

# DiffLoop: Tuning PID Controllers by Differentiating Through the Feedback Loop

Athindran Ramesh Kumar  
Department of Electrical Engineering  
Princeton University  
Princeton, USA  
arkumar@princeton.edu

Peter J. Ramadge  
Department of Electrical Engineering  
Princeton University  
Princeton, USA  
ramadge@princeton.edu

**Abstract**—Since most industrial control applications use PID controllers, PID tuning and anti-windup measures are significant problems. This paper investigates tuning the feedback gains of a PID controller via back-calculation and automatic differentiation tools. In particular, we episodically use a cost function to generate gradients and perform gradient descent to improve controller performance. We provide a theoretical framework for analyzing this non-convex optimization and establish a relationship between back-calculation and disturbance feedback policies. We include numerical experiments on linear systems with actuator saturation to show the efficacy of this approach.

## I. INTRODUCTION

PID controllers are the most popular form of feedback control in industrial applications [5]. In general, the PID gains need to be tuned to obtain good performance. In addition, potential actuator saturation must be taken into account, since saturation can induce integrator wind-up, resulting in unexpected transients during operation. To address the above issues, PID tuning [11] can be performed using classical control based on a model of the plant devised either from prior knowledge or system identification. Model-free [28], [32] PID tuning has also been explored by performing selective experiments to tune the gains gradually. Notably, machine learning approaches such as [21], [23], [31] tune the controller parameters without a system model. In practice, we usually have some prior knowledge of the system and can use it to obtain a coarse system model. The results in [25] show that a coarse model can be useful in obtaining an initial robust controller. Indeed, given a coarse model for a linear system, we can apply several controller design techniques, both classical (PID, loop shaping) and modern (LQR,  $H^\infty$ ). However, when non-linearities induced by actuator saturation are present, simple controller design techniques are not available. In this setting, the back-calculation method [5], [29] is the simplest anti-windup technique.

We explore using the back-calculation method together with a non-convex optimization approach based on differentiating through the system model, actuator, and the feedback loop to tune the feedback gains. This approach is inspired by recent work on differentiable physics engines/models [12], [13]. In particular, we episodically tune the controller parameters by simulating with the current parameters, evaluating the cost, and performing gradient descent on the cost objective. By

propagating the gradients over time for the entire simulation, we capture the controller parameters' long-term dependencies on the system's dynamics. The use of automatic differentiation tools, e.g., in TensorFlow [1], PyTorch [24], or as stand-alone code [8], [22], enable us to easily compute required gradients on-the-fly without using analytical techniques. The computational cost scales with the simulation horizon. For reasonable time horizons, the method can be efficiently implemented on a modern CPU.

The non-convex optimization problem of interest is posed as an output-feedback controller design with augmented state and marginally stabilizable, detectable dynamics. This framework is developed by treating the errors due to actuator saturation as an additional disturbance and predicting future disturbances using disturbance estimates at previous time-steps. Using this approach, we can show an equivalence between the back-calculation method and the class of disturbance feedback controllers introduced in [2]. We leave the convergence analysis of gradient descent for this optimization problem to future work. Here we perform simulations on four different systems with saturation to show the efficacy of the approach.

We discuss our work's relation to the existing literature in section II. Then in section III, we introduce a framework to analyze the optimization's convergence properties. We perform numerical experiments to illustrate our approach in section IV, and conclude in section V.

## II. RELATED WORK

### *Machine learning for PID tuning and anti-windup design*

PID tuning is a well-studied problem, see e.g., [4], [28], [32]. Much of the previous work has focused on black-box optimization, reinforcement learning, and model-free tuning. We briefly discuss these approaches below.

Black-box PID tuning includes approaches such as genetic algorithms [18], [23], particle swarm optimization [10], and model-free decision trees [31]. These do not incorporate an explicit model of the system. A distinct approach to PID tuning uses reinforcement learning. Here, both model-based and model-free approaches are explored. For example, the authors of [14] use a model-based reinforcement learning approach with a Gaussian Process model to tune the PID gains for a seven degree-of-freedom robot arm. However,

their algorithm ignores actuator saturation. Reference [7] tunes fuzzy PD and PI controllers using the Q-learning algorithm. Similarly, [26] uses the model-free Q-learning algorithm to tune the PID gains for a cart-pole system. Also, [21] uses an actor-critic RL algorithm with a neural Q-function to tune the PID parameters with back-calculation.

