2015.
- Marco Forgione, Dario Piga, and Alberto Bemporad. Efficient calibration of embedded MPC. In *Proc. of the 21st IFAC World Congress 2020, Berlin, Germany, July 12-17 2020*, 2020.
- Shlomo Geva and Joaquin Sitte. A cartpole experiment benchmark for trainable controllers. *IEEE Control Systems Magazine*, 13(5):40–51, 1993.
- David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- Tuomas Haarnoja, Anurag Ajay, Sergey Levine, and Pieter Abbeel. Backprop kf: Learning discriminative deterministic state estimators. *Advances in neural information processing systems*, 29, 2016.
- Tuomas Haarnoja, Sehoon Ha, Aurick Zhou, Jie Tan, George Tucker, and Sergey Levine. Learning to walk via deep reinforcement learning. *arXiv preprint arXiv:1812.11103*, 2018.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- Nicklas Hansen, Hao Su, and Xiaolong Wang. Td-mpc2: Scalable, robust world models for continuous control. *arXiv preprint arXiv:2310.16828*, 2023.
- Lukas Hewing, Juraj Kabzan, and Melanie N Zeilinger. Cautious model predictive control using gaussian process regression. *IEEE Transactions on Control Systems Technology*, 28(6):2736–2743, 2019.
- Lukas Hewing, Kim P Wabersich, Marcel Menner, and Melanie N Zeilinger. Learning-based model predictive control: Toward safe learning in control. *Annual Review of Control, Robotics, and Autonomous Systems*, 3:269–296, 2020.

- Tobias Johannink, Shikhar Bahl, Ashvin Nair, Jianlan Luo, Avinash Kumar, Matthias Loskyll, Juan Aparicio Ojea, Eugen Solowjow, and Sergey Levine. Residual reinforcement learning for robot control. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 6023–6029. IEEE, 2019.
- Karl Henrik Johansson. The quadruple-tank process: A multivariable laboratory process with an adjustable zero. *IEEE Transactions on control systems technology*, 8(3):456–465, 2000.
- Jean B Lasserre. Global optimization with polynomials and the problem of moments. *SIAM Journal on optimization*, 11(3):796–817, 2001.
- Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 62, 2022.
- Zhongyu Li, Xuxin Cheng, Xue Bin Peng, Pieter Abbeel, Sergey Levine, Glen Berseth, and Koushil Sreenath. Reinforcement learning for robust parameterized locomotion control of bipedal robots. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2811–2817. IEEE, 2021.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Antonio Loquercio, Elia Kaufmann, René Ranftl, Alexey Dosovitskiy, Vladlen Koltun, and Davide Scaramuzza. Deep drone racing: From simulation to reality with domain randomization. *IEEE Transactions on Robotics*, 36(1):1–14, 2019.
- Yiwen Lu, Bo Yang, Jiayun Li, Yihan Zhou, Hongshuai Chen, and Yilin Mo. Consecutive inertia drift of autonomous rc car via primitive-based planning and data-driven control. In *2023 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2023.
- Michael Lutter, Kim Listmann, and Jan Peters. Deep lagrangian networks for end-to-end learning of energy-based control for under-actuated systems. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7718–7725. IEEE, 2019.
- Gabriel B Margolis, Ge Yang, Kartik Paigwar, Tao Chen, and Pulkit Agrawal. Rapid locomotion via reinforcement learning. *arXiv preprint arXiv:2205.02824*, 2022.
- David Q Mayne, María M Seron, and SV Raković. Robust model predictive control of constrained linear systems with bounded disturbances. *Automatica*, 41(2):219–224, 2005.
- Bhairav Mehta, Manfred Diaz, Florian Golemo, Christopher J Pal, and Liam Paull. Active domain randomization. In *Conference on Robot Learning*, pages 1162–1176. PMLR, 2020.
- Marcel Menner, Peter Worsnop, and Melanie N Zeilinger. Constrained inverse optimal control with application to a human manipulation task. *IEEE Transactions on Control Systems Technology*, 29(2):826–834, 2019.
- Vishal Monga, Yuelong Li, and Yonina C Eldar. Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing. *IEEE Signal Processing Magazine*, 38(2):18–44, 2021.

