

A Comparison of Tracking Step Inputs with a Piezo Stage Using Model Predictive Control and Saturated Linear Quadratic Gaussian Control

Lucy Y. Pao^{a,*}

^a*Department of Electrical, Computer & Energy Engineering, University of Colorado, Boulder, CO 80309 USA*

Abstract

Compressed Sensing for Atomic Force Microscopy is a newer imaging mode that requires the piezo stage be driven rapidly between measurement locations. In contrast to raster scanning applications, this translates to a setpoint tracking problem. This paper considers the setpoint tracking performance of a piezo nano-positioning stage subject to rate-of-change limitations on the control signal, which is derived from the current limit of the power amplifier. To compensate the vibrational dynamics of the stage, a model predictive control scheme (MPC) and a linear quadratic Gaussian (LQG) controller which saturates the control increment are considered. In both cases, hysteresis and drift are compensated via dynamic inversion. To design the weighting matrices required by the MPC and linear feedback designs, an extension to classic reciprocal root locus ideas is proposed. The robustness of both schemes using classical methods like gain margin, phase margin, and gain of the sensitivity function at low frequencies is analyzed. The overall settle times achieved by both controllers (in both simulation and experiment) across a range of control weights where the reference input is a sequence of step inputs of varying amplitudes are compared. The results show that the best simulation settle time is achieved by MPC using the smallest control weight. However under experimental conditions, the best settle time is achieved by a much larger control weight and the performance of MPC becomes comparable with that of saturated linear feedback. This result is explained by showing that robustness increases with larger control weights.

Keywords: tracking step inputs, piezo stage, model predictive control, linear quadratic Gaussian control, hysteresis

1. Introduction

The Atomic Force Microscope (AFM) is a nano-scale imaging instrument which acquires an image of the surface topography of a specimen by mechanically interrogating it with an atomically-sharp probe [1, 2]. Typically, the probe is scanned across a specimen in a raster pattern, sequentially acquiring pixels in an image. Although this process gives the AFM excellent spatial resolution, the serial acquisition of pixels limits the speed of any given instrument, yielding image acquisition rates on the order of minutes per frame for many commercially available instruments.

For static specimens, while slow imaging yields accurate images, the long imaging times are inconvenient. For dynamic specimens, however, slow imaging precludes the ability to capture the dynamics of such specimens; while faster imaging has the potential to allow for the study of the specimen dynamics. Many methods have been proposed to increase AFM frame rates, including better mechanical design [3, 4], using advanced control methods [5, 6, 7, 8], and alternative scanning methods [9, 10, 11, 12, 13, 14].

One newer alternative to raster scanning, and which is of interest here, is the application of Compressive Sensing to

AFM [15, 16, 17]. The central idea of compressive-sensing-based imaging is to leverage the redundancy present in most interesting images such that the number of pixels to be acquired is reduced. For good guarantees on reconstruction quality, measurements in compressive-sensing-based imaging need to be randomly distributed across the specimen. Each measurement might acquire a single pixel [16] or short string of adjacent pixels in a micro-scan [18, 19]. Once a measurement is completed, the AFM probe is retracted from the specimen surface, moved in the XY plane to the next measurement location, and finally re-engaged with the specimen surface before the next measurement is acquired. Details of implementations of this approach can be found in [18, 20].

In this paper, we are concerned with the point-to-point movement in the XY plane between measurement locations. Because the probe is not in contact with the specimen during this operation, it is desirable to minimize the time to move between measurement locations. Thus, in contrast to standard raster scanning where the control goal is to minimize overall tracking error to a triangular reference, the goal here is to minimize the settle time to a step input. Point-to-point movements by AFM are also of interest in other areas like viscoelastic property mapping [21].

One of the primary constraints in setpoint tracking with our piezo stage is the current limit of the power amplifier,

*Corresponding author

Email address: pao@colorado.edu (Lucy Y. Pao)

which roughly translates to a slew-rate limitation on the power amplifier output voltage. In principle, minimizing the settle time of such point-to-point motions is a classic time-optimal control problem. However, for stages with dynamics more complex than a second-order system, including the stage in our own lab, closed-form solutions to the minimum-time problem are intractable.

An enticing alternative to explicitly handle the slew-rate constraint is MPC with a purely quadratic cost. Given a discrete-time state-space system $\{A, B, C, 0\}$ with state x_k , and control input u_k , such an MPC scheme solves, at each time step, the optimal control problem

$$\min_v z_N^T P z_N + \sum_{i=0}^{N-1} z_i^T Q z_i + 2z_i^T S v_i + v_i^T R v_i \quad (1a)$$

$$\text{s.t. } z_{i+1} = A z_i + B v_i \quad (1b)$$

$$z_0 = x_k, \quad (1c)$$

$$v_i \in \mathbb{U} \quad (1d)$$

where N is the control horizon, \mathbb{U} is a polyhedron, P solves the Discrete Algebraic Riccati Equation (DARE), and Q , R , and S are the state, control, and cross weights. Q and R are symmetric matrices, and the triple (Q, R, S) satisfies

$$R > 0 \quad (2)$$

$$Q - S R^{-1} S^T \geq 0. \quad (3)$$

In this paper, \mathbb{U} is restricted to model actuator constraints (e.g., saturation). The solution to the quadratic program (QP) (1) results in a sequence of optimal controls $v_0 \dots v_{N-1}$. One sets $u_k = v_0$ and repeats the process at the next time step. If one eliminates the constraint (1d), then the control action reduces to linear state feedback. That is, $v_0 = -K x_k$ where

$$K = (B^T P B + R)^{-1} (B^T P A + S^T), \quad (4)$$

is the solution to the infinite horizon LQR problem associated with Q , R , and S .

Historically, one of the challenges of applying MPC to systems with fast dynamics is the computational demand needed to solve a QP within a small sample period. However, advances in both computing hardware and algorithms have mitigated this issue. For example, [22] shows that when \mathbb{U} is a simple box (i.e., a saturating constraint), sample rates of up to 1 MHz can be achieved with high-end FPGAs using the Fast Gradient Method (FGM).

In recent work, we applied the FGM formulation of [22] to our piezo stage and showed that, given a particular set of weighting matrices, we could increase the stabilizable range of setpoints compared to simply saturating an equivalent linear feedback [23]. Others have also applied MPC to similar systems [24, 25, 11, 26]. In all these cases, little insight is given into how the cost function was tuned, an issue that is considered here in further depth.

Due to the increased cost and complexity of implementing MPC, it is crucial to characterize how MPC compares

to linear feedback. In some cases, no comparison to linear control is given [11, 26]. In many cases where MPC is compared to a linear feedback law [24, 25], including in our own prior work [23], de-rating the linear feedback to limit constraint violation is never considered. Thus, an important question we seek to answer in this paper is “how much performance is sacrificed by using a de-rated linear feedback compared to MPC?”. In Section 6, we show that, in contrast to simulation results, experiments with de-rated (i.e., large) control weights have *better* performance for both MPC and linear feedback. For control weights where the best experimental performance is attained, linear feedback and MPC yield similar performance.

This somewhat surprising result is explained in Section 6.1 by showing that robustness of the control laws increases as the control weight increases. While many authors have considered robustness in MPC, many of these results assume direct measurements of the state vector [27, 28, 29], which is often impractical. However, as we indicated in (4), it is a well known, though perhaps under exploited, result that when the control trajectory generated by (1) is within the interior of \mathbb{U} , the control action is equivalent to an LQR-based linear feedback law. In the setpoint tracking application considered here, where the constraint limits the rate of change on the control, this will always be the case as the system nears a given setpoint. Thus, within some region around any setpoint, classical ideas like gain and phase margin or the sensitivity function gain are directly applicable. We show that (up to a point) de-rating the design *improves* those metrics. Thus, the nominal performance gains achieved when using a more aggressive MPC (due to constraint handling) are offset by the concomitant decrease in robustness. As we are not aware of publications that have applied these classical metrics in the context of MPC, we present our analysis and results in this paper to provide a useful way to understand and compare the performance of MPC relative to a more classical control approach.

Another related limitation of [23] is that we did not consider the effects of drift and hysteresis and only considered tracking a single setpoint with the stage starting at rest. Other studies also ignore hysteresis and only consider a limited size of inputs [24, 25]. Yet, when tracking a *sequence* of random setpoints across the range of the stage, the effects of hysteresis become much more prominent. Thus, in this paper, we employ inverse drift and hysteresis compensation and test the control laws with a random sequence of steps. Since these inversions are not perfect, good robustness is needed.

The main contributions of this paper are thus to

- Give explicit details on how the cost function is tuned (Sections 5 and 6).
- Compare the experimental performance of MPC to saturated linear feedback (SLF) across a range of control weights varying from aggressive to highly de-rated (Section 6).

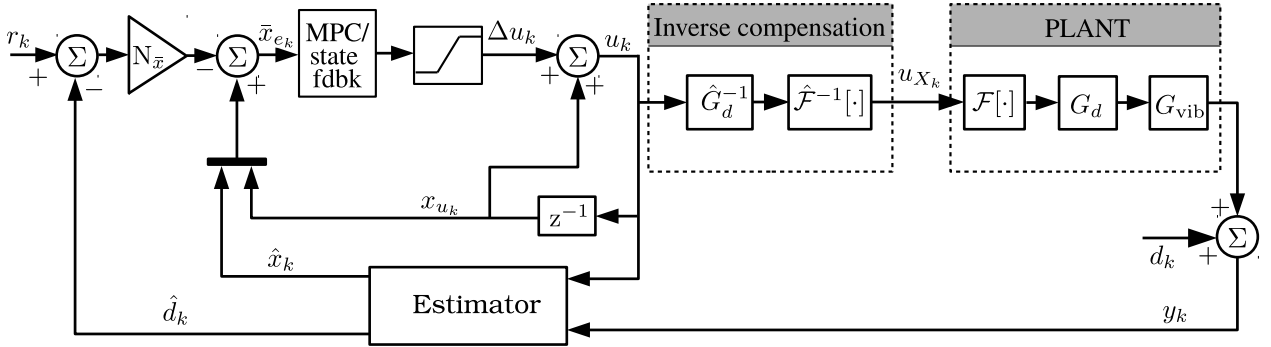


Figure 1: The overall plant model consists of a hysteresis model $\mathcal{F}[\cdot]$, a drift model G_d , and a vibrational model G_{vib} . The effects of drift and hysteresis are compensated for via dynamic inversion.

- Show that for the chosen state weighting schemes, closed-loop robustness plays a more prominent role in ultimate experimental performance than explicitly handling constraints (Section 6.1).

The overall control structure considered in this paper is illustrated in the block diagram in Fig. 1. The plant is considered to be a cascaded model of hysteresis (\mathcal{F}), drift (G_d), and vibrational dynamics (G_{vib}). Modeling these three systems is the subject of Sections 3.3, 3.2, and 3.1, respectively. Section 4 develops the control structure and associated closed-loop equations. Section 5 explores two schemes to design the weighting matrices. The designs are evaluated in simulation and experiment in Section 6, and conclusions are given in Section 7.

2. Experimental Testbed

The AFM in our lab consists of an Agilent 5400 that has been retrofitted with an nPoint NPXY100A piezo stage, which provides lateral movement of the specimen. The NPXY100A, which is the focus of this paper, is driven by an nPoint C300. The C300 amplifies the low voltage (± 10 volts) control inputs to a high voltage signal which drives the piezo actuators and provides signal conditioning for the capacitive position sensors in the stage. Although the C300 can implement a basic PID controller, in this work the C300 is always operated in open-loop mode. Unfortunately, even in open-loop mode, signals in the C300 still run through an internal DSP, which introduces around $360 \mu\text{s}$ of delay. Nominally, the NPXY100A has a range of $100 \mu\text{m}$, though in practice the usable range is about $67.5 \mu\text{m}$ when operated in open-loop mode.