In addition to tuning the PID parameters, we want to address actuator saturation. Anti-windup compensation has been tackled by both classical [5], [6], [29] and modern control [16] techniques. Among these approaches, the back-calculation method [5] is a simple and effective scheme for compensating for integrator wind-up. It requires an actuator model with knowledge of the saturation limits. Even though the scheme lacks the formal stability and robustness guarantees of modern anti-windup design, its simplicity is appealing. Also relevant is the idea of treating the errors due to actuator saturation as a disturbance [19]. We adopt this approach.

### Differentiable models

There is recent work in building physics engines that can be differentiated through to update model parameters or train controllers [3], [9], [12], [13]. Notably, in [12], the differentiable model approach is effective even in learning controllers from raw images. Our work is complementary in two aspects. First, before extending the differentiable physics approach to complex problems, it would be interesting to understand its performance on simple systems with purely saturation non-linearities. If this is effective, then second, it would be of interest to tune PID controllers (with saturation) using standard automatic differentiation tools.

### Theoretical machine learning and control

Recent growth at the intersection of machine learning theory and control theory has enabled better understanding of some long-standing problems. A class of disturbance feedback policies is introduced and optimized using a convex relaxation in [2], [17], [27]. In [25], it has been shown that model-based control approaches are often sample-efficient compared to model-free reinforcement learning for linear systems control. References [15], [30], show the convergence of gradient descent and policy optimization for non-convex state feedback LQR and  $H_2$ ,  $H_\infty$  controller design. However, the output feedback case is less studied. In our work, we seek a model-based approach applicable to systems with saturation. Actuator saturation is one of the most common forms of non-linearity in practical control systems. We explore a relationship between the back-calculation anti-windup method and disturbance feedback policies.

## III. DISTURBANCE FEEDBACK FOR ANTI-WINDUP COMPENSATION

Assume the system to be controlled has a stabilizable and detectable state space representation:

$$\begin{aligned} x_{t+1} &= Ax_t + Bu_t + w_t \\ y_t &= Cx_t + e_t. \end{aligned} \quad (1)$$

Here  $u_t \in \mathbb{R}^m$ ,  $x_t \in \mathbb{R}^n$ , and  $y_t \in \mathbb{R}^p$  are the input, state, and output at time  $t$ , respectively. To model actuator saturation, we modify (1) to:

$$x_{t+1} = Ax_t + B\text{sat}(u_t) + w_t \quad (2)$$

### A. Back-calculation method

In the back-calculation method [5], the errors due to actuator saturation are integrated and fed back to prevent windup. Let  $r_t$  denote the reference signal to be tracked, and  $P_t$ ,  $I_t$ , and  $D_t$  denote the proportional, integral, and derivative signal components of the PID controller, respectively. Then a standard back-calculation PID controller is given by:

$$\begin{aligned} P_t &= k_p (r_t - y_t) \\ D_t &= \alpha D_{t-1} + k_d \Delta y_t \\ I_{t+1} &= I_t + k_i (r_t - y_t) + b(\text{sat}(v_t) - v_t) \\ v_t &= P_t + I_t + D_t \\ \text{sat}(v_t) &= \text{clamp}(v_t, u_{\text{low}}, u_{\text{high}}). \end{aligned} \quad (3)$$

Here  $\Delta$  is the difference operator,  $\alpha$  is a filter parameter,  $k_p$ ,  $k_i$ ,  $k_d$  and  $b$  are the proportional, integral, derivative and back-calculation feedback gains. Finally,  $u_{\text{low}}$  (resp.  $u_{\text{high}}$ ) is the minimum (resp. maximum) actuator output. Here,  $r_t$  is just a reference operating point and can be set to 0 for analysis.

### B. Disturbance feedback policies and back-calculation

We now connect the back-calculation technique to disturbance feedback policies. To do so, we formulate a controller design starting from the linear state-space system in (1).

