

High-gain observer design for multi-output systems: Transformation to a canonical form by dynamic output shaping

Håvard Fjær Grip¹ and Ali Saberi¹

¹*School of Electrical Engineering and Computer Science, Washington State University, Pullman, WA 99164, USA*

SUMMARY

In this paper we consider the observer design problem for a class of observable linear systems perturbed by nonlinear, time-varying terms. Our design methodology is based on a canonical form, similar to canonical forms used elsewhere in the literature, that allows the nonlinearities to be dominated using high gain. We show that linear state and output transformations to this canonical form exist if, and only if, the data of the system satisfies a certain admissibility property. Moreover, the appropriate transformations can easily be constructed using available tools. We furthermore show that, if a system does not satisfy the admissibility property, it may be possible to extend it with an invertible output filter that makes the data of the extended system admissible. We refer to the problem of constructing such a filter as the *output shaping problem* and introduce an algorithm that solves the problem whenever it is solvable. Copyright © 0000 John Wiley & Sons, Ltd.

Received ...

KEY WORDS: High-gain observers; Estimation; Nonlinear systems

1. INTRODUCTION

In estimation problems it is common to encounter systems that are predominantly described by an observable linear time-invariant part, but that also include nonlinear and time-varying terms. Such systems can be described by

$$\dot{x} = Ax + \phi(t, x), \quad (1a)$$

$$y = Cx, \quad (1b)$$

where $x \in \mathbb{R}^n$ is the state, $y \in \mathbb{R}^p$ is the output, and (A, C) is an observable pair. In some cases the nonlinearities may be exploited to enhance the stability properties of an observer; in other cases, the nonlinearities create an undesirable influence that must be dominated. In the latter case, one often looks for a Lyapunov-type formulation that guarantees stability if the observer gains are chosen in a particular way. A typical result is that stability is ensured if some of the gains are chosen sufficiently high.

It is not always possible to dominate the effect of nonlinearities by increasing gains. The theory of high-gain observers aims to classify the types of systems for which such domination is possible, and to specify how it may be achieved. Typically, one assumes that the nonlinearities are globally Lipschitz continuous (or at least locally Lipschitz continuous within some region of interest) with

Contract/grant sponsor: The work of Håvard Fjær Grip is supported by the Research Council of Norway. The work of Ali Saberi is partially supported by NAVY grants ONR KKK777SB001 and ONR KKK760SB0012.

Copyright © 0000 John Wiley & Sons, Ltd.

Prepared using rncauth.cls [Version: 2010/03/27 v2.00]

respect to the state, uniformly in time. Beyond that, the question of whether domination is possible depends on the structural relationship between the nonlinearities and the outputs.

A situation that does allow for domination is when the system can be written as

$$\dot{x} = Ax + E\psi(t, x), \quad (2a)$$

$$y = Cx, \quad (2b)$$

where the triple (A, C, E) is left-invertible and of minimum phase. The high-gain observer design problem for such a triple is dual to the high-gain feedback design problem, for which much of the early high-gain theory was developed; for an overview, we refer to Saberi and Sannuti [1]. High-gain observers were used early on in the context of *loop transfer recovery* [2], and later for nonlinear systems [3, 4]. A recent paper by Khalil [5] gives an overview of high-gain observers used in nonlinear feedback control.

1.1. High-gain without left-invertibility or minimum phase

The conditions of left-invertibility and minimum phase are sensible when using high gain to suppress an uncertainty about which little or nothing is known. However, these conditions are often too stringent when the uncertainty is due to a nonlinearity whose dependency on the states of the system is known. This was demonstrated by Gauthier, Hammouri, and Othman [6] for single-output systems in the lower-triangular form $\dot{x}_i = x_{i+1} + \phi_i(x_1, \dots, x_i)u$, $i = 1, \dots, n-1$, $\dot{x}_n = f(x_1, \dots, x_n) + \phi_n(x_1, \dots, x_n)u$, $y = x_1$, where u is a known input (and earlier by Williamson [7] for bilinear systems). For such systems the linear part of the system is not described by a left-invertible triple, since the number of independent nonlinearities is greater than the number of outputs.

Generalizing single-output designs like that of Gauthier et al. [6] to multiple-output systems is a complicated matter. Many results have appeared on this topic [8]–[18], based on various canonical forms that generalize the chained structure of the single-output case to multiple chains. Among these, the results by Bornard and Hammouri [8, 13], Hammouri, Bornard, and Busawon [17], and Farza, M'Saad, Triki, and Maatoug [18] allow for the most complex interaction between the chains. However, applying these results requires identifying a set of integers that are difficult to find in practice, as pointed out by Liu, Farza, M'Saad, and Hammouri [16].

Two crucial questions that often receive little attention is when and how a given system can be transformed to a relevant canonical form. In some cases the existence of an appropriate coordinate change can be guaranteed if the system satisfies certain nonlinear observability conditions [11, 12, 15]. However, these conditions are typically hard to confirm and provide little insight regarding how one might construct the coordinate change as a practical matter. A natural approach is to define new coordinates by taking repeated Lie derivatives of the output. In addition to the drawback of often producing highly complicated transformations, this approach is generally not successful when applied to multiple-output systems. This problem is demonstrated by Hou, Busawon, and Saif [19], who propose a procedure that consists of taking repeated Lie derivatives of the output and effectively discarding problematic output components. Such a procedure is likely to waste crucial output information, and it may therefore not succeed even for simple, uniformly observable systems [19, Example 3].

1.2. Topic of this paper

In this paper our focus is on designing an observer for the system (1). Our working assumption is that the nonlinearity does not contribute toward stability, and that it must therefore be dominated by the proper selection of observer gains. Our procedure will be based on transforming the system description to a canonical form that is similar to the canonical forms used in several of the papers cited above. In this context we are interested in *linear* state and output transformations, and we shall show that the existence of such transformations depends in a necessary and sufficient manner on a certain admissibility property of the system data. We shall also demonstrate how the appropriate transformations may be constructed using available tools and software.

Many systems do not satisfy the necessary admissibility property that allows for transformation to the canonical form. However, this situation can often be remedied by first extending the system with an invertible output filter, a process that we refer to as *dynamic output shaping*. The goal of dynamic output shaping is to introduce a larger number of inherent integrations between the outputs and certain subspaces of the state space, so as to make the data of the extended system admissible. We shall present an algorithm that systematically shapes the output in order to achieve this goal.*

2. PROBLEM FORMULATION

We assume throughout the paper that (A, C) is observable and that C is of maximal rank p . Furthermore, we assume that $\phi(t, x)$ is globally Lipschitz continuous, uniformly in t , piecewise continuous in t , and continuously differentiable with respect to x .[†] It follows that the elements of the partial derivative matrix $[\partial\phi/\partial x](t, x)$ are uniformly bounded.

For the purpose of this paper, it is necessary to construct $n \times n$ matrices W_1, \dots, W_v , so that the partial derivative matrix can be written as

$$\frac{\partial\phi}{\partial x}(t, x) = \sum_{k=1}^v \zeta_k(t, x) W_k, \quad (3)$$

where $\zeta_k(t, x)$, $k \in 1, \dots, v$, are scalar basis functions composed as linear combinations of the elements of the partial derivative matrix. The basis functions $\zeta_k(t, x)$ do not need to be known; our entire analysis and design will instead be based on the set of matrices (A, C, W_1, \dots, W_v) associated with the given system, which contains crucial information about the relationship between nonlinearities and outputs.

2.1. Constructing W_1, \dots, W_v

The matrices W_1, \dots, W_v can easily be constructed without even computing the partial derivatives. Specifically, for each (i, j) such that $[\partial\phi_i/\partial x_j](t, x) \neq 0$, we can add to our list of matrices a matrix with the number 1 in element (i, j) and zeros elsewhere. Then (3) clearly holds with each basis function $\zeta_k(t, x)$ representing one nonzero element of $[\partial\phi/\partial x](t, x)$. Although this method of constructing W_1, \dots, W_v is straightforward, it may result in an unnecessarily large number of matrices, due to linear dependencies between the basis functions $\zeta_k(t, x)$. To whatever extent possible, it is advantageous to eliminate linear dependencies so that the number of matrices is minimal. A detailed discussion on this topic is given in Section 5.6.

2.2. Observer and error dynamics

We shall construct an observer for the system (1) on the standard form

$$\dot{\hat{x}} = A\hat{x} + \phi(t, \hat{x}) + K(y - C\hat{x}), \quad (4)$$

where K is a constant gain matrix. Defining the estimation error $\tilde{x} = x - \hat{x}$, this leads to the error dynamics

$$\dot{\tilde{x}} = (A - KC)\tilde{x} + \phi(t, x) - \phi(t, \hat{x}). \quad (5)$$

As shown in Appendix A, the nonlinear term in (5) can be rewritten as

$$\phi(t, x) - \phi(t, \hat{x}) = \sum_{k=1}^v \mu_k(t, x, \hat{x}) W_k \tilde{x}, \quad (6)$$

*A preliminary version of the output shaping algorithm was presented at the 2010 *IEEE Conference and Decision and Control*, but without any proof of its effectiveness [20].