All control logic is programmed into a Xilinx Spartan-6 LX150 FPGA in a cRIO 9082 from National Instruments. A sampling frequency of 25 kHz is used, which is based on the system dynamics (see Fig. 2 in Section 3.1). With a 25 kHz sampling frequency, the $360 \mu\text{s}$ of delay translates to 9 samples.

In characterizing the limitations of this system, it is helpful to have a direct measurement of the power amplifier current, I_X , of the C300. This measurement is

obtained by re-routing the C300 drive signal through a low-side current sensing resistor ($R_{\text{sense}} = 0.1\Omega$).

3. System Modeling

Since there is little coupling between the X and Y directions in our AFM [20] and imaging results have been separately presented in [18, 20], the point-to-point motion discussion and analysis in this paper is focused on the X direction. The overall plant for the X-direction is modeled as three cascaded systems: \mathcal{F} which models the hysteresis of the piezo, a drift model G_d , and a vibrational model G_{vib} . This cascaded structure, suggested in [7, 30], is shown in Fig. 1. In general, the effects of hysteresis are most noticeable when moving across wide ranges. Thus, by using relatively small input signals, the drift and vibrational dynamics can be identified separately from the hysteresis [30]. Here, we model both G_{vib} and G_d as linear, time-invariant discrete-time systems. The dynamics of drift are predominantly low frequency while vibrational aspects on the other hand are fast by comparison, which allows the two systems to be easily separated in the identification. Modeling these three components is the subject of the next three subsections.

3.1. Modeling G_{vib}

To obtain an experimental frequency response of G_{vib} , we use a stepped-sines method (single frequency at a time). The amplitude of the driving sinusoid is chosen to be small enough that the effects of hysteresis are minimized. After the system reaches steady-state, the input and output signals are demodulated into their first (complex) Fourier coefficients, the ratio of which yields the frequency response at that frequency. Fig. 2 shows the resulting experimental frequency response function (FRF) as the solid red curve.

Obtaining a parametric model of G_{vib} for control design involves two steps. A preliminary model using an Eigenspace Realization Algorithm (ERA) [31] is obtained. In general, the ERA does not produce a model with poles at $z = 0$, which is what is needed in order to model the delay. Thus, the delay in the frequency response is divided out of the FRF before passing it to the ERA algorithm.

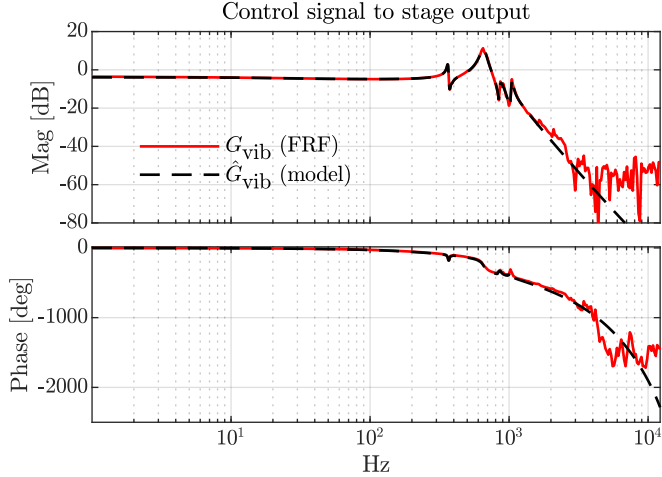


Figure 2: The solid red curve is the frequency response from control input to stage position output in the X direction. The dashed-black curve is the vibrational model, \hat{G}_{vib} .

The second step uses the model generated by the ERA as the initial guess to a non-linear least squares problem which minimizes the logarithm of the ratio of the experimental frequency response to that of the model [32]. Though [32] develops the idea for continuous-time models, their strategy is easily adapted to fit a discrete-time model. In this scenario, the optimization is given by

$$\min_{\theta} \sum_{i=1}^M \left| \log(G_{\text{vib}}(e^{j\omega_i T_s})) - \log(\hat{G}_{\text{vib}}(e^{j\omega_i T_s} | \theta)) \right|^2 \quad (5)$$

where ω_i is each frequency in the experimental frequency response and $\hat{G}_{\text{vib}}(e^{j\omega_i T_s} | \theta)$ is the model parameterized by the vector θ . The model $\hat{G}_{\text{vib}}(z | \theta)$ is composed of first- and second-order factors

$$\hat{G}_{\text{vib}}(z | \theta) = K \frac{\prod_{i=0}^{n_{rz}-1} (z - b_i^r)}{\prod_{\ell=0}^{n_{rp}-1} (z - a_{\ell}^r)} \times \frac{\prod_{j=0}^{n_{cz}-1} (z^2 + b_{2j}^c z + b_{2j+1}^c)}{\prod_{m=0}^{n_{cp}-1} (z^2 + a_{2m}^c z + a_{2m+1}^c)} z^{-\rho} \quad (6)$$

where n_{rz} and n_{cz} (resp., n_{rp} and n_{cp}) are the number of real and complex zeros (resp., poles) in the model generated by the ERA. The parameter vector θ is given by

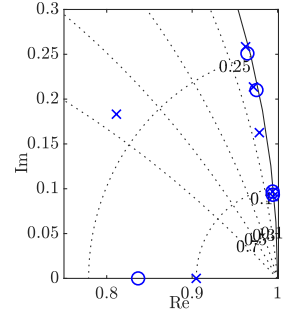
$$\theta = [b_0^r \dots b_{n_{rz}-1}^r, b_0^c \dots b_{2n_{cz}-1}^c, a_0^r \dots a_{n_{rp}-1}^r, a_0^c \dots a_{2n_{cp}-1}^c, K, \rho].$$

Due to the logarithms in (5) and the multiplicative structure (6), the Jacobian of $\log(\hat{G}_{\text{vib}}(z | \theta))$ is surprisingly easy to calculate. Details can be found in [32], though some modifications are required for the discrete-time case.

The model structure (6) includes a fractional delay $z^{-\rho}$. This allows including the delay as a term in the decision variable as an alternative to optimizing over integers. This is beneficial since it is not guaranteed that the latency from

Table 1: Pole and zero locations of \hat{G}_{vib} . The graphic to the right shows their location (excluding the nine poles at $z = 0$ and the non-minimum phase zero) in the Z -plane.

pole	zero
0.8113 ± 0.1832	0.9647 ± 0.2509
0.9625 ± 0.2584	0.9753 ± 0.2101
0.9718 ± 0.2136	0.9945 ± 0.0933
0.9786 ± 0.1626	0.9942 ± 0.0969
0.9946 ± 0.0912	—
0.9942 ± 0.0968	—
0.9047	-1.55
(9) 0	0.8368



input to output is an exact integer multiple of the sample period, and a fractional delay allows the optimization to more accurately match the phase. In the final model, ρ is rounded to the nearest integer, (8 in this case), which gives a sum squared error in the phase match of 0.196. If the optimization was instead performed with ρ fixed at 8, the sum squared phase error is 3.98. In this work, modes above 1100 Hz are not modeled. Thus, the optimization (5) is only done over frequency up to 1100 Hz. The final pole and zero locations for \hat{G}_{vib} are listed in Table 1.

3.2. Drift Modeling

Drift is modeled as the transfer function

$$\hat{G}_d(z | \theta) = \theta_5 \frac{(z - \theta_1)(z - \theta_2)}{(z - \theta_3)(z - \theta_4)},$$

which is equivalent to the model structure used in [30]. An LTI drift model can achieve a better fit than both logarithmic and fractional derivative models [33]. Due to the comparatively slow dynamics of drift, it is more attractive to identify the drift model in the time domain rather than the frequency domain. The stage is given a step input with relatively small amplitude (to minimize the effects of hysteresis). The stage response is shown as the solid-blue curve in Fig. 3, while the simulated response of the vibrational model is shown as the dotted-black curve. The piezo drift is evident in the slow increase of stage position after the vibrational dynamics have decayed.

Let \mathcal{Y}_{exp} be the step response data collected from the stage and \mathcal{Y}_{vib} the response of the model \hat{G}_{vib} to the same input. The goal is to solve the non-linear least squares problem

$$\min_{\theta} \left\| \hat{g}_d(k | \theta) * \mathcal{Y}_{\text{vib}} - \mathcal{Y}_{\text{exp}} \right\|_2 \quad (7)$$

where $\hat{g}_d(k | \theta)$ is the impulse response corresponding to $\hat{G}_d(z | \theta)$, '*' represents the convolution operator, and $\theta = [\theta_1 \dots \theta_5]$ is the vector of parameters. To the extent that G_{vib} accurately models the vibrational dynamics, the inclusion of \mathcal{Y}_{vib} in (7) effectively nullifies the vibrational aspects in the optimization. This is possible because, since the system is SISO, G_{vib} and G_d commute. The non-linear optimization problem (7) is solved with MATLAB's

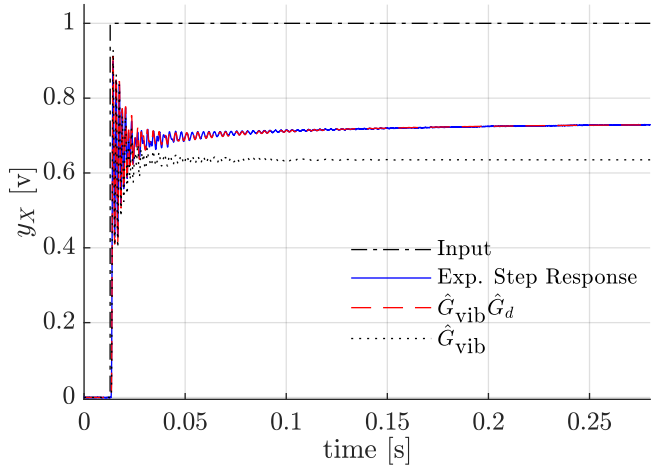


Figure 3: A 1.0 volt amplitude step input (dash-dotted black) is given to the stage, yielding the solid blue output trajectory. The response of the combined G_{vib} and G_d models is shown as the dashed red, while that of the vibrational model alone is the dotted black curve.

`lsqnonlin` and results in the red curve in Fig. 3, which shows the simulated step response of $\hat{G}_d\hat{G}_{\text{vib}}$. Validating for other small amplitude step inputs indicates that the obtained \hat{G}_d is similarly accurate as shown in Fig. 3.

3.3. Hysteresis Modeling

In typical raster scanning applications, hysteresis manifests as a bowing of the trajectory as the stage tries to track the linear ramps in a triangle wave (see, e.g., Fig. 3 of [7]). To motivate the need for hysteresis compensation in a step tracking application, consider Fig. 4, which shows an input signal of various filtered steps applied open-loop to the stage. The solid black curve is the stage response, while the dotted-black curve is the input (scaled by the nominal DC-gain of G_dG_{vib}), which shows good agreement for the first step, but much worse agreement with the later steps, particularly those with large amplitudes. Effectively, the gain of the system depends on the control history.