*PID controller design:* We first pose PID tuning as output-feedback controller design. Since we feedback integral and derivative terms, we append these terms to the state. Let

$$\begin{aligned} i_{t+1} &= \sum_{t'=1}^{t+1} x_{t'} = i_t + x_t \\ d_{t+1} &= x_t - x_{t-1}. \end{aligned} \quad (4)$$

Then, form the augmented state  $X_t = [x_t; x_{t-1}; i_t]$ , and the augmented state-space equations:

$$\begin{aligned} \begin{bmatrix} x_{t+1} \\ x_t \\ i_{t+1} \end{bmatrix} &= \begin{bmatrix} A & \mathbf{0} & \mathbf{0} \\ I & \mathbf{0} & \mathbf{0} \\ I & \mathbf{0} & I \end{bmatrix} \begin{bmatrix} x_t \\ x_{t-1} \\ i_t \end{bmatrix} + \begin{bmatrix} B \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix} u_t + w_t \\ Y_t &= \begin{bmatrix} C & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & C \\ C & -C & \mathbf{0} \end{bmatrix} \begin{bmatrix} x_t \\ x_{t-1} \\ i_t \end{bmatrix} + e_t. \end{aligned} \quad (5)$$

We write these equations more concisely as:

$$\begin{aligned} X_{t+1} &= \tilde{A}X_t + \tilde{B}u_t + w_t \\ Y_t &= \tilde{C}X_t + e_t, \end{aligned} \quad (6)$$

with  $w_t$  and  $e_t$  defined appropriately. If the initial system is stabilizable and detectable, so is the augmented version. This is because the dynamics of the augmented states are not directly controllable. If the system in equation 1 is exponentially stable, then the augmented states cannot grow exponentially.

All PID controllers (with  $\alpha = 0$ ) can be expressed as  $u_t = -KY_t$  for the system. If the state was measurable, i.e.  $\tilde{C} = I$ , the problem would reduce to LQR, and the optimal PID gains can be obtained by both gradient descent and policy optimization [15]. However, the output-feedback controller optimization problem is still open (6). For  $\alpha \neq 0$ , it is straightforward to introduce a controller that stores the filtered derivative as its state. Such a controller would not be purely state/output feedback. Let  $X_t^c$  be the controller state at time  $t$ . The controller with  $\alpha \neq 0$  can be expressed as:

$$\begin{aligned} X_{t+1}^c &= \alpha X_t^c + K_b Y_t \\ u_t &= -K_x X_t^c - K Y_t \end{aligned} \quad (7)$$

We append  $X_t^c$  to  $X_t$  to augment the state-space further:

$$\begin{aligned} \begin{bmatrix} X_{t+1} \\ X_{t+1}^c \end{bmatrix} &= \begin{bmatrix} \tilde{A} & \mathbf{0} \\ \mathbf{0} & \alpha I \end{bmatrix} \begin{bmatrix} X_t \\ X_t^c \end{bmatrix} + \begin{bmatrix} \tilde{B} & \mathbf{0} \\ \mathbf{0} & K_b \end{bmatrix} \begin{bmatrix} u_t \\ Y_t \end{bmatrix} + \begin{bmatrix} w_t \\ \mathbf{0} \end{bmatrix} \\ Y_t &= \tilde{C} X_t + e_t \\ u_t &= -K_x X_t^c - K Y_t \end{aligned} \quad (8)$$

Henceforth, we assume  $\alpha = 0$  with the extension to  $\alpha \neq 0$  possible with the above controller state.

1) *Actuator saturation as a disturbance:* Let  $w_t^a \triangleq B'(\text{sat}(u_t) - u_t)$  denote the saturation error. We treat the saturation error as a disturbance and write

$$\begin{aligned} X_{t+1} &= \tilde{A} X_t + \tilde{B} \text{sat}(u_t) + w_t \\ &= \tilde{A} X_t + \tilde{B} u_t + w_t^a + w_t. \end{aligned} \quad (9)$$

The error  $w_t^a$  is a non-linear function of the input. It can be modeled as adversarial (as in some online learning settings). To handle adversarial disturbances, [2] introduces disturbance feedback policies of the form:

$$u = -K X_t - \sum_{l=1}^h K_d^{[l]} w_{t-l}. \quad (10)$$

If this approach is provided with a stabilizing controller  $K$ , the resulting online optimization is convex and provides tight regret bounds [2]. To draw the connection between disturbance feedback policies and back-calculation, we observe that if  $h$  is the length of the simulation horizon and  $K_d^{[l]} = K_d$  for all  $l$ , this class of controllers reduces to the back-calculation method. In back-calculation we integrate the disturbances due to actuator saturation and feed it back to the input.