[†]As in other places in the literature (see, e.g., [6]), the global Lipschitz assumption can be relaxed to a local one, within some bounded region of interest, by modifying the nonlinearity.

where $\mu_k(t, x, \hat{x}), k \in 1, \dots, v$, are scalar functions that are uniformly bounded by constants $\mu_{k \max}$. Hence, the error dynamics can be written as

$$\dot{\tilde{x}} = (A - KC)\tilde{x} + \sum_{k=1}^v \mu_k(t, x, \hat{x})W_k\tilde{x}. \quad (7)$$

Our goal is to construct the gain matrix K in such a way as to render the origin of (7) globally exponentially stable.

3. CANONICAL FORM

In this section we introduce a canonical form for the set of matrices (A, C, W_1, \dots, W_v) .

Definition 1 (Canonical form)

The set of matrices (A, C, W_1, \dots, W_v) is said to be in the *canonical form* if

1. the matrices A and C have the form

$$A = \begin{bmatrix} A_{q_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & A_{q_p} \end{bmatrix} + \begin{bmatrix} L_1 \\ \vdots \\ L_p \end{bmatrix} C, \quad C = \begin{bmatrix} C_{q_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & C_{q_p} \end{bmatrix}, \quad (8)$$

where for each $i \in 1, \dots, p$, $A_{q_i} \in \mathbb{R}^{q_i \times q_i}$ and $C_{q_i} \in \mathbb{R}^{1 \times q_i}$ have the special form

$$A_{q_i} = \begin{bmatrix} 0 & I_{q_i-1} \\ 0 & 0 \end{bmatrix}, \quad C_{q_i} = [1 \quad 0 \quad \cdots \quad 0] \quad (9)$$

2. the matrices $W_k, i \in 1, \dots, v$, have the form

$$W_k = \begin{bmatrix} W_{k11} & \cdots & W_{k1p} \\ \vdots & \ddots & \vdots \\ W_{kp1} & \cdots & W_{kpp} \end{bmatrix}, \quad (10)$$

where for each $i, j \in 1, \dots, p$, element (r, c) of $W_{kij} \in \mathbb{R}^{q_i \times q_j}$ is zero if $c > r$. That is, W_{kij} has the lower-triangular structure

$$W_{kij} = \begin{bmatrix} \star & 0 & \cdots & 0 \\ \star & \star & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \quad (11)$$

For a system (1) whose associated set of matrices is in the canonical form, it is instructive to partition the state as $x = [x_1^\top, \dots, x_p^\top]^\top$, where for each $i \in 1, \dots, p$, x_i is of dimension q_i . Similarly, the nonlinearity can be partitioned as $\phi(t, x) = [\phi_1(t, x)^\top, \dots, \phi_p(t, x)^\top]^\top$, where for each $i \in 1, \dots, p$, $\phi_i(t, x)$ is of dimension q_i , and the output can be partitioned as $y = [y_1, \dots, y_p]^\top$, where for each $i \in 1, \dots, p$, y_i is scalar. Then the x_i subsystem with output y_i is described by

$$\begin{aligned} \dot{x}_i &= A_{q_i}x_i + L_i y + \phi_i(t, x), \\ y_i &= C_{q_i}x_i. \end{aligned}$$

Due to the special structure of A_{q_i} and C_{q_i} , the linear part of this subsystem consists of a chain of integrators terminating with the output at the top level, plus a term $L_i y$ that depends only on the output. The nonlinearity $\phi_i(t, x)$ causes additional interaction between the states of the integrator chain, as well as influence from other integrator chains. The structure of this interaction

is governed by the matrices $W_k, k \in 1, \dots, v$, which can be seen by using (6) to write $\phi_i(t, x) = \phi_i(t, 0) + \sum_{k=1}^v \sum_{j=1}^p \mu_k(t, x, 0) W_{kij} x_j$. The crucial thing to note is that W_{kij} is lower-triangular, which means that at each level of integrator chain i , $W_{kij} x_j$ injects a linear function of states that are at the same level or further up in integrator chain j . This structure is visualized in Figure 1(a), which shows three chains of integrators with each arrow representing the influence from one integrator state on another according to the structure dictated by W_k . The arrows point only down or horizontally; arrows pointing up, as illustrated in Figure 1(b), break with the canonical form.

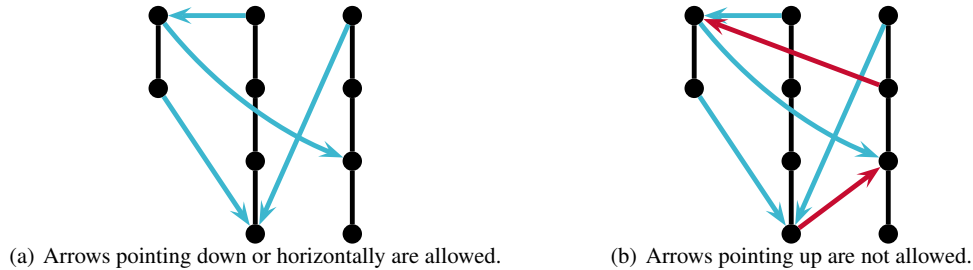


Figure 1. Visualization of canonical form

3.1. Observer design based on canonical form

When (A, C, W_1, \dots, W_v) is in the canonical form, the observer design procedure is straightforward and familiar from the high-gain literature. For each $i = 1, \dots, p$, let $K_i^* = [k_{i1}^*, \dots, k_{iq_i}^*]^T$ be chosen such that the matrix $H_i := A_{q_i} - K_i^* C_{q_i}$ is Hurwitz. Then, define $\tilde{K}_i = [k_{i1}^*/\varepsilon, \dots, k_{iq_i}^*/\varepsilon^{q_i}]^T$, where $\varepsilon \in (0, 1]$ is a high-gain parameter. Finally, define

$$K = \begin{bmatrix} \tilde{K}_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \tilde{K}_p \end{bmatrix} + \begin{bmatrix} L_1 \\ \vdots \\ L_p \end{bmatrix}. \quad (12)$$

Theorem 1

Suppose the gains are chosen as described above. Then there exists an $\varepsilon^* \in (0, 1]$ such that, for all $\varepsilon \in (0, \varepsilon^*]$, the error dynamics (7) is globally exponentially stable.

Proof

If we partition the error state \tilde{x} of (7) as $\tilde{x} = [\tilde{x}_1^T, \dots, \tilde{x}_p^T]^T$, where each $\tilde{x}_i, i = 1, \dots, p$, is of dimension q_i , then the dynamics of \tilde{x}_i is

$$\dot{\tilde{x}}_i = (A_{q_i} - \tilde{K}_i C_{q_i}) \tilde{x}_i + \sum_{k=1}^v \sum_{j=1}^p \mu_k(t, x, \hat{x}) W_{kij} \tilde{x}_j.$$

Define $\xi_i = \Theta_i \tilde{x}_i$, where $\Theta_i = \text{diag}(1, \varepsilon, \dots, \varepsilon^{q_i-1})$. It is easily verified that $\Theta_i (A_{q_i} - \tilde{K}_i C_{q_i}) \Theta_i^{-1} = \frac{1}{\varepsilon} H_i$. Hence, the error dynamics in the new coordinates can be written as

$$\varepsilon \dot{\xi}_i = H_i \xi_i + \varepsilon \sum_{k=1}^v \sum_{j=1}^p \mu_k(t, x, \hat{x}) \Theta_i W_{kij} \Theta_j^{-1} \xi_j.$$

The lower-triangular structure of W_{kij} means that $\Theta_i W_{kij} \Theta_j^{-1}$ has the structure

$$\Theta_i W_{kij} \Theta_j^{-1} = \begin{bmatrix} \star & 0 & 0 & \cdots & 0 \\ \varepsilon \star & \star & 0 & \cdots & 0 \\ \varepsilon^2 \star & \varepsilon \star & \star & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{bmatrix},$$

and thus there exists a bound M_{kij} such that, for all $\varepsilon \in (0, 1]$, $\|\Theta_i W_{kij} \Theta_j^{-1}\| \leq M_{kij}$. For each $i = 1, \dots, p$, let P_i be the unique symmetric positive-definite solution of the Lyapunov equation $P_i H_i + H_i^\top P_i = -I_{q_i}$, and consider the Lyapunov function candidate $V = \varepsilon \sum_{i=1}^p \xi_i^\top P_i \xi_i$. The time derivative of V is

$$\begin{aligned} \dot{V} &= - \sum_{i=1}^p \left(\|\xi_i\|^2 - 2\varepsilon \xi_i^\top P_i \sum_{k=1}^v \sum_{j=1}^p \mu_k(t, x, \hat{x}) \Theta_i W_{kij} \Theta_j^{-1} \xi_j \right) \\ &\leq - \left(1 - 2\varepsilon \sum_{i=1}^p \|P_i\| \sum_{k=1}^v \sum_{j=1}^p \mu_k \max M_{kij} \right) \|\xi\|^2. \end{aligned}$$