There are many models for hysteresis [30, 34, 35]. Here, the simple (and computationally fast) Modified Prandtl-Ishlinski Hysteresis model from [36] is used. This hysteresis model is composed of a linear combination of saturation operators cascaded with a linear combination of classic hysteretic play¹ operators. The overall input-output relationship of the modified hysteresis operator $\mathcal{F}[\cdot]$ is

$$\mathcal{F}(u_X) = w_s^T \mathbf{S} [w_H^T \mathbf{H}[u_X, z]]$$

where \mathbf{S} and \mathbf{H} are vectors of elementary saturation and play operators, respectively, and where w_s and w_H are vectors of weights. The i th elementary play operator with associated threshold d_H^i , output ξ_k^i , and input ν_k is defined by the recursive relationship

$$\xi_k^i = H^i(\nu_k | d_H^i) = \max\{\xi_{k-1}^i - d_H^i, \min\{\xi_{k-1}^i + d_H^i, \nu_k\}\}.$$

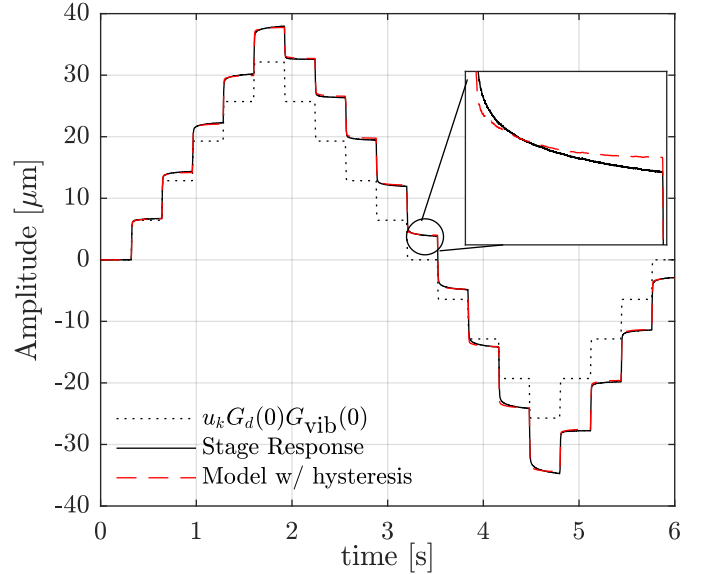


Figure 4: The stage is driven by a sequence of filtered step inputs shown in the dotted black curve. The resulting stage response is the solid black curve, showing good agreement at the first step, but much worse agreement for later steady-state values. The same steady-state input does not produce the same steady-state output, which can be seen, e.g., by comparing the step centered at 4.3 s to the step centered at 4.95 s. The dashed red curve is the response of the overall combined model of G_{vib} , G_d , and the complex hysteresis model \mathcal{F} .

In contrast, the saturation operator has no memory. The input-output relationship of the i th elementary saturation operator with associated threshold d_S^i is

$$\mu_k = S^i(p_k | d_S^i) = \begin{cases} \max\{p_k - d_S^i, 0\} & d_S^i > 0 \\ p_k & d_S^i = 0 \\ \min\{p_k + d_S^i, 0\} & d_S^i < 0. \end{cases}$$

for an input p_k and output μ_k .

If the thresholds d_H^i and d_S^i are pre-defined, it is possible to solve for the weights w_S and w_H from a quadratic program [36]. Using the input and output data from Fig. 4, the weights w_H and w_S are fit to a model with 7 saturation and 7 hysteresis operators. Details on model order selection can be found in [36]. The resulting fit is shown in Fig. 4 as the dashed red curve (which also includes the drift and vibrational models). Again, validating for other sequences of filtered step inputs shows that the obtained drift model yields similar accuracies as shown in Fig. 4.

3.4. Power Amplifier Characterization and Limitations

The high voltage output of the C300 is current limited to 100 mA. The solid black curve in Fig. 5 shows the transfer function, $G_{I_X, u_X}(z)$, from the low voltage input u_X of the C300 to the current I_X flowing through the stage. The current is measured via the sensing resistor mentioned in

¹The term “play” is derived from the operator’s use in modeling mechanical slop.

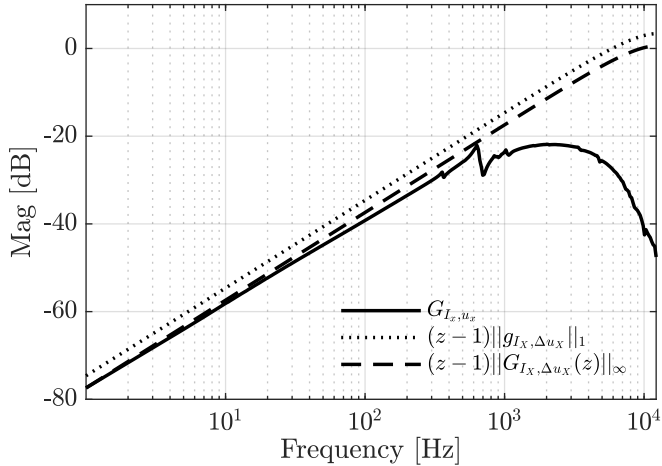


Figure 5: The solid curve is the transfer function, $G_{I_X, u_X}(z)$, from low voltage control to power amplifier output current, which is upper-bounded by the dotted and dashed curves, representing a pure discrete derivative multiplied by the bounds (11) and (12), respectively.

Section 2. Because the piezo actuators are highly capacitive, at frequencies below about 600 Hz, $G_{I_X, u_X}(z)$ looks like a pure derivative. To derive a box constraint on the control increment $\Delta u_X(z)$, we first factor $G_{I_X, u_X}(z)$ as

$$I_X(z) = G_{I_X, u_X}(z)u_X(z) \quad (8)$$

$$\begin{aligned} &= (z-1)G_o(z)u_X(z) \\ &= G_o(z)\Delta u_X(z). \end{aligned} \quad (9)$$

Ideally, the current limit is enforced via a state-like constraint using (8) and a parametric model of $G_{I_X, u_X}(z)$. While it is likely possible to solve such a problem with a high-end FPGA using, e.g., the Alternating Direction Method of Multipliers [22], it is not possible to solve that problem on our hardware. Instead, we approximate $G_o(z)$ as a constant, which leads to a box constraint on Δu_X . Thus, we need a bound $(\Delta u_X)_{\max}$ such that

$$|\Delta u_{X_k}| < (\Delta u_X)_{\max} \implies |I_{X_k}| < I_{\max}. \quad (10)$$

It is easy to show that such a bound is given by

$$(\Delta u_X)_{\max} = \frac{I_{\max}}{\|g_o\|_1} \quad (11)$$

where g_o is the impulse response of $G_o(z)$ and $\|g_o\|_1 = \sum_{k=0}^{\infty} |g_o(k)|$. The frequency response of this bound the dotted-black curve in Fig. 5. In practice, this bound is conservative, and an alternative is to choose²

$$(\Delta u_X)_{\max} = \|G_o(z)\|_{\infty} \approx 0.1980, \quad (12)$$

which results in the dashed-black curve in Fig. 5. Although (12) is only sufficient to guarantee (10) for sinusoidal inputs, in practice we find that enforcing (12) does lead to the current staying under 100 mA.

²For an arbitrary transfer matrix $G(z)$, we follow [37] and define $\|G(z)\|_{\infty} \triangleq \sup_{\phi \in [0, 2\pi]} \bar{\sigma}(G(e^{j\phi}))$, where $\bar{\sigma}(\cdot)$ yields the maximum singular value of its argument.

Finally, the slew-rate limit used in the MPC/linear feedback controller must be discounted from (12) to account for the inverse drift compensator. This adjustment for the inverse drift operator follows essentially the same argument as above. We have

$$(\Delta u)_{\max} \leq \frac{(\Delta u_X)_{\max}}{\|G_d^{-1}(z)\|_{\infty}} \approx 0.167 \quad (13)$$

4. Control Setup

The constraint (13) can be remodeled as a pure saturating constraint if an incremental form of G_{vib} is used that has as its input $\Delta u_k := u_k - u_{k-1}$, rather than u_k . This not only allows us to directly penalize the rate of change in the optimal control problem but also renders the constraint (13) as a box constraint on Δu , enabling the use of the computationally efficient FGM. Details on the form of the FGM we use can be found in [22, 38] while specifics about our implementation are discussed in [23].

4.1. The Incremental Form

To develop the required incremental form, the dynamics of $\hat{G}_{\text{vib}} = \{A, B, C, 0\}$ are augmented with a state x_{u_k} :

$$x_{u_k} = u_{k-1}.$$

It follows that

$$\bar{x}_{k+1} = \begin{bmatrix} A & B \\ 0 & 1 \end{bmatrix} \bar{x}_k + \begin{bmatrix} B \\ 1 \end{bmatrix} \Delta u_k \quad (14a)$$

$$y_k = \begin{bmatrix} C & 0 \end{bmatrix} \bar{x}_k \quad (14b)$$

$$\bar{x}_k := \begin{bmatrix} x_k \\ x_{u_k} \end{bmatrix}. \quad (14c)$$

We call this system $\bar{G} = \{\bar{A}, \bar{B}, \bar{C}, 0\}$, which has $\bar{n}_s = 23$ states, 9 of which model delay. To solve the setpoint tracking problem, we work in the error coordinates of \bar{G} . For a constant reference r_{ss} , in steady state we have $\Delta u_{ss} = 0$ and $\bar{x}_{ss} = N_{\bar{x}} r_{ss}$ where $N_{\bar{x}} \in \mathbb{R}^{\bar{n}_s}$ is found by solving

$$\begin{bmatrix} N_{\bar{x}} \\ N_u \end{bmatrix} = \begin{bmatrix} I - \bar{A} & -\bar{B} \\ \bar{C} & 0 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ I \end{bmatrix}, \quad (15)$$

which, due to the augmented pole at $z = 1$, will give $N_u \equiv 0$. The error state, $\bar{x}_{e_k} = \bar{x}_k - \bar{x}_{ss}$ has dynamics

$$\begin{aligned} \bar{x}_{e_{k+1}} &= \bar{A}\bar{x}_k + \bar{B}\Delta u_k - \bar{x}_{ss} \\ &= \bar{A}\bar{x}_{e_k} + \bar{B}\Delta u_k \end{aligned}$$

because \bar{x}_{ss} is in the nullspace of $(I - \bar{A})$.

4.2. Observer Design

To achieve zero-offset tracking (to constant disturbances), a disturbance estimator [39] is employed. The

disturbance dynamics are modeled as a pure integrating disturbance. The estimator dynamics are then given by

$$\begin{bmatrix} \hat{x}_{k+1} \\ \hat{d}_{k+1} \end{bmatrix} = A_m \begin{bmatrix} \hat{x}_k \\ \hat{d}_k \end{bmatrix} + B_m u_k + L_m (y_k - \hat{y}_k) \quad (16)$$

$$\hat{y}_k = C_m \begin{bmatrix} \hat{x}_k \\ \hat{d}_k \end{bmatrix} \quad (17)$$

where \hat{x} is the estimate of x_k (not \bar{x}_k), \hat{d}_k is the disturbance estimate, and

$$A_m = \begin{bmatrix} A & B_d \\ 0 & 1 \end{bmatrix}, \quad B_m = \begin{bmatrix} B \\ 0 \end{bmatrix} \\ C_m = [C \quad C_d], \quad L_m = \begin{bmatrix} L_x \\ L_d \end{bmatrix}. \quad (18)$$

It is shown in [39] that the gains L_x and L_d may be designed separately such that the closed-loop poles $A_m - L_m C_m$ are the same as $\sigma(A - L_x C) \cup \sigma(1 - L_d C_d)$. L_x is set equal to the steady-state solution of the discrete LQR problem applied to the dual of \hat{G}_{vib} , where $R = 1$, $Q = \alpha B B^T$ and α is a tuning parameter. L_d is designed such that the disturbance pole is placed at $z = 0.8$.