2) *Disturbance feedback policies in episodic learning:* Here, we focus on an episodic setting. We run an episode with the PID parameters and tune the parameters episodically. In this case, additional modeling assumptions are required for disturbance feedback controllers to be meaningful. We postulate that the disturbance has marginally stable dynamics that we want to learn. Hence, we introduce a predictor for  $w_t^a$ :

$$w_t^a = \sum_{i=1}^h M^{[i]} w_{t-i}^a. \quad (11)$$

We want the state to be sufficient for selecting the control action at any given time. Hence we augment the state to

$$Z_t = [\tilde{X}_t; w_t^a; w_{t-1}^a; w_{t-2}^a \dots; w_{t-h}^a].$$

When the disturbance is purely stochastic (no internal dynamics), state/output feedback is optimal. Here we model disturbance dynamics and use the model to obtain a class of disturbance feedback policies. Using the augmented state  $Z_t$ , we can write the dynamics as:

$$\begin{aligned} Z_{t+1} &= \begin{bmatrix} \tilde{A} & I & 0 & 0 \\ 0 & M^{[1]} & M^{[2:h-1]} & M^{[h]} \\ 0 & I & I & 0 \end{bmatrix} Z_t + \begin{bmatrix} \tilde{B} \\ 0 \\ 0 \end{bmatrix} u_t + w_t^r \\ Y_t^z &= \begin{bmatrix} \tilde{C} & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix} Z_t + e_t^r \end{aligned}$$

with  $w_t^r$  the unmodeled disturbance. The dynamics of  $w_t^a$  has to be at least marginally stable in order for a controller to stabilize the system, but it need not be asymptotically stable. The class of output-feedback controllers  $u_t = -K Y_t^z$  can be expressed as:

$$\begin{aligned} u_t &= -K_c Y_t - K_d' w_{t:t-h}^a \\ &= -K_c Y_t - K_d' \begin{bmatrix} M^{[1:h]} & I \end{bmatrix} w_{t-1:t-h}^a \\ &= -K_c Y_t - K_d w_{t-1:t-h}^a. \end{aligned}$$

This gives a general class of controllers that includes disturbance-feedback controllers and the back-calculation method. Further, for the specific scenario of anti-windup compensation, we can recover  $w_t^a$  required in the policy using the assumed actuator model.

3) *Optimization for Parameter Tuning:* In order to tune  $K_c$  and  $K_d$ , we perform gradient descent with the objective function

$$\min_{K_c, K_d} \sum_t y_t^T Q y_t + u_t^T R u_t. \quad (12)$$

Solving this optimization problem using gradient-descent does not require explicit learning of the disturbance dynamics since the controller parameters can be directly optimized. Note that the dynamics are only marginally stabilizable. However, the uncontrollable and marginally stable components of the state do not appear in the cost function. A theoretical understanding of gradient descent for marginally stabilizable systems with output feedback is beyond the scope of this work. We relegate this study to future work and for the present paper perform an empirical evaluation of this non-convex optimization problem.

#### IV. NUMERICAL RESULTS

We illustrate the proposed approach of tuning the controller parameters by differentiating through the model around the feedback loop using four different systems. The setup for the different experiments is summarized in Table I.