- Manfred Morari and Jay H Lee. Model predictive control: past, present and future. *Computers & chemical engineering*, 23(4-5):667–682, 1999.
- Tianwei Ni, Benjamin Eysenbach, and Ruslan Salakhutdinov. Recurrent model-free rl can be a strong baseline for many pomdps. *arXiv preprint arXiv:2110.05038*, 2021.
- Pavel Osinenko, Dmitrii Dobriborsci, and Wolfgang Aumer. Reinforcement learning with guarantees: a review. *IFAC-PapersOnLine*, 55(15):123–128, 2022.
- S Joe Qin and Thomas A Badgwell. A survey of industrial model predictive control technology. *Control engineering practice*, 11(7):733–764, 2003.
- Ugo Rosolia and Francesco Borrelli. Learning model predictive control for iterative tasks. a data-driven control framework. *IEEE Transactions on Automatic Control*, 63(7):1883–1896, 2017.
- Nikita Rudin, David Hoeller, Philipp Reist, and Marco Hutter. Learning to walk in minutes using massively parallel deep reinforcement learning. In *Conference on Robot Learning*, pages 91–100. PMLR, 2022.
- Ernest K Ryu and Wotao Yin. *Large-scale convex optimization: algorithms & analyses via monotone operators*. Cambridge University Press, 2022.
- Yahya Sattar and Samet Oymak. Quickly finding the best linear model in high dimensions via projected gradient descent. *IEEE Transactions on Signal Processing*, 68:818–829, 2020.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Roland Schwan, Colin N Jones, and Daniel Kuhn. Stability verification of neural network controllers using mixed-integer programming. *IEEE Transactions on Automatic Control*, 2023.
- Max Schwenzer, Muzaffer Ay, Thomas Bergs, and Dirk Abel. Review on model predictive control: An engineering perspective. *The International Journal of Advanced Manufacturing Technology*, 117(5-6):1327–1349, 2021.
- Raffaele Soloperto, Matthias A Müller, Sebastian Trimpe, and Frank Allgöwer. Learning-based robust model predictive control with state-dependent uncertainty. *IFAC-PapersOnLine*, 51(20):442–447, 2018.
- Mario Srouji, Jian Zhang, and Ruslan Salakhutdinov. Structured control nets for deep reinforcement learning. In *International Conference on Machine Learning*, pages 4742–4751. PMLR, 2018.
- Bartolomeo Stellato, Goran Banjac, Paul Goulart, Alberto Bemporad, and Stephen Boyd. Osqp: An operator splitting solver for quadratic programs. *Mathematical Programming Computation*, 12(4):637–672, 2020.
- Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE, 2017.

- Zhaoming Xie, Glen Berseth, Patrick Clary, Jonathan Hurst, and Michiel van de Panne. Feedback control for cassie with deep reinforcement learning. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1241–1246. IEEE, 2018.
- Bo Yang, Yiwen Lu, Xu Yang, and Yilin Mo. A hierarchical control framework for drift maneuvering of autonomous vehicles. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 1387–1393. IEEE, 2022.
- Dongjie Yu, Haitong Ma, Shengbo Li, and Jianyu Chen. Reachability constrained reinforcement learning. In *International Conference on Machine Learning*, pages 25636–25655. PMLR, 2022.
- Hao Zheng, Zhanlei Yang, Wenju Liu, Jizhong Liang, and Yanpeng Li. Improving deep neural networks using softplus units. In *2015 International joint conference on neural networks (IJCNN)*, pages 1–4. IEEE, 2015.
- Ruikun Zhou, Thanin Quartz, Hans De Sterck, and Jun Liu. Neural lyapunov control of unknown nonlinear systems with stability guarantees. *Advances in Neural Information Processing Systems*, 35:29113–29125, 2022.
- Vrushabh Zinage and Efstathios Bakolas. Neural koopman lyapunov control. *Neurocomputing*, 527:174–183, 2023.