To achieve zero-offset tracking, disturbance estimators re-compute the steady-state target \bar{x}_{ss} at each time step. Here, an output disturbance model is used, so $B_d = 0$ and $C_d = I$, and the reference is adjusted by subtracting \hat{d}_k . In other words, at each time step, we need to compute

$$\bar{x}_e = \begin{bmatrix} \hat{x}_k \\ x_{u_k} \end{bmatrix} - N_{\bar{x}} (r_k - \hat{d}_k).$$

This is simpler than the case for an input disturbance model ($C_d = 0$ and $B_d \neq 0$), which involves an additional vector-scalar multiplication and a vector-vector addition (see (21) in [39]). It is shown in [39] that output disturbance and input disturbance estimators are equivalent provided 1 is not an eigenvalue of A , which is the case here, because the state x_u is not estimated. Thus, the output disturbance formulation is used for computational savings.

4.3. Closed-Loop Equations

When interpreting the experimental and simulation results in Section 6, it will be helpful to discuss the closed-loop behavior in terms of various transfer functions. To this end, the closed-loop equations are derived here for the block diagram of Fig. 1 with the drift and hysteresis operators set to the identity and with the controller given by a partitioned feedback gain $K = [K_x \quad K_u]$. Similarly, the observer gain is partitioned as in (18). In closed-loop x_u is not estimated because it is perfectly known, and the state-feedback portion of the control with \hat{d} is not computed because the disturbance is uncontrollable from u_k . Equations (14), (16), and (17) yield

$$\hat{x}_{k+1} = A \hat{x}_k + B u_k + L_x (y_k - \hat{y}_k) \quad (19)$$

$$\hat{d}_{k+1} = \hat{d}_k + L_d (y_k - \hat{y}_k) \quad (20)$$

$$x_{u_{k+1}} = x_{u_k} + \Delta u_k, \quad (21)$$

where y_k is the plant output. The control increment Δu_k and control u_k are given by

$$\Delta u_k = - [K_x \quad K_u] \begin{bmatrix} \hat{x}_k \\ x_{u_k} \end{bmatrix} + \bar{N} (r_k - \hat{d}_k). \quad (22)$$

$$u_k = \Delta u_k + x_{u_k}. \quad (23)$$

In (22), $\bar{N} \triangleq K N_{\bar{x}}$ is the feedforward control gain, where $N_{\bar{x}}$ is defined by (15). Equations (19)-(23) can be written as the combined state-space system

$$\tilde{x}_{k+1} = \tilde{A} \tilde{x}_k + \tilde{L} y_k + \tilde{B} \bar{N} r_k \quad (24)$$

$$u_k = -\tilde{K} \tilde{x}_k + \bar{N} r_k \quad (25)$$

where

$$\tilde{A} = \begin{bmatrix} A - BK_x - L_x C & B(1 - K_u) & -B\bar{N} - L_x C_d \\ -K_x & 1 - K_u & -\bar{N} \\ -L_d C & 0 & 1 - L_d C_d \end{bmatrix} \\ \tilde{L} = \begin{bmatrix} L_x \\ 0 \\ L_d \end{bmatrix}, \quad \tilde{B} = \begin{bmatrix} B \\ 1 \\ 0 \end{bmatrix}, \quad \tilde{x}_k = \begin{bmatrix} \hat{x}_k \\ x_{u_k} \\ \hat{d}_k \end{bmatrix}, \\ \tilde{K} = [K_x \quad K_u - 1 \quad \bar{N}].$$

Taking the \mathcal{Z} -transform of (24) and (25) yields

$$u(z) = \bar{N} (1 - \tilde{K} (zI - \tilde{A})^{-1} \tilde{B}) r(z) - \tilde{K} (zI - \tilde{A})^{-1} \tilde{L} y(z). \quad (26)$$

If $G(z)$ (which need not be the same as the model G_{vib}) is the transfer function of the plant, then the \mathcal{Z} -transform of the output y subject to an output disturbance d and control input u is $y(z) = G(z)u(z) + d(z)$. Combining this expression with (26) yields

$$y(z) = \frac{G(z)\bar{N}(1 - D_2(z))}{1 + G(z)D_1(z)} r(z) + \frac{1}{1 + G(z)D_1(z)} d(z) \quad (27)$$

where

$$D_1(z) = \tilde{K} (zI - \tilde{A})^{-1} \tilde{L} \\ D_2(z) = \tilde{K} (zI - \tilde{A})^{-1} \tilde{B}.$$

Thus, the loop gain is given by $L(z) = G(z)D_1(z)$ and the sensitivity function is defined as

$$\mathcal{S}(z) \triangleq \frac{1}{1 + L(z)}. \quad (28)$$

Due to the disturbance estimator, $\mathcal{S}(z)$ will always have a zero at $z = 1$ so that its DC-gain is zero. In Section 6, it will be helpful to quantify how large the sensitivity function gain is at small but non-zero frequencies. To this end, the ‘‘integrated sensitivity’’ is defined as

$$\mathcal{S}_I(z) \triangleq \mathcal{S}(z) \frac{1}{z - 1}. \quad (29)$$

The DC-gain of $\mathcal{S}_I(z)$ can then be used to quantify the low-frequency gain of $\mathcal{S}(z)$.

Recall that the closed-loop poles are the transmission zeros of $1 + L(z)$ and are the union of the controller and observer poles, which can be seen through the separation principle or manipulation of the matrix pencil describing the transmission zeros. Moreover, in the first term of (27), the observer poles are canceled by the transmission zeros of $(1 - D_2(z))$. Of course, these properties only hold when the observer uses a perfect model of the plant. The advantage in representing the closed-loop dynamics as (27) is that (i) it is valid when the observer and plant dynamics do not match and (ii) it exposes how to analyze the robustness of our designs in terms of using, e.g., the sensitivity function $\mathcal{S}(z)$ and stability margins derived from the loop gain $L(z)$.

5. Control Designs

Consider the optimal control problem

$$\min_v z_N^T P z_N + \sum_{i=0}^{N-1} z_i^T Q z_i + 2z_i^T S v_i + v_i^T R v_i \quad (30a)$$

$$\text{s.t. } z_{i+1} = \bar{A} z_i + \bar{B} v_i \quad (30b)$$

$$z_0 = [\hat{x}_k^T, x_{u_k}^T]^T - N_{\bar{x}}(r_k - \hat{d}_k) \quad (30c)$$

$$|v_i| \leq (\Delta u)_{\max}. \quad (30d)$$

where Q and R are symmetric matrices and the matrices Q, R, S satisfy (2) and (3). The terminal cost P is the solution of the DARE.

Two control strategies are considered based on (30):

Constrained Model Predictive Control (MPC): Here, one solves (30a) online, which results in a sequence of N optimal controls, $\{v_i\}_{i=0}^{N-1}$. One then sets $\Delta u_k = v_0$ and discards the remaining v_i . The process is repeated at the next time step. With MPC, the saturator in Fig. 1 is superfluous because the optimal control satisfies the constraints by design. The goal with MPC is to directly account for the slew-rate constraint as part of the control law itself and hope that this results in increased performance.

Saturated Linear state feedback (SLF): Here, the constraint (30d) is eliminated. Thus,

$$\Delta u_k = v_0 = -K z_0$$

where $K = (\bar{B}^T P \bar{B} + R)^{-1}(\bar{B}^T P \bar{A} + S^T)$, which is the LQR feedback gain associated with Q, R , and S . In contrast to MPC, with SLF, the saturator in Fig. 1 is necessary to avoid exceeding the current limit. While this scheme is much more computationally efficient than MPC, the fact that the constraint is not directly accounted for typically means that the performance requirements must be relaxed to maintain stability (e.g., by increasing the control weight R), due to the saturator.

In the next two subsections, the problem of designing the weighting matrices Q, R , and S is addressed.

5.1. Control Weight Selection

Both the MPC and SLF control schemes require the selection of weighting matrices Q, R , and S . Here, we propose a method which can produce well damped closed-loop pole locations (to mitigate residual vibration) and is also easily de-ratable. The method is a small extension to classic symmetric root locus (SRL) or reciprocal root locus (RRL) methods [40, 41]. Let $Q = f f^T$ for some vector $f \in \mathbb{R}^{\bar{n}_s}$, $S = f$, and $R = 1 + \gamma$ where γ is a scalar parameter. Consider the fictitious output

$$\eta_k = f^T \bar{x}_k + v_k. \quad (31)$$

Then the unconstrained version of (30a) can be written as

$$\sum_{i=0}^{\infty} \eta_i^T \eta_i + \gamma v_i^T v_i. \quad (32)$$

As the control weight γ decreases, the closed-loop poles of the unconstrained LQR move to the zeros of $\{A, B, f^T, 1\}$. In the standard RRL, $S = 0$ and there are at most $\bar{n}_s - 1$ finite zeros in the fictitious system and one zero at infinity.

Lemma 1. *Let $\{A, B, f^T, 1\}$ represent a state-space discrete-time system with n states. Let $Q = f f^T$, $S = f$ and $R = 1 + \gamma$ where $f \in \mathbb{R}^n$ and γ is a scalar. Let $K_\gamma = (B^T P B + R)^{-1}(B^T P A + S^T)$ solve the LQR problem associated with Q, R, S for a particular γ , where P solves the DARE*

$$Q = P - A^T P A - (A^T P B + S)(B^T P B + R)^{-1}(B^T P A + S^T). \quad (33)$$

Then as γ approaches zero, the closed-loop eigenvalues of $A - B K_\gamma$ approach the zeros of $\{A, B, f^T, 1\}$.

Proof: Using the same spectral factorization technique as [42, pg. 135] and [41, pg. 97] (which both prove the claim when $S = 0$), it can be shown that

$$\begin{aligned} & B^T (z^{-1}I - A^T)^{-1} Q (zI - A)^{-1} B \\ & + S^T (zI - A)^{-1} B + B^T (z^{-1}I - A^T) S + R \\ & = M + M K_\gamma (zI - A)^{-1} B + B^T (z^{-1}I - A^T)^{-1} K_\gamma^T M \\ & + B^T (z^{-1}I - A^T)^{-1} K_\gamma^T M K_\gamma (zI - A)^{-1} B \end{aligned} \quad (34)$$

where $M = B^T P B + R$. Define $g_o(z) = f^T (zI - A)^{-1} B$ and $\bar{g}_o(z^{-1}) = B^T (z^{-1}I - A^T)^{-1} f$. Using this in (34) yields

$$\begin{aligned} & (\bar{g}_o(z^{-1}) + 1)(g_o(z) + 1) + \gamma = \\ & [1 + B^T (z^{-1}I - A^T)^{-1} K_\gamma^T] M [1 + K_\gamma (zI - A)^{-1} B]. \end{aligned} \quad (35)$$

The zeros of the right-hand side are the n stable poles of the closed-loop system, together with their reflections about the unit circle. As γ approaches 0, the zeros of the left-hand side approach the zeros of $\{A, B, f^T, 1\}$ (and their reflections), which establishes the result. ■