By choosing ε sufficiently small, the second term inside the parenthesis can be made smaller than 1. Hence \dot{V} becomes negative definite, and global exponential stability follows. \square

3.2. Transformation to the canonical form

In general, one cannot expect the set of matrices (A, C, W_1, \dots, W_v) to be in the canonical form. Suppose, however, that there exist nonsingular transformations Γ_x and Γ_y such that the set of matrices $(\bar{A}, \bar{C}, \bar{W}_1, \dots, \bar{W}_v)$ is in the canonical form, where $\bar{A} := \Gamma_x^{-1} A \Gamma_x$, $\bar{C} := \Gamma_y^{-1} C \Gamma_x$, and $\bar{W}_k := \Gamma_x^{-1} W_k \Gamma_x$, $k \in 1, \dots, v$. Then we can perform a state transformation $x = \Gamma_x \bar{x}$ and an output transformation $y = \Gamma_y \bar{y}$, which yields the system

$$\begin{aligned} \dot{\bar{x}} &= \bar{A} \bar{x} + \Gamma_x^{-1} \phi(t, \Gamma_x \bar{x}), \\ \bar{y} &= \bar{C} \bar{x}. \end{aligned}$$

Performing the same transformation $\hat{x} = \Gamma_x \bar{x}$ on the observer state yields

$$\dot{\hat{x}} = \bar{A} \hat{x} + \Gamma_x^{-1} \phi(t, \Gamma_x \hat{x}) + \bar{K} (\bar{y} - \bar{C} \hat{x}), \quad (13)$$

where $\bar{K} := \Gamma_x^{-1} K \Gamma_y$. The corresponding error dynamics becomes

$$\begin{aligned} \dot{\tilde{x}} &= (\bar{A} - \bar{K} \bar{C}) \tilde{x} + \Gamma_x^{-1} (\phi(t, \Gamma_x \bar{x}) - \phi(t, \Gamma_x \hat{x})) \\ &= (\bar{A} - \bar{K} \bar{C}) \tilde{x} + \Gamma_x^{-1} \sum_{k=1}^v \mu_k(t, \Gamma_x \bar{x}, \Gamma_x \hat{x}) W_k \Gamma_x \tilde{x} \\ &= (\bar{A} - \bar{K} \bar{C}) \tilde{x} + \sum_{k=1}^v \bar{\mu}_k(t, \bar{x}, \hat{x}) \bar{W}_k \tilde{x}, \end{aligned} \quad (14)$$

where $\bar{\mu}_k(t, \bar{x}, \hat{x}) := \mu_k(t, \Gamma_x \bar{x}, \Gamma_x \hat{x})$. This is the same error dynamics as in (7), but with respect to the set of matrices $(\bar{A}, \bar{C}, \bar{W}_1, \dots, \bar{W}_v)$ in the canonical form; thus, selecting \bar{K} as described in Section 3.1 and implementing the observer (13) with $K = \Gamma_x \bar{K} \Gamma_y^{-1}$ ensures exponential stability of the error dynamics in accordance with Theorem 1.

4. ADMISSIBILITY

We now investigate the existence of nonsingular transformations Γ_x and Γ_y that take a set of matrices (A, C, W_1, \dots, W_v) to the canonical form, and how to construct them if they exist.

We start by noting that it is *always* possible to construct transformations to ensure that $\bar{A} = \Gamma_x^{-1} A \Gamma_x$ and $\bar{C} = \Gamma_y^{-1} C \Gamma_x$ have the form given by (8), (9). For example, if we construct Γ_x and Γ_y to transform the observable triple $(A, 0, C)$ into the *special coordinate basis* (SCB) of Sannuti and Saberi [21], then \bar{A} and \bar{C} will have the desired form. The question is therefore whether we can simultaneously make the matrices $\bar{W}_k = \Gamma_x^{-1} W_k \Gamma_x$, $k \in 1, \dots, v$, satisfy (10), (11). To answer this question, we define a formal property of *admissibility*.

Definition 2 (Admissibility)

A set of matrices (A, C, W_1, \dots, W_v) is said to be *admissible* if for each $i \in 1, \dots, n$ the subspace

$$\mathcal{S}_i(A, C) := \ker \begin{bmatrix} C \\ \vdots \\ CA^{i-1} \end{bmatrix} \quad (15)$$

is W_k -invariant for all $k \in 1, \dots, v$. That is, $W_k \mathcal{S}_i(A, C) \subset \mathcal{S}_i(A, C)$.

Remark 1

The property of admissibility remains unchanged under application of nonsingular transformations Γ_x and Γ_y . To see this, suppose that (A, C, W_1, \dots, W_v) is admissible and note that

$$\mathcal{S}_i(\bar{A}, \bar{C}) = \ker \begin{bmatrix} \bar{C} \\ \vdots \\ \bar{C} \bar{A}^{i-1} \end{bmatrix} = \ker \begin{bmatrix} \Gamma_y^{-1} C \Gamma_x \\ \vdots \\ \Gamma_y^{-1} C (\Gamma_x^{-1} A \Gamma_x)^{i-1} \end{bmatrix} = \ker \begin{bmatrix} C \\ \vdots \\ CA^{i-1} \end{bmatrix} \Gamma_x.$$

Hence, $\bar{x} \in \mathcal{S}_i(\bar{A}, \bar{C}) \Leftrightarrow \Gamma_x \bar{x} \in \ker \mathcal{S}_i(A, C)$. For any $\bar{x} \in \mathcal{S}_i(\bar{A}, \bar{C})$, we therefore have that $W_k \Gamma_x \bar{x} = \Gamma_x \Gamma_x^{-1} W_k \Gamma_x \bar{x} = \Gamma_x \bar{W}_k \bar{x} \in \mathcal{S}_i(A, C) \implies \bar{W}_k \bar{x} \in \mathcal{S}_i(\bar{A}, \bar{C})$.

The next theorem shows that admissibility is necessary and sufficient for the existence of nonsingular transformations to the canonical form. Moreover, it shows that if the set of matrices (A, C, W_1, \dots, W_v) is admissible, then it can be transformed to the canonical form by using any transformations Γ_x and Γ_y that put \bar{A} and \bar{C} in the required form defined by (8), (9).

Theorem 2

Let Γ_x and Γ_y be nonsingular transformations such that \bar{A} and \bar{C} have the form given by (8), (9). Then $(\bar{A}, \bar{C}, \bar{W}_1, \dots, \bar{W}_v)$ is in the canonical form if, and only if, (A, C, W_1, \dots, W_v) is admissible.

Proof

As explained above, there always exist transformations Γ_x and Γ_y that give \bar{A} and \bar{C} the required form defined by (8), (9). Thus, we must show that (i) if, for such a pair of transformations, \bar{W}_k also has the required form defined by (10), (11) for all $k \in 1, \dots, v$, then (A, C, W_1, \dots, W_v) is admissible; and (ii) if (A, C, W_1, \dots, W_v) is admissible then such a transformation always gives \bar{W}_k the required form defined by (10), (11).

We can write $\bar{A} = A_q + L\bar{C}$, where $A_q = \text{diag}(A_{q_1}, \dots, A_{q_p})$, $L = [L_1^T, \dots, L_p^T]^T$, and $\bar{C} = \text{diag}(C_{q_1}, \dots, C_{q_p})$. It follows from Lemma 5 in Appendix B that

$$\mathcal{S}_i(\bar{A}, \bar{C}) = \mathcal{S}_i(A_q, \bar{C}) = \ker \begin{bmatrix} \text{diag}(C_{q_1}, \dots, C_{q_p}) \\ \vdots \\ \text{diag}(C_{q_1} A_{q_1}^{i-1}, \dots, C_{q_p} A_{q_p}^{i-1}) \end{bmatrix}.$$

Let $\bar{x} \in \mathbb{R}^n$ be partitioned as $\bar{x} = [\bar{x}_1^T, \dots, \bar{x}_p^T]^T$, where for each $\rho \in 1, \dots, p$, $\bar{x}_\rho \in \mathbb{R}^{q_\rho}$. Suppose that $\bar{x} \in \mathcal{S}_i(\bar{A}, \bar{C})$ for some i . From the expression for $\mathcal{S}_i(\bar{A}, \bar{C})$ above and the special structure of A_{q_ρ} and C_{q_ρ} , it is easy to see that this is equivalent to the first i components of \bar{x}_ρ being zero (or all components if $q_\rho \leq i$) for all $\rho \in 1, \dots, p$. Defining the vector $\bar{x}_k^* = \bar{W}_k \bar{x}$ and partitioning it in the same way, we have that for each $\rho \in 1, \dots, p$, $\bar{x}_{k\rho}^* = \sum_{j=1}^p \bar{W}_{k\rho j} \bar{x}_j$.