In each of the experiments, we use a linear system with actuator saturation. Systems 1,3 are unstable, and systems 2,4 are stable. The systems are progressively higher order and complex. The only non-linearity in the system is actuator saturation. The saturation limits were chosen to induce windup in the absence of any anti-windup strategy. The simulation horizon is chosen to balance the steady-state cost with the transient costs to obtain a good step response. If the simulation

	Plant	Actuator limits	Limits on the step reference	Initial feedback gains
System 1	$P(s) = \frac{2e^{-0.02s}}{s-0.995}$	$\pm 3.3$	$\pm 4$	$k_p = 4, k_i = 10, b = 0.5$
System 2	$P(s) = \frac{1}{(20s^2+10s+1)}$	$\pm 7.0$	$\pm 4$	$k_p = 10, k_i = 1.5, k_d = 8, b = 0.4$
System 3	$P(s) = \frac{1}{(s+0.1)(s-0.1)}$	$\pm 3.0$	$\pm 2.9$	$k_p = 20, k_i = 2, k_d = 5, b = 1$
System 4	$P(s) = \frac{1}{(s+0.1)(s+0.2)(s+0.4)(s+0.6)}$	$\pm 4.0$	$\pm 3$	$k_p = 20, k_i = 8, k_d = 10, b = 0.2$

TABLE I  
SETUP OF THE EXPERIMENTS

horizon is too short, then the costs in the loss function due to the transients will dominate. The controllers have to ensure precise convergence to the step after the transients. We convert the continuous-time linear system into its discrete-time counterpart using zero-order hold (ZOH) and simulate the system in discrete time. We tune the feedback gains episodically by performing gradient descent on the squared-error cost function  $\sum_{t=1}^T (y_t - r_t)^2$ . We generate 30 reference signals and segregate them into 20 signals used for training and 10 used for testing.

For each system, we compare four distinct controllers: (1) The initial PI/PID controller is tuned using classical control techniques to work well without actuator saturation. Integrator windup impacts performance when saturation is present; (2) We initialize a back-calculation constant manually to decrease windup; (3) The PID and back-calculation parameters are then optimized using gradient descent with the Adam optimizer [20] to obtain a third controller; (4) This controller has dynamically changing feedback gains modeled as a small neural network with tracking error and actuation error as inputs.

Simulation 1 performs step-reference tracking on system 1. After optimization, the parameters converged to  $k_p^* = 16.58$ ,  $k_i^* = 11.07$  and  $b^* = 0.87$ . The cost on the training and testing reference signals of the four different controllers are summarized in Table II. Optimization of the feedback gains using gradient descent is effective in improving performance. Using a neural-network to change the feedback gains dynamically does not further improve performance. This suggests that simple controllers are sufficient to control linear systems with saturation. From Figure 1, we see that both the dynamic and static optimized PI controllers are effective in preventing windup and tracking the reference accurately. In Figure 2, we plot the variation of the dynamic PI controller's feedback gains with time. It is interesting that the feedback gains switch during the transition and return to the initial values. This indicates that switching between multiple controllers depending on the input is an effective strategy to improve performance. However, for the simple systems used here, the static controller is sufficient for good performance.

Simulations 2 and 3 repeat the above procedure on the second-order systems 2 and 3. In simulation 2, the parameters converged to  $k_p^* = 11.31$ ,  $k_i^* = 1.71$ ,  $k_d^* = 4.26$  and  $b^* = 0.23$ . In simulation 3, the parameters converged to  $k_p^* = 7.72$ ,  $k_i^* = 1.47$ ,  $k_d^* = 3.01$  and  $b^* = 0.66$ . The performance of the controllers in terms of squared error cost is summarized in table III. Figures 3 and 4, indicate that

Method	Training cost	Test cost
Initial PI	$304.4 \pm 432.3$	$349.0 \pm 503.3$
Initial PI with backcalculation	$178.2 \pm 153.0$	$189.0 \pm 164.8$
PI+backcalculation optimized	$110.2 \pm 79.8$	$114.7 \pm 82.3$
Dynamic PI+backcalculation optimized	$109.5 \pm 79.9$	$114.0 \pm 82.2$

TABLE II  
SQUARED ERROR COST OF THE FOUR CONTROLLERS ON SYSTEM 1. SEE SECTION IV FOR INTERPRETATION.