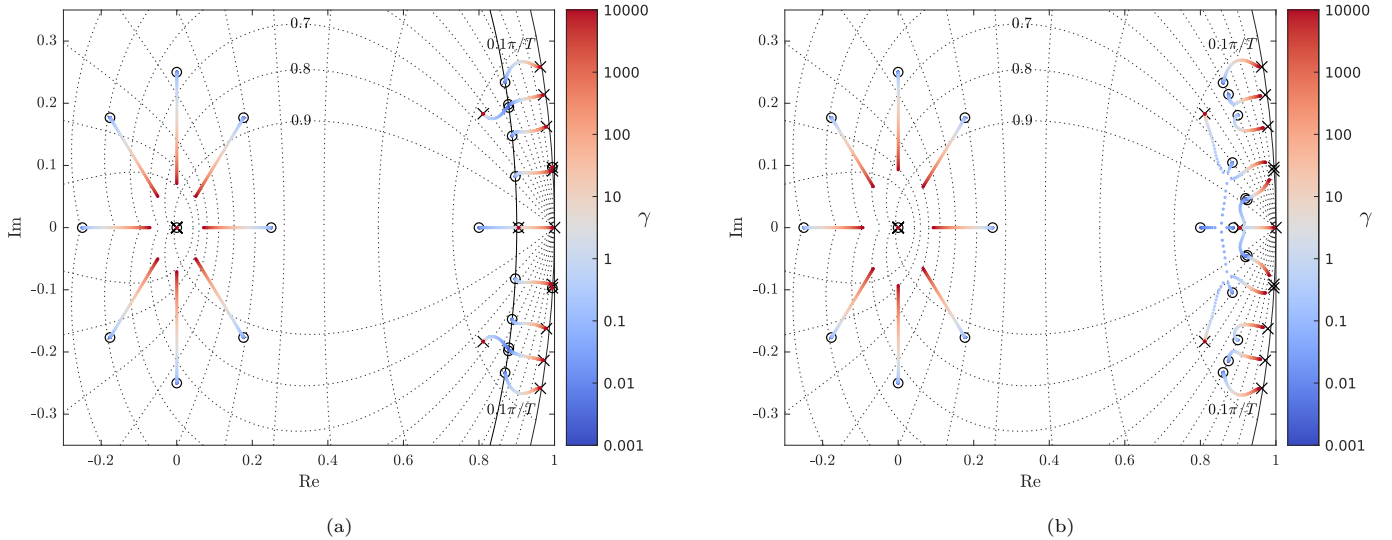


Figure 6: Root locus of closed-loop poles as a function of γ . For clarity, the plant zeros are not shown. The \times 's denote the open-loop plant poles. The \circ 's denote the fictitious zeros, which are at the desired pole locations. (a) Constant- σ scheme with $\sigma = 0.9$. (b) Chosen- ζ scheme.

This behavior is illustrated in Fig. 6 for two different f vectors (discussed below). Thus, pole-placement can be achieved through proper design of f and by taking γ to be small. With this selection of Q and S , (3) becomes

$$ff^T \left(1 - \frac{1}{1 + \gamma} \right) \geq 0 \quad (36)$$

which holds for all $\gamma \geq 0$. Numerical difficulties can arise in computing P when γ is too close to zero, though γ can usually be chosen small enough that the difference between the desired pole locations and their actual locations is negligible. The direct feedthrough in (31) leads to the cross-weighting term S and is necessary to endow the fictitious system with \bar{n}_s zeros in order to place all \bar{n}_s poles.

Through elementary block row and column operations, it is easy to show that the zeros of $\{A, B, f^T, 1\}$ are the same as the solutions of the generalized eigenvalue problem

$$\begin{bmatrix} (A - Bf^T) & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} zI & 0 \\ 0 & 0 \end{bmatrix} = 0.$$

Thus, f may be found via pole placement techniques.

Two methods are considered to choose a set of desired pole locations. The first method, called ‘‘constant- σ ’’ (CS), moves all complex poles such that they are projected radially inward to lie on a circle with a specified radius, which endows them all with the same time constant. Here, the radius is selected as $\sigma = 0.9$. The continuous-time equivalent (which may provide better intuition) would move the poles onto a vertical line at 420 Hz, which is slightly faster than the first open-loop mode at 350 Hz.

The second scheme, called ‘‘choose- ζ ’’ (CZ), keeps the natural frequency of each complex pole unchanged but specifies a damping ratio ζ . In order of increasing natural frequency, the damping ratios are specified as $\{0.85, 0.85, 0.7, 0.4, 0.4, 0.4\}$, so that the slowest poles

have the highest damping. In the next section, this approach is shown to generally yield slower settle times than the constant- σ method; however it tends to result in slightly smaller residual oscillations and is the same scheme we considered in [43, 23], and is included here so that results can be evaluated in light of our prior work.

For both cases, 8 of the 9 poles corresponding to delay are placed at roots of unity with a radius chosen more or less arbitrarily at $\sigma = 0.25$, with the remaining pole left at the origin. Fig. 6 shows root locus-like plots for each scenario as a function of γ . As γ approaches zero, the closed-loop poles approach the designed fictitious zeros, which are indicated by black circles.

Clearly, there is room for variation on these two schemes. For example, one could try to speed things up by decreasing σ (in CS) or by also increasing the natural frequencies in the CZ scheme. However, given that the control horizon is limited to $N = 22$ (due to hardware limitations), these more aggressive strategies generally require much larger control weights to achieve stability through the entire desired operating range.

Certainly, in the case of the SLF state feedback controller, one could use a pole-placement design to start with. However, the method here has two advantages: (i) it permits a straightforward comparison to the MPC design (which, in the present formulation, requires weighting matrices, not pole locations) and (ii) the design becomes parameterized by the scalar parameter γ . Increasing γ increases the control penalty, which makes de-rating the design easy. In contrast, with standard pole-placement, it is less clear how to ‘‘back-off’’ the design in a systematic way if the slew-rate constraint is violated to the extent that instability results. Intuitively, one could e.g., try increasing σ in this situation. However, this could require a different σ for each pole (since the open-loop poles do not all have the same σ). It is also unclear if this method

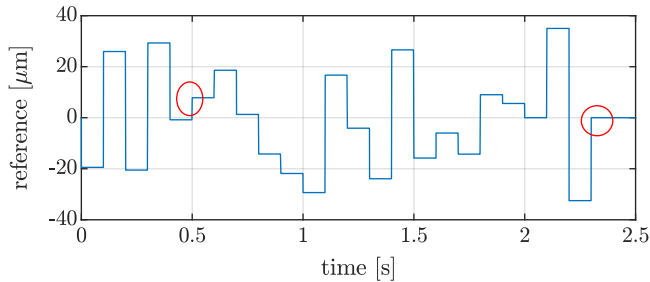


Figure 7: The sequence of step commands used to test the control laws. The red circles indicate the 6th and 24th reference; responses to these references are shown in closer detail in Figs. 10a and 10b.

guarantees a reduction in control action.

5.2. Selecting γ

The choice of γ has a significant effect on system performance, which differs between simulation and experiment. This aspect of the tuning is explored in the next section.

6. Simulated vs. Experimental Performance and Dependence on γ

The goal of this section is to explore how the simulation and experimental performances of MPC and SLF (for both pole placement schemes) depend on the control weight γ and to determine to what extent, if any, MPC provides a benefit. In all simulations and experiments, the MPC control horizon is $N = 22$.

The experimental and simulation performances of the MPC and SLF control schemes were tested using a sequence of 24 reference commands. The first 20 references were selected randomly within the range $[-32.5 \mu\text{m}, 35 \mu\text{m}]$. The final four references were chosen as $0 \mu\text{m}$, $35 \mu\text{m}$, $-32.5 \mu\text{m}$, and $0 \mu\text{m}$, to exercise the full usable range of the stage.³ This sequence of references is shown in Fig. 7. Using a sequence of randomly generated references (rather than, e.g., a single reference or several references beginning from rest) has several benefits: (i) it is representative of the type of references seen in a compressive sensing imaging scenario [18, 20]; (ii) it eliminates the possibility of inadvertently picking reference values where one control law does better; and (iii) using a *sequence* of references will help to draw out the deleterious effects of imperfections in the hysteresis and drift inversion.

To explore performance dependence on γ , a series of experiments and simulations were run across a grid of γ 's from $\gamma = 10^{-5}$ to $\gamma = 400$. Each simulation was run once and each experiment was conducted 8 times. The settle time is defined in an absolute sense: the settle boundary is $\frac{70}{512} \mu\text{m}$, corresponding to settling within one pixel of a 512

³For references larger than $35 \mu\text{m}$, the control signal saturates; for references smaller than $-32.5 \mu\text{m}$, the sensor saturates due to a bias in the stage.

by 512 pixel image for the given range. Let $t_i^j(\gamma_\ell)$ represent the settle time of the i th reference for the j th experimental run using the ℓ th γ in the grid. Then the total settle time of the j th experiment is $T_j = \sum_{i=1}^{24} t_i^j(\gamma_\ell)$. The sample mean of the total settle times for a specific γ_ℓ is

$$\bar{T}(\gamma_\ell) = \frac{1}{8} \sum_{j=1}^8 T_j(\gamma_\ell). \quad (37)$$

Fig. 8 plots the mean of the total experimental settle times (i.e., $\bar{T}(\gamma_\ell)$) vs. γ as the red (MPC) and black (SLF) dots with error bars. The results of the CS and CZ schemes are shown on the left and right, respectively. The simulation settle times are shown in Fig. 9. The MPC simulations are denoted by the red circles and the SLF simulations are denoted by the black dots. For reference, the values of the plotted data are given in Tables 2a and 2b.

Note that the γ values used for CS and CZ differ. Simulations indicated that, with a control horizon⁴ of $N = 22$, the MPC would fail to stabilize the system for the largest setpoints if γ was taken much smaller than 10^{-3} (resp., 10^{-5}) for CS (resp., CZ). Similarly, due to saturating Δu_k , the SLF control laws are unstable for γ smaller than about 7.5 (resp., 3.5) for CS (resp., CZ). The γ values 46.4 (for CS) and 50.9 (for CZ) are discussed in Section 6.1.

For the smallest tested values of γ , Fig. 10a shows zoomed-in experimental trajectories of the four controllers for references 6 and 24 (circled in Fig. 7). The fastest experimental settle times occur for $\gamma = 100$ for the CS scheme and $\gamma = 25$ for the CZ scheme. Zoomed-in experimental trajectories for these values of γ are shown in Fig. 10b, also for references 6 and 24.

For all of the experimental trajectories, the largest measured power amplifier current was 98.5 mA, indicating success in respecting the 100 mA current limit.

6.1. Discussion

In both simulation and experiment, MPC is able to utilize a much smaller γ than SLF. In the simulation results, increasing γ results in an increased total settle time. This is as one would expect, because nominal closed-loop bandwidth decreases as γ increases, as shown by the solid black curves in Fig. 9. Interestingly, this trend does not hold in the experimental results. Up to about $\gamma = 100$ for CS and $\gamma = 25$ for CZ, total settle time *decreases* as γ increases in the experimental results. When $\gamma = 10^{-5}$ in the CZ MPC experiment, the total settle time is nearly 3 times slower than the simulation; when $\gamma = 25$ the experiment is only 1.3 times slower than simulation.

These trends can be explained by analyzing how robustness depends on γ . First, recall from Sections 3.2 and 3.3 that the overall control loop inverts imperfect models of

⁴ $N = 22$ is the largest control horizon our FPGA can sustain, and the stabilized region generally improves with a larger N . See [44] for a detailed discussion of the horizon length for this system.

Table 2: Total settle times for (a) the constant- σ and (b) choose- ζ state weighting schemes. The second and third columns of each table are the simulation results while the fourth and fifth columns are the average of 8 experimental runs. SLF was not run for the smallest γ values, as indicated by the dashes. All times are in milliseconds.