To prove statement (i) above, suppose that \bar{W}_k has the required form for all $k \in 1, \dots, v$, so that $(\bar{A}, \bar{C}, \bar{W}_1, \dots, \bar{W}_v)$ is in the canonical form. Then $\bar{W}_{k\rho j}$ is lower-triangular for all $k \in 1, \dots, v$ and $\rho, j \in 1, \dots, p$, so the first i components of $\bar{x}_{k\rho}^*$ are zero (or all components if $q_\rho \leq i$) for all $\rho \in 1, \dots, p$. Hence $\bar{x}_k^* \in \mathcal{S}_i(\bar{A}, \bar{C})$, and it follows that for each $i \in 1, \dots, n$, $\bar{W}_k \mathcal{S}_i(\bar{A}, \bar{C}) \subset \mathcal{S}_i(\bar{A}, \bar{C})$ for all $k \in 1, \dots, v$. By Remark 1, this implies that $W_k \mathcal{S}_i(A, C) \subset \mathcal{S}_i(A, C)$, so (A, C, W_1, \dots, W_v) is admissible.

To prove statement (ii) above, suppose that (A, C, W_1, \dots, W_v) is admissible, which implies by Remark 1 that $(\bar{A}, \bar{C}, \bar{W}_1, \dots, \bar{W}_v)$ is admissible. Hence $\bar{W}_k \mathcal{S}_i(\bar{A}, \bar{C}) \subset \mathcal{S}_i(\bar{A}, \bar{C})$, which means that for each $k \in 1, \dots, v$ and each $\rho \in 1, \dots, p$, the first i components of $\bar{x}_{k\rho}^*$ must be zero (or all components of $q_\rho \leq i$). Hence, for all $k \in 1, \dots, v$ and all $\rho \in 1, \dots, p$, we must have $\sum_{j=1}^p \bar{W}_{k\rho j} \bar{x}_j = 0$, where $\bar{W}_{k\rho j}$ consists of the upper right-hand $i \times (q_j - i)$ block of $\bar{W}_{k\rho j}$ and \bar{x}_j consists of the last $q_j - i$ elements of \bar{x}_j (i.e., the elements that may be nonzero). Since the vectors \bar{x}_j , $j \in 1, \dots, p$, are arbitrary, this implies that for each $k \in 1, \dots, v$ and each $\rho, j \in 1, \dots, p$, $\bar{W}_{k\rho j} \bar{x}_j = 0$, which in turn implies that $\bar{W}_{k\rho j} = 0$. It follows that for each $k \in 1, \dots, v$ and each $\rho, j \in 1, \dots, p$, the upper right-hand $i \times (q_j - i)$ block of the matrix $\bar{W}_{k\rho j} \in \mathbb{R}^{q_\rho \times q_j}$ must be zero, and this must be true for $i \in 1, \dots, q_\rho$. Hence, $\bar{W}_{k\rho j}$ must have the lower-triangular structure shown in (11), which means that $(\bar{A}, \bar{C}, \bar{W}_1, \dots, \bar{W}_v)$ is in the canonical form. \square

A significant implication of Theorem 2 is that we can simultaneously test for admissibility and find appropriate transformations to the canonical form by simply constructing Γ_x and Γ_y so that $\bar{A} = \Gamma_x^{-1} A \Gamma_x$ and $\bar{C} = \Gamma_y^{-1} C \Gamma_x$ have the required form defined by (8), (9). If the resulting matrices $\bar{W}_k = \Gamma_x^{-1} W_k \Gamma_x$, $k \in 1, \dots, v$, satisfy (10), (11), then $(\bar{A}, \bar{C}, \bar{W}_1, \dots, \bar{W}_v)$ is in the canonical form. If, on the other hand, \bar{W}_k , $k \in 1, \dots, v$, do not satisfy (10), (11), then (A, C, W_1, \dots, W_v) is inadmissible and cannot be transformed to the canonical form by any transformations. As mentioned above, finding transformations so that \bar{A} and \bar{C} have the required form defined by (8), (9) can be done by transforming the triple $(A, 0, C)$ to the SCB. Software is available for accomplishing this task both numerically [22] and symbolically [23].

Example 1

Consider a linear time-varying system with state vector $x = [x_{11}, x_{12}, x_{21}, x_{22}, x_{23}]^T$, given by

$$\begin{aligned} \dot{x}_{11} &= x_{12} + \mu(t)(-x_{12} + x_{22}), & \dot{x}_{21} &= x_{22} + \mu(t)(-x_{12} + x_{22}), \\ \dot{x}_{12} &= \mu(t)x_{23}, & \dot{x}_{22} &= x_{23} + \mu(t)(x_{11} + x_{23}), \\ y_1 &= x_{11}, & \dot{x}_{23} &= \mu(t)x_{12}, \\ & & y_2 &= x_{21}. \end{aligned}$$

We have

$$\phi(t, x) = \begin{bmatrix} \mu(t)(-x_{12} + x_{22}) \\ \mu(t)x_{23} \\ \mu(t)(-x_{12} + x_{22}) \\ \mu(t)(x_{11} + x_{23}) \\ \mu(t)x_{12} \end{bmatrix} \implies \frac{\partial \phi}{\partial x}(t, x) = \mu(t) \begin{bmatrix} 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & -1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}.$$

Hence an associated set of matrices is (A, C, W_1) , where

$$A = \left[\begin{array}{c|ccc} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{array} \right], \quad C = \left[\begin{array}{c|ccc} 1 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 & 0 \end{array} \right], \quad W_1 = \left[\begin{array}{c|ccc} 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ \hline 0 & -1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 \end{array} \right].$$

The matrices A and C are already on the required form defined by (8), (9). However, W_1 is not on the required form (10), (11). Theorem 2 therefore implies that (A, C, W_1) is inadmissible and that no nonsingular transformations can transform it to the canonical form. Indeed, it is easily verified that the condition $W_1 \mathcal{S}_1(A, C) \subset \mathcal{S}_1(A, C)$ fails to hold.

5. DYNAMIC OUTPUT SHAPING

The above analysis shows that, given a system on the form (1) with an associated set of matrices (A, C, W_1, \dots, W_v) , we can construct an observer using high gain if (A, C, W_1, \dots, W_v)

can be transformed to the canonical form. Such a transformation is possible if, and only if, (A, C, W_1, \dots, W_v) is admissible. Even if the admissibility property fails to hold, however, it may be possible to extend the system (1) by adding filters to the outputs, so that the set of matrices associated with the *extended system* becomes admissible. This is demonstrated in the following example by adding integrators to part of the output.

Example 2

Consider the system with state vector $[x_{11}, x_{12}, x_{21}, x_{22}, x_{23}]^T$, given by

$$\begin{aligned} \dot{x}_{11} &= x_{12} + \phi_{11}(x_{23}), & \dot{x}_{21} &= x_{22}, \\ \dot{x}_{12} &= 0, & \dot{x}_{22} &= x_{23}, \\ y_1 &= x_{11}, & \dot{x}_{23} &= 0, \\ & & y_2 &= x_{21}. \end{aligned}$$

This system is used by Hou and Busawon [19, Example 3] as an example of a system that cannot be handled by their proposed method, because no injective map exists to take the system to the relevant canonical form. We have

$$\phi(t, x) = \begin{bmatrix} \phi_{11}(x_{23}) \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \implies \frac{\partial \phi}{\partial x}(t, x) = \frac{\partial \phi_{11}}{\partial x_{23}}(x_{23}) \begin{bmatrix} 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

Hence, a set of matrices associated with the system is given by (A, C, W_1) , where

$$A = \left[\begin{array}{ccc|ccc} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right], \quad C = \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 & 0 & 0 \end{array} \right], \quad W_1 = \left[\begin{array}{ccc|ccc} 0 & 0 & 0 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right].$$

The matrices A and C are already in the required form defined by (8), (9), but W_1 is not in the required form (10), (11). Theorem 2 therefore implies that the set of matrices (A, C, W_1) is inadmissible and that no nonsingular transformation can put it in the canonical form. Suppose, however, that the system is extended by replacing y_1 and y_2 with y_{f1} and y_{f2} , where y_{f1} is a twice-integrated version of y_1 and $y_{f2} = y_2$. Then the extended system equations can be written as

$$\begin{aligned} \dot{z}_1 &= z_2, & \dot{x}_{21} &= x_{22}, \\ \dot{z}_2 &= x_{11}, & \dot{x}_{22} &= x_{23}, \\ \dot{x}_{11} &= x_{12} + \phi(x_{23}), & \dot{x}_{23} &= 0, \\ \dot{x}_{12} &= 0, & y_{f2} &= x_{21}, \\ y_{f1} &= x_{11}, & & \end{aligned}$$

where z_1 and z_2 are the states of the two integrators. It is then easy to see that a set of matrices associated with the extended system is (A_e, C_e, W_{e1}) , given by

$$A_e = \left[\begin{array}{ccc|ccc} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right], \quad C_e = \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 1 & 0 \end{array} \right], \quad W_{e1} = \left[\begin{array}{ccc|ccc} 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right], \quad (16)$$

which does satisfy the canonical form. Hence, an observer can be designed for the extended system.