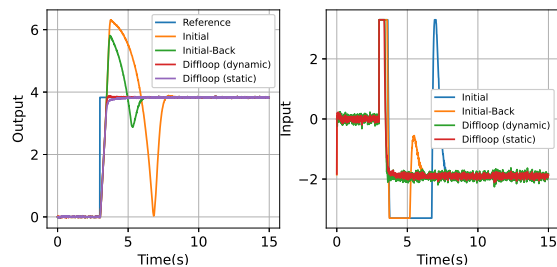


Fig. 1. Performance of the four controllers on a difficult test reference for system 1. Left: Output of the SISO system with the four controllers. Right: Input to the system with the four controllers. See section IV for interpretation.

optimization of the PID controllers is effective.

Finally, in simulation 4, we track a rapidly switching reference using a PID controller with back-calculation. The saturation limits constrain the controller from switching the outputs rapidly to achieve accurate tracking. However, the optimization minimizes the squared error cost. Interestingly, for rapidly switching reference signals, the feedback gains converges to  $k_p^* = 11.78$ ,  $k_i^* = -0.47$ ,  $k_d^* = 3.46$  and  $b^* = 0.22$ . The integrator gain consistently converges to negative values. For a rapidly switching reference, the transient costs outweigh the cost due to lack of precise convergence and noise. Hence,

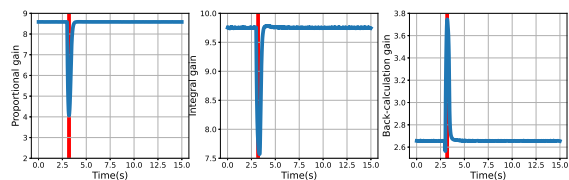
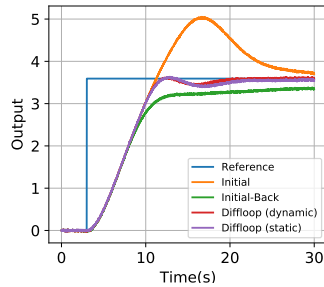
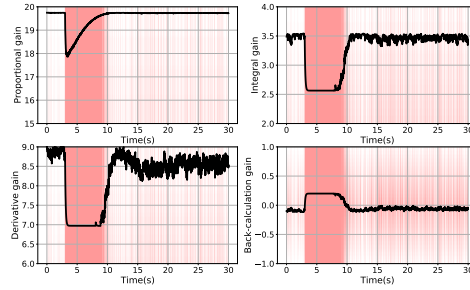


Fig. 2. Variation of the tuning constants with time for the dynamic PI controller for system 1. The red shaded region indicates periods of saturation. The controller gains switch during saturation.

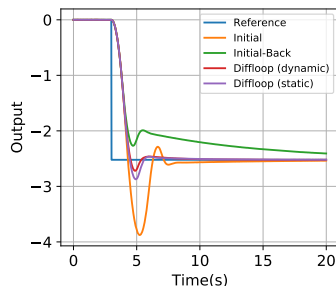


(a) System 2

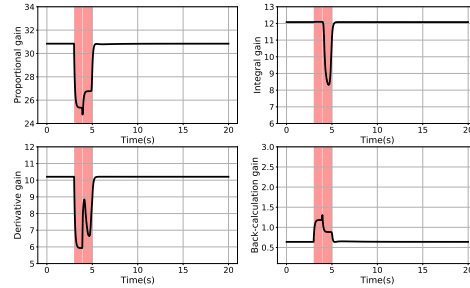


(b) System 2 feedback gains

Fig. 3. (a) Output of the four controllers on a step input for system 2. We can see that optimizing the PID gains both in the static and dynamic case is effective in providing a good step response. (b) Variation of the feedback gains with time for the Dynamic PID controller. The gains switch during the transition and return to their initial values. The red shaded regions denote periods of saturation

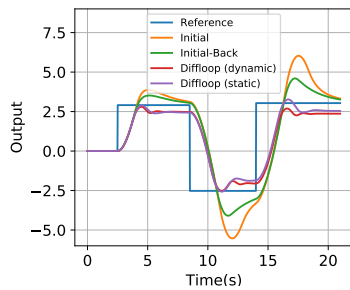


(a) System 3

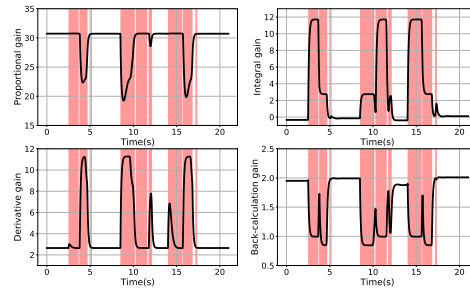


(b) System 3 feedback gains

Fig. 4. (a) Output of the four controllers on a step input for system 3. We can see that optimizing the PID gains both in the static and dynamic case is effective in providing a good step response. (b) Variation of the feedback gains with time for the Dynamic PID controller. The gains switch during the transition and return to their initial values. The red shaded regions denote periods of saturation