(a) Constant- σ					(b) Choose- ζ				
γ	MPC-sim	SLF-sim	MPC-exp	SLF-exp	γ	MPC-sim	SLF-sim	MPC-exp	SLF-exp
10^{-3}	112.6	—	252.6	—	10^{-5}	141.8	—	423.9	—
10^{-2}	112.7	—	254.2	—	10^{-2}	141.9	—	420.3	—
10^{-1}	112.7	—	252.4	—	10^{-1}	142.1	—	420.2	—
1	113.3	—	250.3	—	1	144.5	—	409.6	—
7.5	116.5	118.6	243.5	242.4	3.5	159.9	161.3	368.6	396.7
10	117.6	119.2	242.4	241.7	10	190.9	191.7	319.6	320.0
25	129.4	130.2	237.5	236.6	25	229.3	229.7	305.6	305.9
46.4	138.2	138.2	232.3	232.7	50.9	279.9	279.9	315.6	315.2
75	151.0	151.1	229.2	229.3	75	318.4	318.7	337.9	337.7
100	159.2	159.3	228.4	227.9	100	353.1	353.2	362.8	362.9
200	186.5	186.6	243.3	242.1	200	464.3	464.6	446.1	443.3
300	215.2	216.1	256.9	257.3	300	552.3	552.3	524.2	519.8
400	236.8	237.0	268.5	268.3	400	627.3	627.6	582.7	581.0

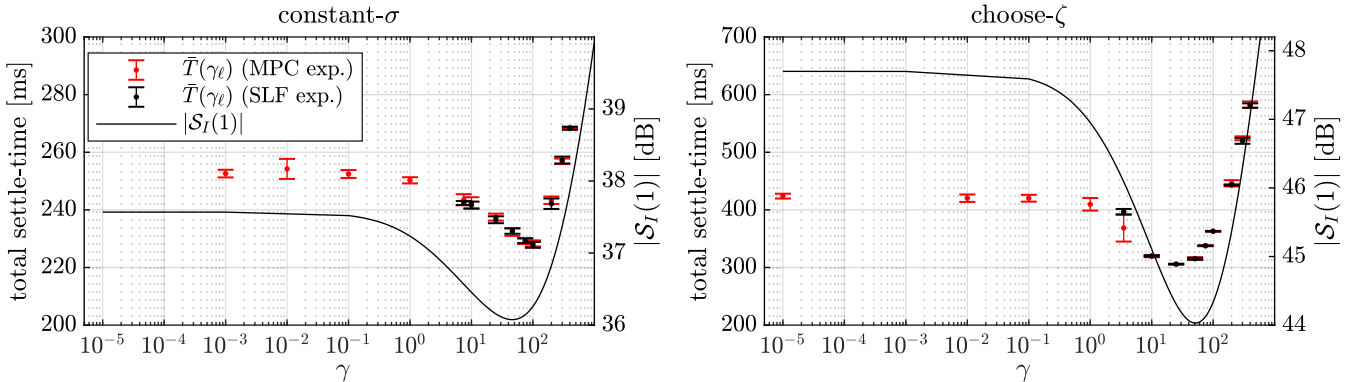


Figure 8: Experimental results, (left) constant- σ , (right) choose- ζ : The red (MPC) and black (SLF) dots with error bars are the sample means of the total settle times for different values of γ and are plotted against the left axes. The error bars represent one standard deviation. The solid black curve is the DC-gain of the integrated sensitivity function $S_I(1)$ evaluated at different values of γ , plotted against the right axes.

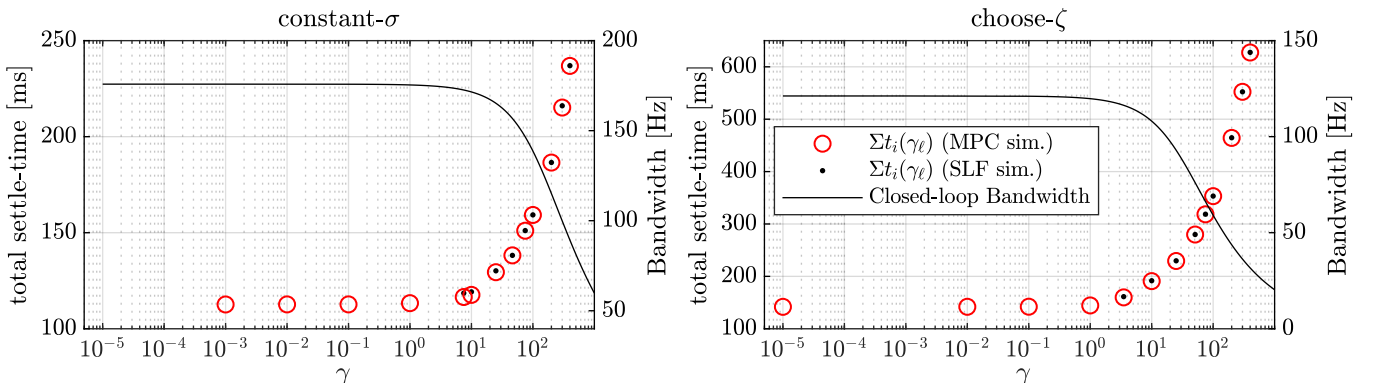


Figure 9: Simulation results, (left) constant- σ , (right) choose- ζ : The red circles (MPC) and black dots (SLF) are the total simulated settle times for different values of γ and are plotted against the left axes. The solid black curve is the closed-loop 3 dB bandwidth evaluated for different values of γ , plotted against the right axes.

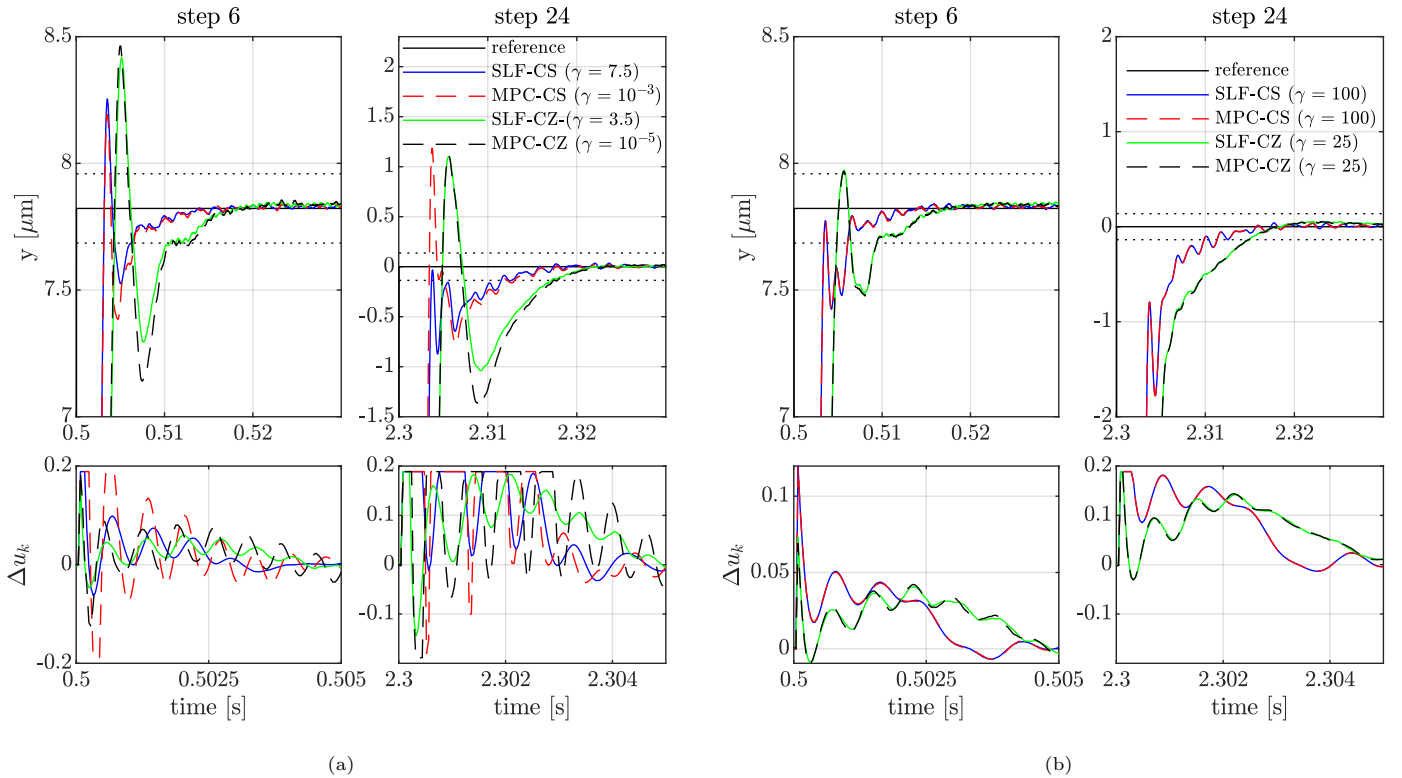


Figure 10: Stage output (top row) and control increment (bottom row) experimental trajectories for references 6 and 24 (circled in Fig. 7). The dotted black lines indicate the settle boundary. (a) Trajectories for the smallest tested values of γ . (b) Trajectories for the values of γ resulting in the fastest total experimental settle time. As indicated in (12), $(\Delta u_X)_{\max} = 0.1980$.

hysteresis and drift. Errors in these inversions will show up as model uncertainty at low frequencies. When the gain of the sensitivity function at low frequencies is small, then the effect of these uncertainties will be reduced. To quantify how the low-frequency gain of \mathcal{S} (which is always zero at DC) depends on γ , the DC gain of the integrated sensitivity function \mathcal{S}_I , defined in (29), is computed at a grid of γ values. The resulting parametric plot is shown as the black curves in Fig. 8, which are plotted against the right axes. The experiments using $\gamma = 46.4$ (for CS) and $\gamma = 50.9$ (for CZ) correspond to the minimum of the respective $|\mathcal{S}_I(1)|$ curves. Although the decrease and subsequent increase in total settle time roughly follows the $|\mathcal{S}_I(1)|$ vs. γ curve, the fastest experimental settle times do not exactly match to the minima of $|\mathcal{S}_I(1)|$.

A complementary analysis is to consider how the gain margin (GM) and phase margin (PM) change with γ . Fig. 11 shows parametric plots of these metrics (computed from the loop gain $L(z)$). For both pole-placement scenarios, robustness in terms of GM and PM increases monotonically with γ and is quite poor for the smallest γ 's.

Ultimately, the fastest experimental settle time which is achieved for either pole-placement scenario is a trade-off between the robustness metrics GM, PM, and $|\mathcal{S}_I(1)|$ and the decrease in nominal closed-loop bandwidth as γ increases. Determining the precise nature of this trade-off would require precise knowledge of the model uncertainty,

which is unknown. However, qualitatively similar results to those shown in Fig. 8 can be obtained by running simulations with small perturbations to the plant gain. This is shown in Fig. 12 for the CZ scheme (the results for the CS scheme are similar, but not shown for space reasons). Even for the 2.5% plant perturbation, the best settle time occurs at $\gamma = 10$, rather than 10^{-5} .

For the fastest experimental settle times (i.e., $\gamma = 100$ for CS and $\gamma = 25$ for CZ), the difference in trajectories and thus total settle times between MPC and SLF is negligible. The close correspondence of MPC and SLF trajectories for large γ is illustrated in Fig. 10b for references 6 and 24. The similarity in total settle times between MPC and SLF holds for both simulation and experiment as can be seen in the $\gamma = 100$ row of Table 2a and the $\gamma = 25$ row of Table 2b. We conclude that, under the present state weighting schemes and our system with modeling uncertainties, input constrained MPC does not reduce the overall settle time compared to SLF. For systems with less model uncertainty, the input constrained MPC may provide more advantages. Comparisons with robust MPC formulations [45, 46], while outside the scope of the current work, would be of interest to explore in the future.