Motivated by Example 2, we introduce the general idea of dynamically shaping the output y of (1) by adding an invertible output filter. Such a filter is defined by the system equations

$$\begin{aligned}\dot{z} &= A_z z + B_z y, \\ y_f &= C_z z + D_z y,\end{aligned}$$

where $z \in \mathbb{R}^{n_f}$, $u \in \mathbb{R}^p$, and $y_f \in \mathbb{R}^p$. After application of the filter, we can describe the overall system in terms of an extended system vector $x_e = [z^\top, x^\top]^\top$ and write the extended system equations as

$$\dot{x}_e = A_e x_e + \phi_e(t, x_e), \quad (17a)$$

$$y_f = C_e x_e. \quad (17b)$$

It is easily verified that

$$A_e = \begin{bmatrix} A_z & B_z C \\ 0 & A \end{bmatrix}, \quad C_e = [C_z \quad D_z C] \quad (18)$$

and that

$$\phi_e(t, x_e) = \begin{bmatrix} 0 \\ \phi(t, x) \end{bmatrix}.$$

Moreover

$$\frac{\partial \phi_e}{\partial x_e}(t, x_e) = \begin{bmatrix} 0 & 0 \\ 0 & \frac{\partial \phi}{\partial x}(t, x) \end{bmatrix} = \sum_{k=1}^v \zeta_k(t, x) \begin{bmatrix} 0 & 0 \\ 0 & W_k \end{bmatrix}.$$

Hence the set of *extended matrices* associated with the extended system (17) is $(A_e, C_e, W_{e1}, \dots, W_{ev})$, where A_e and C_e are given in (18) and

$$W_{ek} = \begin{bmatrix} 0 & 0 \\ 0 & W_k \end{bmatrix}, \quad k \in 1, \dots, v. \quad (19)$$

We would like to find a filter that makes $(A_e, C_e, W_{e1}, \dots, W_{ev})$ admissible, and we define the following formal problem.

Problem 1 (Output shaping)

Given the set of matrices (A, C, W_1, \dots, W_v) , the *output shaping problem* is to find an invertible filter quadruple (A_z, B_z, C_z, D_z) such that the set of extended matrices $(A_e, C_e, W_{e1}, \dots, W_{ev})$ defined by (18), (19) is admissible and (A_e, C_e) is observable.

Remark 2

Because C is presumed to be of maximal rank p and the filter is invertible, the matrix C_e is also of maximal rank p . Thus, if a solution to the output shaping problem is found, all the conditions for carrying out the observer design in Section 3.1 are satisfied.

Remark 3

The solvability of the output shaping problem remains unchanged under application of nonsingular transformations Γ_x and Γ_y . To see this, suppose that the output shaping problem is solvable for the set of matrices (A, C, W_1, \dots, W_v) by applying the filter quadruple (A_z, B_z, C_z, D_z) to yield the admissible set of extended matrices $(A_e, C_e, W_{e1}, \dots, W_{ev})$. Then the filter quadruple $(A_z, B_z \Gamma_y, C_z, D_z \Gamma_y)$ applied to the transformed set of matrices $(\bar{A}, \bar{C}, \bar{W}_1, \dots, \bar{W}_v)$ yields the extended matrices

$$\begin{aligned}\bar{A}_e &= \begin{bmatrix} A_z & B_z \Gamma_y \bar{C} \\ 0 & \bar{A} \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & \Gamma_x^{-1} \end{bmatrix} \begin{bmatrix} A_z & B_z C \\ 0 & A \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & \Gamma_x \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & \Gamma_x \end{bmatrix}^{-1} A_e \begin{bmatrix} I & 0 \\ 0 & \Gamma_x \end{bmatrix}, \\ \bar{C}_e &= [C_z \quad D_z \Gamma_y \bar{C}] = [C_z \quad D_z C] \begin{bmatrix} I & 0 \\ 0 & \Gamma_x \end{bmatrix} = C_e \begin{bmatrix} I & 0 \\ 0 & \Gamma_x \end{bmatrix},\end{aligned}$$

$$\bar{W}_{ek} = \begin{bmatrix} I & 0 \\ 0 & \Gamma_x^{-1} \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & W_k \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & \Gamma_x \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & \Gamma_x \end{bmatrix}^{-1} W_{ek} \begin{bmatrix} I & 0 \\ 0 & \Gamma_x \end{bmatrix}.$$

Hence, the extended set of matrices $(\bar{A}_e, \bar{C}_e, \bar{W}_{e1}, \dots, \bar{W}_{ev})$ is obtained by a state transformation $\text{diag}(I, \Gamma_x)$ applied to $(A_e, C_e, W_{e1}, \dots, W_{ev})$. By Remark 1 it is therefore admissible.

5.1. Output shaping algorithm

In this section we present an algorithm that solves the output shaping problem whenever it is solvable. We begin by defining a sub-algorithm called **extend**, which takes matrices $\tilde{A} \in \mathbb{R}^{\tilde{n} \times \tilde{n}}$, $\tilde{C} \in \mathbb{R}^{p \times \tilde{n}}$, and $\tilde{W}_k \in \mathbb{R}^{\tilde{n} \times \tilde{n}}$, $k \in 1, \dots, v$, and an integer $m \geq 1$ as parameters. The algorithm returns system matrices $\tilde{A}_z \in \mathbb{R}^{\ell \times \ell}$, $\tilde{B}_z \in \mathbb{R}^{\ell \times p}$, $\tilde{C}_z \in \mathbb{R}^{p \times \ell}$, and $\tilde{D}_z \in \mathbb{R}^{p \times p}$ describing an invertible filter.

$$[\tilde{A}_z, \tilde{B}_z, \tilde{C}_z, \tilde{D}_z] = \text{extend}(\tilde{A}, \tilde{C}, \tilde{W}_1, \dots, \tilde{W}_v, m)$$

Define the matrices

$$R = \begin{bmatrix} \tilde{C} \\ \vdots \\ \tilde{C} \tilde{A}^{m-2} \end{bmatrix}, \quad \tilde{R} = \begin{bmatrix} R \\ \tilde{C} \tilde{A}^{m-1} \end{bmatrix},$$

and let r and \tilde{r} denote their ranks (if $m = 1$, R is an empty matrix and $r = 0$). Let $S \in \mathbb{R}^{pm \times pm}$ be a nonsingular matrix such that

$$S \tilde{R} = \begin{bmatrix} R \\ \tilde{R}^* \\ 0 \end{bmatrix}, \quad S = \begin{bmatrix} I_{(m-1)p} & 0 \\ 0 & S_{22} \\ S_{31} & S_{32} \end{bmatrix},$$

where $\tilde{R}^* \in \mathbb{R}^{(\tilde{r}-r) \times \tilde{n}}$ is of maximal rank $\tilde{r} - r$, $S_{22} \in \mathbb{R}^{(\tilde{r}-r) \times p}$, $S_{31} \in \mathbb{R}^{(p+r-\tilde{r}) \times (m-1)p}$, and $S_{32} \in \mathbb{R}^{(p+r-\tilde{r}) \times p}$. Then $S_{22} \tilde{C} \tilde{A}^{m-1} = \tilde{R}^*$ and $S_{32} \tilde{C} \tilde{A}^{m-1} = -S_{31} R$. Note that the choice of S is in general not unique.

Next, let the columns of E_0 be a linearly independent basis for the $(\tilde{n} - \tilde{r})$ -dimensional kernel of \tilde{R} . For $i = 1, \dots, \sigma$, let the columns of E_i be a linearly independent basis for $\text{im}[E_{i-1}, \tilde{W}_1 E_{i-1}, \dots, \tilde{W}_v E_{i-1}]$, where σ is the smallest integer such that $\rho := \text{rank } E_\sigma = \text{rank } E_{\sigma-1}$. Let $\ell = \text{rank } \tilde{R}^* E_\sigma$.