(a) System 4 output



(b) System 4 feedback gains

Fig. 5. (a) Output of the four controllers on a rapidly switching reference for system 4. We can see that optimizing the PID gains both in the static and dynamic case is effective to some extent to cope up with the hard saturation limit. (b) Variation of the feedback gains with time for the Dynamic PID controller. The gains switch during the transitions and return to their initial values. The red shaded regions denote periods of saturation

Method	System 2 Test cost	System 3 Test cost
Initial PID	$2115.2 \pm 1822.9$	$264.8 \pm 220.0$
Initial PID with backcalculation	$1788.1 \pm 1506.7$	$247.6 \pm 197.0$
PID+backcalculation optimized	$1688.9 \pm 1405.9$	$198.2 \pm 148.7$
Dynamic PID+backcalculation optimized	$1686.7 \pm 1410.7$	$198.1 \pm 150.0$

TABLE III

PERFORMANCE OF THE FOUR CONTROLLERS IN TERMS OF SQUARED ERROR FOR SYSTEM 2 AND SYSTEM 3. SEE SECTION IV FOR INTERPRETATION.

it is better to choose feedback gains that reduce overshoot and windup. This observation gives rise to an interesting phenomenon. In figure 5, we plot the system's outputs with the four different controllers. Even though actuator saturation is too limiting to switch rapidly, the optimized controllers achieve lower squared error costs and perform better.

## V. CONCLUSIONS AND FUTURE WORK

We outline a PID tuning approach for linear systems with input saturation. This approach differentiates through the model and around the feedback loop to tune the controller parameters. The numerical experiments demonstrate the efficacy of this approach. We also propose a theoretical framework to analyze the convergence properties for this optimization. However, a convergence proof is beyond the scope of this work. We noted that the framework shows the equivalence of the backcalculation method and disturbance feedback policies.

Future work can extend this technique for generating robust controllers for MIMO systems using robust optimization. Further, the automatic differentiation technique could also be used for tuning PID controllers in robotic systems. Output feedback controller optimization also warrants further theoretical study.