6.2. Possible Objections

We envision three main criticisms to these results and analysis. First, one might point out that both MPC and

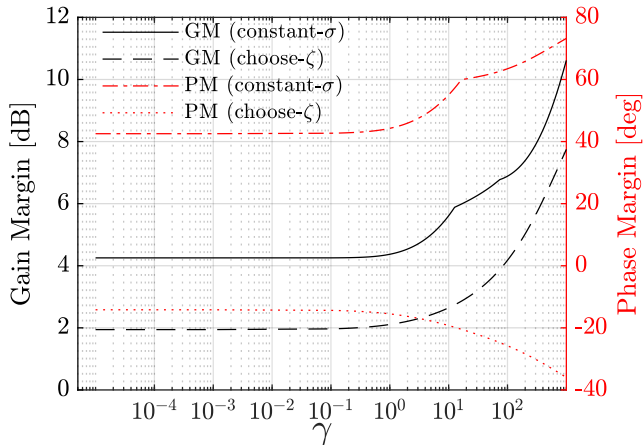


Figure 11: GM and PM dependence on γ . PM is computed as the smallest absolute difference between $\angle L(e^{j\omega T_s})$ and 180° for ω such that $|L(e^{j\omega T_s})| = 1$. Note that both loops are nominally stable: negative PM means that a positive perturbation in phase larger than |PM| at the crossover frequency will lead to instability.

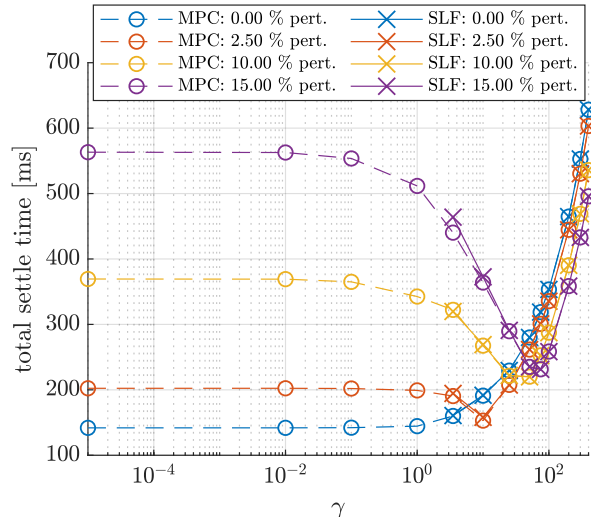


Figure 12: Total simulation settle times using the CZ control scheme with perturbations to the plant gain as indicated by the legend.

SLF are in general, non-linear, making the classical metrics like GM, PM, and $|\mathcal{S}_I(1)|$ not applicable. However, for any setpoint, both SLF and MPC behave linearly in a neighborhood around that setpoint. In other words, while the beginning of a setpoint tracking maneuver may saturate the SLF controller or put the MPC on its constraint boundary, the latter part of that same maneuver will be governed by linear dynamics. In that sense, the metrics GM and PM and the DC gain of the integrated sensitivity still provide insight, which we believe is manifested, e.g., in the correlation between overall settle time and $|\mathcal{S}_I(1)|$.

A related objection is that the relatively poor GMs and PMs could be improved through a loop transfer recovery (LTR) method. However, the classic method as proposed in [47] does not translate directly to discrete-time. Most efforts at obtaining LTR-like results for discrete-time assume that (i) a current estimator is used and (ii) $CB \neq 0$ [48, 49, 50]. When considering MPC, especially at high sample rates, a prediction estimator must be used, since computing the control action requires a significant portion of the sample period. Also, for our system model, $CB = 0$, because the relative degree is 12. These two factors preclude discrete-time LTR. Even if these issues could be circumvented, in contrast to the continuous-time case, the full-state feedback discrete-time LQR controller does not have a guaranteed infinite gain margin nor a guaranteed phase margin greater than 60° . Rather, the guaranteed bounds for GM and PM depend on the DARE solution P [42, p. 136], which in turn depends on γ .

Finally, one might argue that time-optimal MPC (TOMPC) [51] might show some benefit over MPC with a purely quadratic cost. Unfortunately, TOMPC requires the solution of a *sequence* of optimal control problems with long horizons at each time step [51]. Despite advances in hardware and algorithms, such methods remain too slow for systems with fast dynamics.

7. Conclusions, Discussions, and Future Work

This paper has developed, compared, and analyzed setpoint tracking control laws for a piezo stage, subject to a current limitation in the piezo power amplifier. Using an incremental form with Δu_k as the control input and conservatively constraining Δu_k enabled the power amplifier current to be kept within the 100 mA limitation. Model predictive control (MPC) and saturated linear feedback (SLF) control laws were developed, and two separate state weighting schemes were outlined that allow the aggressiveness of the controllers to be easily adjusted. In agreement with earlier work [23], the MPC controllers were able to utilize more aggressive control weights.

The controllers were implemented and evaluated on our piezo stage across a wide range of control weights in both simulations and experiments. The experimental results for both MPC and SLF generally yielded settle times significantly slower than the simulations predicted, and the most aggressive control weights did not yield the best experimental settle times. These results highlight the importance of robustness in experimental performance. Analyzing both the MPC and SLF using the classical notions of gain margin, phase margin, and gain of the sensitivity function at low frequencies showed their usefulness in assessing the robustness of the controllers as a function of the control weight. Given the model uncertainty in the system, and the de-rating of both the MPC and SLF controllers needed to improve robustness, the best settle times achieved experimentally are comparable between the MPC and SLF approaches. Given the lower computational cost of the SLF controller, it is the preferred method based upon the results and analysis in this paper. Further exploration, development, and comparison of classical robust controller formulations [52, 53] with robust MPC approaches [45, 46] are an area of future work. Moreover, comparisons with other control approaches such

as H^∞ [54, 55, 56, 57, 58], repetitive [59, 60, 61, 62], and resonant [63, 64, 65, 66] controllers would provide even greater understanding of the advantages and disadvantages of these control approaches relative to one another for piezo actuators in AFMs. The use of time-optimal minimum jerk trajectories with a rate constraint, using trajectory generation methods such as [67], with the linear LQG controller should also be further considered and evaluated in the future.

MPC could possibly yield improved results if, rather than constraining the slew rate of the low-voltage input, the output of a power amplifier current model was constrained. Such an MPC problem could no longer be solved with the FGM, though with a larger FPGA, it may be possible that the problem could be solved with the Alternating Direction Method of Multipliers (ADMM) [22]. Such a formulation changes the constraints but not the cost function. Thus, the robustness properties explored here would remain unchanged, though it is possible they would play a smaller role. Since embedded platforms with FPGAs capable of implementing the ADMM at the required sample rate are still expensive, while re-formulating MPC in this way may improve performance over SLF, a holistic systems-level optimization should weigh this improvement against the cost of upgrading the power amplifier.

Declaration of Competing Interest

None.

Acknowledgements

The author thanks Roger A. Braker for the significant work he did on the material documented here while he was a PhD student. This work was supported in part by the US National Science Foundation (CMMI-1234980), Agilent Technologies, Inc., a Palmer Endowed Chair Professorship, and the Hanse Wissenschaftskolleg Institute for Advanced Study.

References

- [1] D. Y. Abramovitch, S. B. Andersson, L. Y. Pao, G. Schitter, A tutorial on the mechanisms, dynamics, and control of atomic force microscopes, in: Proc. American Control Conf., 2007, pp. 3488–3502. doi:10.1109/ACC.2007.4282300.
- [2] M. S. Rana, H. R. Pota, I. R. Petersen, Improvement in the imaging performance of atomic force microscopy: A survey, IEEE Trans. Automation Science and Eng. 14 (2017) 1265–1285. doi:10.1109/TASE.2016.2538319.
- [3] G. Schitter, K. J. Astrom, B. E. DeMartini, P. J. Thurner, K. L. Turner, P. K. Hansma, Design and modeling of a high-speed AFM-scanner, IEEE Trans. Control Systems Tech. 15 (2007) 906–915. doi:10.1109/TCST.2007.902953.
- [4] B. J. Kenton, K. K. Leang, Design and control of a three-axis serial-kinematic high-bandwidth nanopositioner, IEEE/ASME Trans. Mechatronics 17 (2012) 356–369. doi:10.1109/TMECH.2011.2105499.
- [5] J. Butterworth, L. Y. Pao, D. Abramovitch, Dual-adaptive feed-forward control for raster tracking with applications to AFMs, in: Proc. IEEE Int. Conf. Control Applications, 2011, pp. 1081–1087. doi:10.1109/CCA.2011.6044387.
- [6] Y. Li, J. Bechhoefer, Feedforward control of a closed-loop piezoelectric translation stage for atomic force microscope, Rev. Scientific Instruments 78 (2007) 013702. doi:10.1063/1.2403839.
- [7] K. K. Leang, Q. Zou, S. Devasia, Feedforward control of piezoactuators in atomic force microscope systems, IEEE Control Systems Magazine 29 (2009) 70–82. doi:10.1109/MCS.2008.930922.
- [8] Y. K. Yong, S. O. R. Moheimani, Collocated z-axis control of a high-speed nanopositioner for video-rate atomic force microscopy, IEEE Trans. Nanotech. 14 (2015) 338–345. doi:10.1109/TNANO.2015.2394327.
- [9] I. A. Mahmood, S. O. R. Moheimani, Fast spiral-scan atomic force microscopy, Nanotechnology 20 (2009) 365503. doi:10.1088/0957-4484/20/36/365503.
- [10] T. Tuma, J. Lygeros, V. Kartik, A. Sebastian, A. Pantazi, High-speed multiresolution scanning probe microscopy based on Lissajous scan trajectories, Nanotechnology 23 (2012) 185501. doi:10.1088/0957-4484/23/18/185501.
- [11] M. S. Rana, H. R. Pota, I. R. Petersen, Spiral scanning with improved control for faster imaging of AFM, IEEE Trans. Nanotech. 13 (2014) 541–550. doi:10.1109/TNANO.2014.2309653.
- [12] A. Fleming, B. Kenton, K. Leang, Bridging the gap between conventional and video-speed scanning probe microscopes, Ultramicroscopy 110 (2010) 1205–1214. doi:10.1016/j.ultramic.2010.04.016.
- [13] P. Huang, S. B. Andersson, Note: Fast imaging of DNA in atomic force microscopy enabled by a local raster scan algorithm, Rev. Scientific Instruments 85 (2014) 066101. doi:10.1063/1.4881682.
- [14] B. Hartman, S. B. Andersson, Feature tracking for high-speed AFM: Experimental demonstration, in: Proc. American Control Conf., 2017, pp. 773–778. doi:10.23919/ACC.2017.7963046.
- [15] C. S. Oxvig, T. Arildsen, T. Larsen, Structure assisted compressed sensing reconstruction of undersampled AFM images, Ultramicroscopy 172 (2017) 1–9. doi:10.1016/j.ultramic.2016.09.011.
- [16] S. Andersson, L. Pao, Non-raster sampling in atomic force microscopy: A compressed sensing approach, in: Proc. American Control Conf., 2012, pp. 2485–2490. doi:10.1109/ACC.2012.6315406.
- [17] B. Song, N. Xi, R. Yang, K. W. C. Lai, C. Qu, Video rate atomic force microscopy (AFM) imaging using compressive sensing, in: Proc. IEEE Int. Conf. Nanotech., 2011, pp. 1056–1059. doi:10.1109/NANO.2011.6144587.
- [18] R. A. Braker, Y. Luo, L. Y. Pao, S. B. Andersson, Hardware demonstration of atomic force microscopy imaging via compressive sensing and μ -path scans, in: Proc. American Control Conf., 2018, pp. 6037–6042. doi:10.23919/ACC.2018.8431873.
- [19] B. D. Maxwell, S. B. Andersson, A compressed sensing measurement matrix for atomic force microscopy, in: Proc. American Control Conf., 2014, pp. 1631–1636. doi:10.1109/ACC.2014.6858710.
- [20] R. A. Braker, Y. Luo, L. Y. Pao, S. B. Andersson, Improving the image acquisition rate of an atomic force microscope through spatial sub-sampling and reconstruction, IEEE/ASME Trans. Mechatronics 25 (2020) 570–580. doi:10.1109/TMECH.2020.2974251.
- [21] J. P. Killgore, D. G. Yablon, A. H. Tsou, A. Gannepalli, P. A. Yuya, J. A. Turner, R. Proksch, D. C. Hurlley, Viscoelastic property mapping with contact resonance force microscopy, Langmuir 27 (2011) 13983–13987. doi:10.1021/la203434w.
- [22] J. L. Jerez, P. J. Goulart, S. Richter, G. A. Constantinides, E. C. Kerrigan, M. Morari, Embedded online optimization for