If $\ell = 0$, define the return matrices as $\tilde{A}_z = 0_{0 \times 0}$, $\tilde{B}_z = 0_{0 \times p}$, $\tilde{C}_z = 0_{p \times 0}$, and $\tilde{D}_z = I_p$. Otherwise, let $T \in \mathbb{R}^{(\tilde{r}-r) \times (\tilde{r}-r)}$ be a nonsingular matrix such that

$$T \tilde{R}^* E_\sigma = \begin{bmatrix} 0 \\ U \end{bmatrix}, \quad T = \begin{bmatrix} T_1 \\ T_2 \end{bmatrix}, \quad (20)$$

where $U \in \mathbb{R}^{\ell \times \rho}$ is of maximal rank ℓ , $T_1 \in \mathbb{R}^{(\tilde{r}-r-\ell) \times (\tilde{r}-r)}$, and $T_2 \in \mathbb{R}^{\ell \times (\tilde{r}-r)}$. Note that the choice of T is in general not unique. Define

$$\mathcal{B} = T_2 S_{22}, \quad \mathcal{D} = \begin{bmatrix} S_{32} \\ T_1 S_{22} \end{bmatrix},$$

and then define the return matrices $\tilde{A}_z \in \mathbb{R}^{\ell \times \ell}$, $\tilde{B}_z \in \mathbb{R}^{\ell \times p}$, $\tilde{C}_z \in \mathbb{R}^{p \times \ell}$, and $\tilde{D}_z \in \mathbb{R}^{p \times p}$ as

$$\tilde{A}_z = 0, \quad \tilde{B}_z = \mathcal{B}, \quad \tilde{C}_z = \begin{bmatrix} I_\ell \\ 0 \end{bmatrix}, \quad \tilde{D}_z = \begin{bmatrix} 0 \\ \mathcal{D} \end{bmatrix}.$$

Remark 4

The quadruple returned by **extend** describes a $p \times p$ filter

$$G(s) = \begin{bmatrix} \frac{1}{s} I_\ell & 0 \\ 0 & I_{p-\ell} \end{bmatrix} \begin{bmatrix} \mathcal{B} \\ \mathcal{D} \end{bmatrix},$$

where

$$\begin{bmatrix} \mathcal{B} \\ \mathcal{D} \end{bmatrix} = \begin{bmatrix} 0 & I_\ell & 0 \\ 0 & 0 & I_{p+r-\bar{r}} \\ I_{\bar{r}-r-\ell} & 0 & 0 \end{bmatrix} \begin{bmatrix} T_1 & 0 \\ T_2 & 0 \\ 0 & I_{p+r-\bar{r}} \end{bmatrix} \begin{bmatrix} S_{22} \\ S_{32} \end{bmatrix}.$$

Since each of the matrices on the right-hand side are invertible, it is clear that the filter is invertible.

We can now describe the complete output shaping algorithm. This algorithm is iterative and maintains a filter quadruple (A_z, B_z, C_z, D_z) that is initialized as an identity filter, and an integer m that is initialized as 1. At each iteration, the extended system matrices are computed according to (18), (19) based on the current filter quadruple, and they are then passed as arguments to **extend** along with the integer m . The return values from **extend** are used to update the current filter quadruple before the next iteration, and to either increment m or reset it to 1.

$$[A_z, B_z, C_z, D_z] = \text{shape}(A, C, W_1, \dots, W_v)$$

1. Initialize the algorithm by defining $A_z = 0_{0 \times 0}$, $B_z = 0_{0 \times p}$, $C_z = 0_{p \times 0}$, $D_z = I_p$, and $m = 1$.
2. Update the extended system matrices as

$$A_e := \begin{bmatrix} A_z & B_z C \\ 0 & A \end{bmatrix}, \quad C_e := [C_z \quad D_z C], \quad W_{ek} := \begin{bmatrix} 0 & 0 \\ 0 & W_k \end{bmatrix}, \quad k \in 1, \dots, v,$$

and let n_e denote the order of the extended system.

3. If $m = n_e + 1$, terminate the algorithm successfully.
4. Execute $[\tilde{A}_z, \tilde{B}_z, \tilde{C}_z, \tilde{D}_z] = \text{extend}(A_e, C_e, W_{e1}, \dots, W_{ev}, m)$ and let ℓ denote the order of the returned matrix quadruple.
5. If $\ell = 0$, increment m by 1. Otherwise reset m to 1.
6. Update the filter quadruple A_z, B_z, C_z , and D_z respectively as

$$\begin{bmatrix} \tilde{A}_z & \tilde{B}_z C_z \\ 0 & A_z \end{bmatrix}, \quad \begin{bmatrix} \tilde{B}_z D_z \\ B_z \end{bmatrix}, \quad [\tilde{C}_z \quad \tilde{D}_z C_z], \quad \tilde{D}_z D_z.$$

7. Repeat from Step 2.

Remark 5

Note that at Step 6, the filter quadruple (A_z, B_z, C_z, D_z) is updated by creating a cascade of the filter *before* the update and the filter $(\tilde{A}_z, \tilde{B}_z, \tilde{C}_z, \tilde{D}_z)$ returned by **extend**.

To get an intuitive understanding of how the output shaping algorithm works, it is helpful to visualize the system as consisting of integrator chains interconnected in a structure dictated by W_k , as illustrated in Figure 1. One can think of the output shaping algorithm as starting at the top level and moving downward, looking for problematic incoming interconnections from lower levels of the integrator chains. The level currently being inspected corresponds to the integer m in the algorithm **shape**. If a problematic interconnection is found, new integrators are added to the output to eliminate the problem at this particular level, which corresponds to the sub-algorithm **extend** returning a filter of order $\ell > 0$. Then the process starts anew from the top level (corresponding to the reset of m to 1) and continues until all levels have been examined without encountering problems (corresponding to the termination condition $m = n_e + 1$).

As an example, consider the structure illustrated in Figure 1(b), which contains two problematic interconnections. The extensions that would take place during output shaping are illustrated in Figure 2. The top level is first examined, and a problematic interconnection is found. To eliminate this problem, an integrator is added to one of the chains. Next, levels 1 and 2 are examined without encountering problematic interconnections, before a new problematic interconnection is found at level 3. Thus, another integrator is added. The process starts again at the top by examining level 1 without finding any problematic interconnections. At level 2, another problematic interconnection is found, so yet another integrator is added to the output. Finally, the process is started over again at

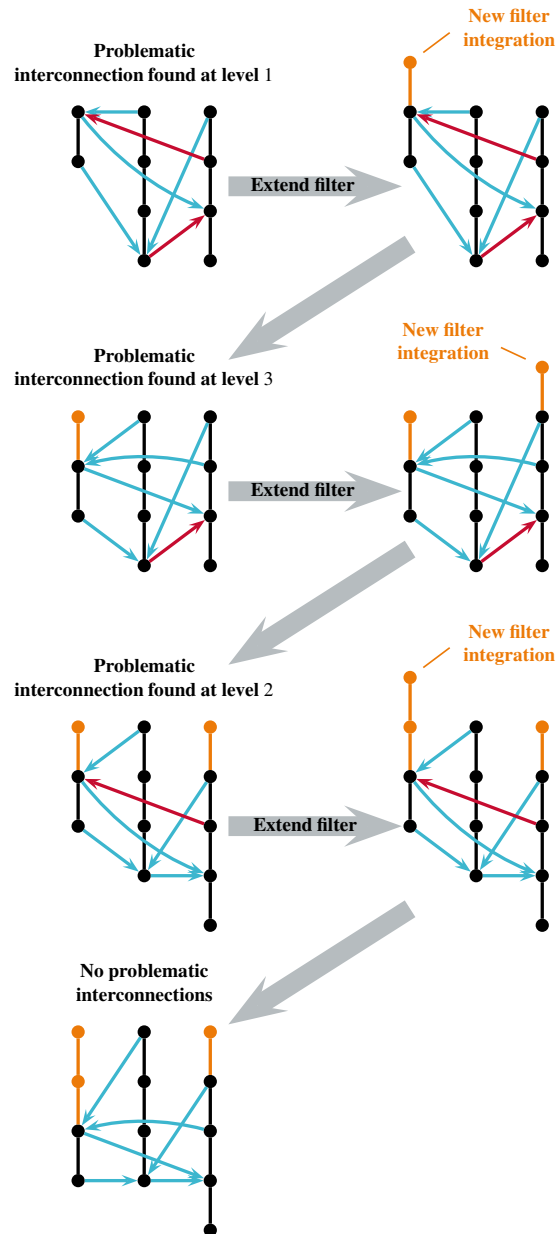


Figure 2. Visualization of output shaping algorithm

the top, and levels 1, . . . , 5 are examined without encountering problematic interconnections. Thus, the algorithm terminates successfully.

Although this visualization is helpful in order to gain an intuitive picture of what happens during the execution of **shape**, it also hides much of the complexity of the process, because the structure in Figure 2 is divided into discrete chains already before the process starts. When dealing with an arbitrary system, we do not know *a priori* how we should transform the linear part into integrator chains, because there are infinitely many ways to do so in general. These are not equivalent for the purpose of carrying out the process illustrated in Figure 2, which accounts for the relative complexity of the algorithm.

5.2. Validity of the output shaping algorithm

Our theoretical result on the validity of the algorithm **shape** is given in the next theorem, which is proven in Section 6.

Theorem 3

Suppose that the output shaping problem for the set of matrices (A, C, W_1, \dots, W_v) is solvable by a filter quadruple of order n_f . Then **shape** (A, C, W_1, \dots, W_v) terminates in finite time and returns a filter quadruple (A_z, B_z, C_z, D_z) that solves the output shaping problem. Moreover, the filter represented by (A_z, B_z, C_z, D_z) is of order less than or equal to n_f .