## REFERENCES

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [2] Naman Agarwal, Brian Bullins, Elad Hazan, Sham M Kakade, and Karan Singh. Online control with adversarial disturbances. *arXiv preprint arXiv:1902.08721*, 2019.
- [3] Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks. *arXiv preprint arXiv:1703.00443*, 2017.
- [4] Karl Johan Åström and Tore Hägglund. *PID controllers: theory, design, and tuning*, volume 2. Instrument society of America Research Triangle Park, NC, 1995.
- [5] Karl Johan Åström and Richard M Murray. *Feedback systems: an introduction for scientists and engineers*. Princeton university press, 2010.
- [6] C Bohn and DP Atherton. An analysis package comparing pid anti-windup strategies. *IEEE Control Systems Magazine*, 15(2):34–40, 1995.
- [7] Hamid Boubertakh, Mohamed Tadjine, Pierre-Yves Glorennec, and Salim Labiod. Tuning fuzzy pd and pi controllers using reinforcement learning. *ISA transactions*, 49(4):543–551, 2010.
- [8] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. *JAX: composable transformations of Python+NumPy programs*, 2018. <http://github.com/google/jax>.
- [9] Michael B Chang, Tomer Ullman, Antonio Torralba, and Joshua B Tenenbaum. A compositional object-based approach to learning physical dynamics. *arXiv preprint arXiv:1612.00341*, 2016.
- [10] Shih-Feng Chen. Particle swarm optimization for pid controllers with robust testing. In *2007 International Conference on Machine Learning and Cybernetics*, volume 2, pages 956–961. IEEE, 2007.
- [11] P Cominos and N Munro. Pid controllers: recent tuning methods and design to specification. *IEE Proceedings-Control Theory and Applications*, 149(1):46–53, 2002.
- [12] Filipe de Avila Belbute-Peres, Kevin Smith, Kelsey Allen, Josh Tenenbaum, and J Zico Kolter. End-to-end differentiable physics for learning and control. In *Advances in Neural Information Processing Systems*, pages 7178–7189, 2018.
- [13] Jonas Degraeve, Michiel Hermans, Joni Dambre, et al. A differentiable physics engine for deep learning in robotics. *Frontiers in neurorobotics*, 13:6, 2019.
- [14] Andreas Doerr, Duy Nguyen-Tuong, Alonso Marco, Stefan Schaal, and Sebastian Trimpe. Model-based policy search for automatic tuning of multivariate pid controllers. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5295–5301. IEEE, 2017.
- [15] Maryam Fazel, Rong Ge, Sham M Kakade, and Mehran Mesbahi. Global convergence of policy gradient methods for linearized control problems. 2018.
- [16] Sergio Galeani, Sophie Tarbouriech, Matthew Turner, and Luca Zaccarian. A tutorial on modern anti-windup design. *European Journal of Control*, 15(3-4):418–440, 2009.
- [17] Elad Hazan, Sham Kakade, and Karan Singh. The nonstochastic control problem. In *Algorithmic Learning Theory*, pages 408–421. PMLR, 2020.
- [18] JM Herrero, X Blasco, M Martinez, and JV Salcedo. Optimal pid tuning with genetic algorithms for non-linear process models. In *15th Triennial World Congress, Barcelona, Spain*, 2002.
- [19] Payam Kheirkhahan. Robust anti-windup control design for pid controllers. In *2017 17th International Conference on Control, Automation and Systems (ICCAS)*, pages 1622–1627. IEEE, 2017.
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [21] Nathan P Lawrence, Gregory E Stewart, Philip D Loewen, Michael G Forbes, Johan U Backstrom, and R Bhushan Gopaluni. Optimal pid and antiwindup control design as a reinforcement learning problem. *arXiv preprint arXiv:2005.04539*, 2020.
- [22] Dougal Maclaurin, David Duvenaud, Matt Johnson, and Jamie Townsend. *Autograd*. <https://github.com/HIPS/autograd>.
- [23] Yasue Mitsukura, T Yamamoto, and M Kaneda. A genetic tuning algorithm of pid parameters. In *1997 IEEE International Conference on Systems, Man, and Cybernetics. Computational Cybernetics and Simulation*, volume 1, pages 923–928. IEEE, 1997.
- [24] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [25] Benjamin Recht. A tour of reinforcement learning: The view from continuous control. *Annual Review of Control, Robotics, and Autonomous Systems*, 2:253–279, 2019.
- [26] Qian Shi, Hak-Keung Lam, Bo Xiao, and Shun-Hung Tsai. Adaptive pid controller based onq-learning algorithm. *CAAI Transactions on Intelligence Technology*, 3(4):235–244, 2018.
- [27] Max Simchowitz, Karan Singh, and Elad Hazan. Improper learning for non-stochastic control. *arXiv preprint arXiv:2001.09254*, 2020.
- [28] Sigurd Skogestad. Probably the best simple pid tuning rules in the world. In *AICHE Annual Meeting, Reno, Nevada*, volume 77, 2001.
- [29] Erik Torstensson. Comparison of schemes for windup protection. *ISSN 0280-5316*, 2013.
- [30] Kaiqing Zhang, Bin Hu, and Tamer Basar. Policy optimization for  $\mathcal{H}_2$  linear control with  $\mathcal{H}_\infty$  robustness guarantee: Implicit regularization and global convergence. In *Learning for Dynamics and Control*, pages 179–190, 2020.

- [31] G Zhou and J Douglas Birdwell. Pid autotuner design using machine learning. In *IEEE Symposium on Computer-Aided Control System Design*, pages 173–179. IEEE, 1992.
- [32] John G Ziegler, Nathaniel B Nichols, et al. Optimum settings for automatic controllers. *trans. ASME*, 64(11), 1942.