- model predictive control at megahertz rates, *IEEE Trans. Automatic Control* 59 (2014) 3238–3251. doi:10.1109/TAC.2014.2351991.
- [23] R. A. Braker, L. Y. Pao, An application of the fast gradient method to model predictive control of an atomic force microscope X-Y stage, in: *Proc. IEEE Conf. Control Tech. and Applications*, 2017, pp. 111–116. doi:10.1109/CCTA.2017.8062449.
- [24] A. Wills, D. Bates, A. Fleming, B. Ninness, R. Moheimani, Application of MPC to an active structure using sampling rates up to 25 kHz, in: *Proc. IEEE Conf. Decision and Control*, 2005, pp. 3176–3181. doi:10.1109/CDC.2005.1582650.
- [25] C. Y. Lin, Y. C. Liu, Precision tracking control and constraint handling of mechatronic servo systems using model predictive control, *IEEE/ASME Trans. Mechatronics* 17 (2012) 593–605. doi:10.1109/TMECH.2011.2111376.
- [26] M. S. Rana, H. R. Pota, I. R. Petersen, H. Habibullah, Effect of improved tracking for atomic force microscope on piezo nonlinear behavior, *Asian J. Control* 17 (2015) 747–761. doi:10.1002/asjc.924.
- [27] L. O. Santos, L. T. Biegler, J. A. Castro, A tool to analyze robust stability for constrained MPC, *IFAC Proceedings Volumes* 37 (2004) 487–492. doi:10.1016/S1474-6670(17)38779-7.
- [28] M. A. Rodrigues, D. Odloak, MPC for stable linear systems with model uncertainty, *Automatica* 39 (2003) 569–583. doi:10.1016/S0005-1098(02)00176-0.
- [29] M. V. Kothare, V. Balakrishnan, M. Morari, Robust constrained model predictive control using linear matrix inequalities, *Automatica* 32 (1996) 1361–1379. doi:10.1016/0005-1098(96)00063-5.
- [30] D. Croft, G. Shed, S. Devasia, Creep, hysteresis, and vibration compensation for piezoactuators: Atomic force microscopy application, *J. Dynamic Systems, Measurement, and Control* 123 (1999) 35–43.
- [31] R. N. Jacques, D. W. Miller, Multivariable model identification from frequency response data, in: *Proc. IEEE Conf. Decision and Control*, 1993, pp. 3046–3051. doi:10.1109/CDC.1993.325762.
- [32] M. D. Sidman, F. E. DeAngelis, G. C. Verghese, Parametric system identification on logarithmic frequency response data, *IEEE Trans. Automatic Control* 36 (1991) 1065–1070. doi:10.1109/9.83539.
- [33] Y. Liu, J. Shan, N. Qi, Creep modeling and identification for piezoelectric actuators based on fractional-order system, *Mechatronics* 23 (2013) 840–847. doi:10.1016/j.mechatronics.2013.04.008.
- [34] M. Rakotondrabe, Bouc-Wen modeling and inverse multiplicative structure to compensate hysteresis nonlinearity in piezoelectric actuators, *IEEE Trans. Automation Science and Eng.* 8 (2011) 428–431. doi:10.1109/TASE.2010.2081979.
- [35] Y. Liu, J. Shan, U. Gabbert, N. Qi, Hysteresis and creep modeling and compensation for a piezoelectric actuator using a fractional-order Maxwell resistive capacitor approach, *Smart Materials and Structures* 22 (2013) 115020.
- [36] K. Kuhnen, Modeling, identification and compensation of complex hysteretic nonlinearities: A modified Prandtl-Ishlinskii approach, *European J. Control* 9 (2003) 407–418. doi:10.3166/ejc.9.407-418.
- [37] K. M. Grigoriadis, J. T. Watson, Reduced-order H_∞ and L_2 - L_∞ filtering via linear matrix inequalities, *IEEE Trans. Aerospace and Electronic Sys.* 33 (1997) 1326–1338. doi:10.1109/7.625133.
- [38] J. L. Jerez, P. J. Goulart, S. Richter, G. A. Constantinides, E. C. Kerrigan, M. Morari, Embedded Predictive Control on an FPGA using the Fast Gradient Method, in: *Proc. European Control Conf.*, 2013, pp. 3614–3620.
- [39] U. Maeder, M. Morari, Offset-free reference tracking for predictive controllers, in: *Proc. IEEE Conf. Decision and Control*, 2007, pp. 5252–5257. doi:10.1109/CDC.2007.4434696.
- [40] G. F. Franklin, J. D. Powell, M. Workman, *Digital Control of Dynamic Systems*, Ellis-Kagle Press, Half Moon Bay, CA, 1998.
- [41] F. L. Lewis, D. L. Vrabie, V. L. Syrmos, *Optimal Control*, John Wiley and Sons, Hoboken, N.J., 2012.
- [42] B. D. Anderson, J. B. Moore, *Optimal Control: Linear Quadratic Methods*, Prentice-Hall, Englewood Cliffs, N.J., 1989.
- [43] R. A. Braker, L. Y. Pao, Fast setpoint tracking of an atomic force microscope x-y stage via optimal trajectory tracking, in: *Proc. American Control Conf.*, 2017, pp. 2875–2881.
- [44] R. A. Braker, *Control Methods for Compressive Sensing in Atomic Force Microscopy*, Ph.D. thesis, University of Colorado, Boulder, CO, 2019.
- [45] A. Bemporad, M. Morari, Robust model predictive control: A survey, *Robustness in Identification and Control*, A. Garulli, A. Tesi, and A. Vicino, Eds., number 245 in *Lecture Notes in Control and Information Sciences*, Springer-Verlag (1999) 207–226.
- [46] M. B. Saltik, L. Ozkan, J. H. A. Ludlage, S. Weiland, P. M. J. Van den Hof, An outlook on robust model predictive control algorithms: Reflections on performance and computational aspects, *J. Process Control* 61 (2018) 77–102.
- [47] J. Doyle, G. Stein, Robustness with observers, *IEEE Trans. Automatic Control* 24 (1979) 607–611. doi:10.1109/TAC.1979.1102095.
- [48] J. Maciejowski, Asymptotic recovery for discrete-time systems, *IEEE Trans. Automatic Control* 30 (1985) 602–605. doi:10.1109/TAC.1985.1104010.
- [49] F. H. Guaracy, L. H. Ferreira, C. A. Pinheiro, The discrete-time controller for the H_∞ /LTR problem with mixed-sensitivity properties, *Automatica* 58 (2015) 28–31. doi:10.1016/j.automatica.2015.04.030.
- [50] T. Ishihara, H. Takeda, Loop transfer recovery techniques for discrete-time optimal regulators using prediction estimators, *IEEE Trans. Automatic Control* 31 (1986) 1149–1151. doi:10.1109/TAC.1986.1104197.
- [51] L. V. den Broeck, M. Diehl, J. Swevers, Time optimal MPC for mechatronic applications, in: *Proc. IEEE Conf. Decision and Control and Chinese Control Conf.*, 2009, pp. 8040–8045. doi:10.1109/CDC.2009.5400667.
- [52] K. Zhou, J. Doyle, *Essentials of Robust Control*, Prentice Hall, 1998.
- [53] S. Skogestad, I. Postlethwaite, *Multivariable Feedback Control: Analysis and Design*, 2 ed., Wiley, 2005.
- [54] G. Schitter, A. Stemmer, F. Allgöwer, Robust two-degree-of-freedom control of an atomic force microscope, *Asian J. Control* 6 (2004) 156–163. doi:10.1111/j.1934-6093.2004.tb00194.x.
- [55] J. A. Butterworth, L. Y. Pao, D. Y. Abramovitch, Architectures for tracking control in atomic force microscopes, *IFAC Proceedings Volumes* 41 (2008) 8236–8250. doi:10.3182/20080706-5-KR-1001.01394.
- [56] H. Ladjal, J.-L. Hanus, A. Ferreira, H_∞ robustification control of existing piezoelectric-stack actuated nanomanipulators, in: *Proc. IEEE Int. Conf. Robotics and Automation*, 2009, pp. 3353–3358. doi:10.1109/ROBOT.2009.5152701.
- [57] N. Chuang, I. R. Petersen, H. R. Pota, Robust H_∞ control in fast atomic force microscopy, in: *Proc. American Control Conf.*, 2011, pp. 2258–2265. doi:10.1109/ACC.2011.5991155.
- [58] C. Kammer, A. P. Nievergelt, G. E. Fantner, A. Karimi, Data-driven controller design for atomic-force microscopy, *IFAC-PapersOnLine* 50 (2017) 10437–10442. doi:10.1016/j.ifacol.2017.08.1972.
- [59] U. Aridogan, Y. Shan, K. K. Leang, Design and analysis of discrete-time repetitive control for scanning probe microscopes, *J. Dynamic Systems, Measurement, and Control* 131 (2009) 061103. doi:10.1115/1.4000068.
- [60] S. Necipoglu, S. A. Cebeci, Y. E. Has, L. Guvenc, C. Basdogan, Robust repetitive controller for fast AFM imaging, *IEEE Trans. Nanotechnology* 10 (2011) 1074–1082. doi:10.1109/TNANO.2011.2106797.
- [61] R. J. E. Merry, M. J. C. Ronde, R. van de Molengraft, K. R. Koops, M. Steinbuch, Directional repetitive control of a metro-

- logical AFM, *IEEE Trans. Control Systems Tech.* 19 (2011) 1622–1629. doi:10.1109/TCST.2010.2091642.
- [62] M. Loganathan, D. A. Bristow, Quasi-repetitive control for fast and accurate atomic force microscopy, in: *Proc. American Control Conf.*, 2016, pp. 360–365. doi:10.1109/ACC.2016.7524941.
- [63] H. R. Pota, S. O. R. Moheimani, M. Smith, Resonant controllers for smart structures, *Smart Materials and Structures* 11 (2002) 1–8.
- [64] M. Fairbairn, S. O. R. Moheimani, Resonant control of an atomic force microscope micro-cantilever for active Q control, *Rev. Scientific Instruments* 83 (2012) 083708. doi:10.1063/1.4746277.
- [65] S. K. Das, H. R. Pota, I. R. Petersen, Resonant control of atomic force microscope scanner: A “mixed” negative-imaginary and small-gain approach, in: *Proc. American Control Conf.*, 2013, pp. 5476–5481. doi:10.1109/ACC.2013.6580694.
- [66] N. Nikooienejad, S. O. R. Moheimani, Frequency domain-based integral resonant control design for a MEMS nanopositioner, in: *Proc. IEEE Conf. Control Tech. and Applications*, 2021, pp. 874–879.
- [67] P. Besset, R. Béarée, FIR filter-based online jerk-constrained trajectory generation, *Control Engineering Practice* 66 (2017) 169–180. doi:10.1016/j.conengprac.2017.06.015.