5.3. Observer implementation

After carrying out the output shaping algorithm, one can construct an observer by first implementing the output filter $\dot{z} = A_z z + B_z y$, $y_f = C_z z + D_z y$, and then implementing an observer for the resulting extended system (17) as

$$\dot{\hat{x}}_e = A_e \hat{x}_e + \phi_e(t, \hat{x}_e) + K(y_f - C_e \hat{x}_e).$$

The resulting error dynamics is then

$$\dot{\tilde{x}}_e = (A_e - K C_e) \tilde{x}_e + \phi_e(t, x_e) - \phi_e(t, \hat{x}_e) \quad (21)$$

The gain matrix K can be chosen to achieve exponential stability of the error dynamics, as described in Section 3.

The same result can also be achieved without the overhead of first implementing the filter states and subsequently estimating them. Consider the observer implementation

$$\dot{\tilde{z}} = (A_z - K_z C_z) \tilde{z} + (B_z - K_z D_z)(y - C \hat{x}), \quad (22a)$$

$$\dot{\hat{x}} = A \hat{x} + \phi(t, \hat{x}) + K_x (C_z \tilde{z} + D_z (y - C \hat{x})). \quad (22b)$$

Defining $\tilde{x}_e = [\tilde{z}^\top, x^\top - \hat{x}^\top]^\top$ and $K = [K_z^\top, K_x^\top]^\top$, we recover precisely the dynamics (21). Hence, by the proper selection of gains, \tilde{x}_e converges exponentially to the origin, and hence $\tilde{x} \rightarrow 0$. Notice that the observer (22) takes the form of a standard observer for x with gain matrix K_x , except that the residual term $y - C \hat{x}$ that would normally be injected has been replaced by the filtered residual $C_z \tilde{z} + D_z (y - C \hat{x})$. Notice also that, because $A_z - K_z C_z$ is always Hurwitz, the implemented filter has no internal instabilities.

5.4. Examples

Although the output shaping algorithm is too complicated to be carried out by hand in most cases, it can be implemented in software. The next two examples illustrate the solutions obtained using *Maple*.

Example 3

Consider again the system in Example 1. As already discussed, the associated set of matrices is inadmissible and can therefore not be transformed to the canonical form through nonsingular transformations of the state and output spaces. Moreover, the system does not satisfy any of the canonical forms from the literature referenced in the introduction. During iteration 1 of **shape**, we have $m = 1$, and the extended system defined in Step 2 is equal to the original system. The sub-algorithm **extend** then returns a filter of order $\ell = 1$, described by

$$\tilde{A}_z = 0, \quad \tilde{B}_z = \begin{bmatrix} 0 & 1 \end{bmatrix}, \quad \tilde{C}_z = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \tilde{D}_z = \begin{bmatrix} 0 & 0 \\ -1 & 1 \end{bmatrix}. \quad (23)$$

The integer m is therefore reset to 1 before iteration 2. During iteration 2, the extended system is the cascade of the original system and the filter described by the quadruple in (23). In this case, **extend**

returns the identity filter of order $\ell = 0$. Hence, m is incremented to 2 before iteration 3. During iteration 3, the extended system stays the same as before. Again **extend** returns the identity filter of order $\ell = 0$, and the same happens for $m = 3, \dots, 6$, so that, before iteration 8, m is incremented to 7. This causes the algorithm to terminate successfully at Step 3 of iteration 8, with the resulting filter being given by the matrices in (23). To select gains for observer implementation, we first transform the extended set of system matrices (A_e, C_e, \bar{W}_{e1}) to the canonical form, using the transformation matrices

$$\Gamma_x = \begin{bmatrix} -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \quad \Gamma_y = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}.$$

This yields the canonical form

$$\bar{A}_e = \left[\begin{array}{ccc|ccc} 0 & 1 & 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right], \quad \bar{C}_e = \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 & 0 & 0 \end{array} \right], \quad \bar{W}_{e1} = \left[\begin{array}{ccc|ccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & -1 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 & 0 \end{array} \right].$$

In the following gain-selection procedure based on Section 3.1, we use a bar above the matrix variables to indicate when we are dealing with the transformed system in the canonical form. We choose $\bar{K}_1^* = \bar{K}_2^* = [0.9, 0.28, 0.04]^T$ to place the poles of \bar{H}_1 and \bar{H}_2 at $-0.2 \pm 0.2j$ and -0.5 . We then compute $\bar{\tilde{K}}_1$ and $\bar{\tilde{K}}_1$ for a given value of $\varepsilon \in (0, 1]$ and assemble the gain matrix \bar{K} according to (12). Finally, we compute the gain matrix $K = \Gamma_x \bar{K} \Gamma_y^{-1}$ in the original coordinate basis, as explained in Section 3.2, and separate it into K_z and K_x for implementation according to Section 5.3. Simulating the system for $\mu(t) = \sin(t)$, we find that the observer error dynamics is unstable for $\varepsilon = 1$. For $\varepsilon = 0.3$, we obtain the gains

$$K_z = [2 \quad 0], \quad K_x \approx \begin{bmatrix} 3.11 & 0 \\ 1.48 & 0 \\ 3.11 & 3.00 \\ 1.48 & 3.11 \\ 0 & 1.48 \end{bmatrix},$$

which results in the a stable error response. Figure 3(a) shows a simulated example of the error response, and Figure 3(b) shows the corresponding trajectory of the internal variable \tilde{z} .

Example 4

Consider again the system in Example 2. During iteration 1, we have $m = 1$, and the extended system defined in Step 2 is the same as the original system. The sub-algorithm **extend** returns a filter of order $\ell = 1$, given by the quadruple

$$\tilde{A}_z = 0, \quad \tilde{B}_z = [1 \quad 0], \quad \tilde{C}_z = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \tilde{D}_z = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}. \quad (24)$$

The integer m is therefore reset to 1 before iteration 2. During iteration 2, the extended system is the cascade of the original system and the filter in (24). In this case **extend** returns the identity filter of order $\ell = 0$, so m is incremented to 2 before iteration 3. During iteration 3, the extended system stays the same as before. The sub-algorithm **extend** now returns a filter order $\ell = 1$, given by exactly the same quadruple as in (24), and hence m is reset to 1. During iteration 4, the extended system is the original system in cascade with two filters given by the quadruple (23). From here on out, **extend** only returns identity filters of order $\ell = 0$ for $m = 1, \dots, 7$, so that the algorithm terminates during iteration 11, when $m = 8$. It can now be confirmed that the result of the output

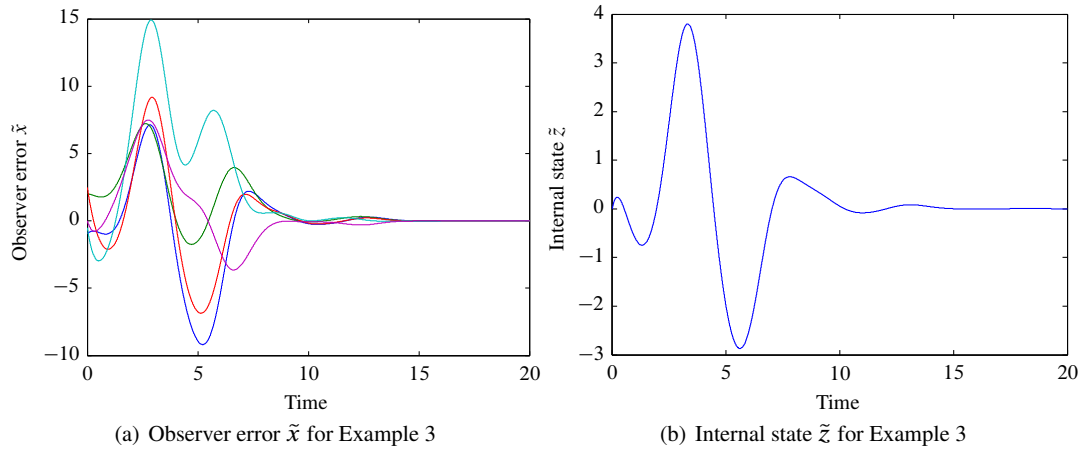


Figure 3. Simulation results for Example 3

shaping algorithm is precisely the extension carried out by hand in Example 2. Hence, the extended set of matrices are in the canonical form, given by (16).

We remark that, although we have chosen examples for which an internal structure is visible to the naked eye, our implementation of the algorithm executes with the same success when these and other systems are transformed to alternative, random coordinate bases, where no such structure is visible.

5.5. Unsuccessful termination

According to Theorem 3, **shape** terminates successfully if the output shaping problem is solvable. What we have not addressed is the case when the output shaping problem is not solvable. In this case the algorithm will continue running *ad infinitum*, building up a larger and larger filter, because it contains no criterion for unsuccessful termination.

One way to create a criterion for unsuccessful termination is to look at the order of the filter maintained by the algorithm. Since, according to Theorem 3, successful termination will always result in a minimal-order filter, we can conclude that no solution exists if the order grows larger than some upper bound on the minimal-order filter. Finding such an upper bound turns out to be highly complicated, however, and we have not been successful in finding a bound that is low enough to be of practical value in most cases. As of now, the decision to terminate the algorithm without a solution should be dictated by practical considerations; that is, the search should be terminated if the computational effort required to continue is too great, or if the order of the filter is too large for practical implementation.

5.6. On the selection of W_1, \dots, W_v

As indicated in Section 2.1, the choice of matrices W_1, \dots, W_v is not unique, because the partial derivative matrix can be decomposed based on different basis functions $\zeta_k(t, x)$ that are linear combinations of the elements of the matrix. We shall show that it is always optimal to choose a linearly independent basis, and that all such bases are equivalent with respect to solvability of the output shaping problem.

Suppose that $\zeta_k(t, x)$, $k \in 1, \dots, v$, is a linearly independent basis, and suppose that $\bar{\zeta}_{\bar{k}}(t, x)$, $\bar{k} \in 1, \dots, \bar{v}$, is an alternative basis. Then we must have

$$\sum_{k=1}^v \zeta_k(t, x) W_k = \sum_{\bar{k}=1}^{\bar{v}} \bar{\zeta}_{\bar{k}}(t, x) \bar{W}_{\bar{k}},$$

for some $\bar{W}_1, \dots, \bar{W}_v$. Linear independence of the basis $\zeta_k(t, x)$, $k \in 1, \dots, v$, implies that for each $k \in 1, \dots, v$, we must have

$$W_k = \sum_{\bar{k}=1}^{\bar{v}} \alpha_{k\bar{k}} \bar{W}_{\bar{k}}$$

for some sets of coefficients $\alpha_{k1}, \dots, \alpha_{k\bar{v}}$. Suppose that the output shaping problem is solvable for $(A, C, \bar{W}_1, \dots, \bar{W}_v)$. Then there exists a filter quadruple that gives rise to a set of extended matrices $(A_e, C_e, \bar{W}_{e1}, \dots, \bar{W}_{e\bar{v}})$, such that $\bar{W}_{e\bar{k}} \mathcal{S}_i(A_e, C_e) \subset \mathcal{S}_i(A_e, C_e)$. Using the same filter for (A, C, W_1, \dots, W_v) gives rise to a set of extended matrices $(A_e, C_e, W_{e1}, \dots, W_{ev})$ where for each $k \in 1, \dots, v$, $W_{ek} = \sum_{\bar{k}=1}^{\bar{v}} \alpha_{k\bar{k}} \bar{W}_{e\bar{k}}$. Hence, we have $W_{ek} \mathcal{S}_i(A_e, C_e) = (\sum_{\bar{k}=1}^{\bar{v}} \alpha_{k\bar{k}} \bar{W}_{e\bar{k}}) \mathcal{S}_i(A_e, C_e) \subset \sum_{\bar{k}=1}^{\bar{v}} (\bar{W}_{e\bar{k}} \mathcal{S}_i) \subset \mathcal{S}_i(A_e, C_e)$, which shows that the output shaping problem is also solvable for (A, C, W_1, \dots, W_v) .

The above analysis shows that if the output shaping problem is solvable for $\bar{W}_1, \dots, \bar{W}_v$ arising from some set of basis functions, then, in particular, it is solvable for all W_1, \dots, W_v arising from a set of linearly independent basis functions. The converse is not true, as can be seen by revisiting Example 1. By using 8 identical basis functions $\zeta_1(t) = \dots = \zeta_8(t) = \mu(t)$, we could split the matrix W_1 into 8 separate matrices, each with only one non-zero element. However, the resulting set of matrices would also be associated with the system

$$\begin{aligned} \dot{x}_{11} &= x_{12} - \mu_1(t)x_{12} + \mu_2(t)x_{22}, & \dot{x}_{21} &= x_{22} - \mu_4(t)x_{12} + \mu_5(t)x_{22}, \\ \dot{x}_{12} &= \mu_3(t)x_{23}, & \dot{x}_{22} &= x_{23} + \mu_6(t)x_{11} + \mu_7(t)x_{23}, \\ y_1 &= x_{11}, & \dot{x}_{23} &= \mu_8(t)x_{12}, \\ & & y_2 &= x_{21}, \end{aligned}$$

which is unobservable for some functions $\mu_1(t), \dots, \mu_8(t)$ (e.g., for the special case $u_1(t) = 1$ and $\mu_2(t) = \dots = \mu_8(t) = 0$). Hence, the output shaping problem cannot be solvable. This shows why methods based on the pattern of the partial derivative matrix without consideration of linear dependence (e.g., [8, 13, 17]), cannot be applied to Example 1.

5.6.1. Effect of state transformation before determining W_1, \dots, W_v According to Remark 3, the solvability of the output shaping problem is not affected by applying nonsingular transformations Γ_x and Γ_y to the set of matrices (A, C, W_1, \dots, W_v) . Suppose, however, that such a transformation were applied to the system (1) before constructing the matrices W_1, \dots, W_v . In this case, we would have a system described by the matrices $\bar{A} = \Gamma_x^{-1} A \Gamma_x$, $\bar{C} = \Gamma_y^{-1} C \Gamma_x$ and the nonlinearity $\bar{\phi}(t, \bar{x}) = \Gamma_x^{-1} \phi(t, \Gamma_x \bar{x})$. It is pertinent to ask whether this might affect the solvability of the output shaping problem. To answer this question, suppose that the output shaping problem is solvable for (A, C, W_1, \dots, W_v) and note that

$$\frac{\partial \bar{\phi}}{\partial \bar{x}}(t, \bar{x}) = \Gamma_x^{-1} \frac{\partial \phi}{\partial x}(t, \Gamma_x \bar{x}) \Gamma_x.$$

If $\zeta_k(t, x)$, $k \in 1, \dots, v$, is a linearly independent basis such that (3) holds, then

$$\frac{\partial \bar{\phi}}{\partial \bar{x}}(t, \bar{x}) = \Gamma_x^{-1} \sum_{k=1}^v \zeta_k(t, \Gamma_x \bar{x}) W_k \Gamma_x = \sum_{k=1}^v \zeta_k(t, \Gamma_x \bar{x}) \Gamma_x^{-1} W_k \Gamma_x = \sum_{k=1}^v \bar{\zeta}_k(t, \bar{x}) \Gamma_x^{-1} W_k \Gamma_x,$$

where $\bar{\zeta}_k(t, \bar{x}) := \zeta_k(t, \Gamma_x \bar{x})$, $k \in 1, \dots, v$. Hence, a set of matrices associated with the transformed system would be $(\bar{A}, \bar{C}, \bar{W}_1, \dots, \bar{W}_v)$, where $\bar{W}_k = \Gamma_x^{-1} W_k \Gamma_x$, $k \in 1, \dots, v$, and it follows from Remark 3 that the output shaping problem is solvable for this set of matrices. By the above analysis, this implies that the output shaping problem is solvable for any linearly independent partitioning of $[\partial \bar{\phi} / \partial \bar{x}](t, \bar{x})$.

6. PROOF OF THEOREM 3

We start by proving that termination of the algorithm implies that the resulting filter quadruple solves the output shaping problem. To this end we need the following lemma.

Lemma 1

For a set of matrices $(\tilde{A}, \tilde{C}, \tilde{W}_1, \dots, \tilde{W}_v)$ and an integer $m \geq 1$, suppose that for each $i \in 1, \dots, m-1$, $\tilde{W}_k \mathfrak{S}_i(\tilde{A}, \tilde{C}) \subset \mathfrak{S}_i(\tilde{A}, \tilde{C})$ for all $k \in 1, \dots, v$. Then, if **extend** $(\tilde{A}, \tilde{C}, \tilde{W}_1, \dots, \tilde{W}_v, m)$ returns a filter quadruple of order $\ell = 0$, then we also have $\tilde{W}_k \mathfrak{S}_m(\tilde{A}, \tilde{C}) \subset \mathfrak{S}_m(\tilde{A}, \tilde{C})$ for all $k \in 1, \dots, v$.

Proof

Let $R, \tilde{R}, \tilde{R}^*$, and E_0, \dots, E_σ refer to the internal values from the execution of **extend**. The subspace $\text{im } E_\sigma$ is the smallest subspace containing $\ker \tilde{R}$ that is \tilde{W}_k -invariant for all $k \in 1, \dots, v$. This can be seen by noting that the definition of E_σ implies $\tilde{W}_k \text{im } E_{\sigma-1} \subset \text{im } E_\sigma$, and that $\text{im } E_\sigma = \text{im } E_{\sigma-1}$. Thus $\tilde{W}_k \text{im } E_\sigma \subset \text{im } E_\sigma$. Since $\text{im } E_0 = \ker \tilde{R}$, it is clear that $\ker \tilde{R} \subset \text{im } E_\sigma$. Finally, if \mathcal{N} is a \tilde{W}_k -invariant subspace for all $k \in 1, \dots, v$ that contains $\ker \tilde{R}$, then we must have $\text{im } E_0 + \sum_{k=1}^v (\tilde{W}_k \text{im } E_0) = \text{im } E_1 \subset \mathcal{N}$, $\text{im } E_1 + \sum_{k=1}^v (\tilde{W}_k \text{im } E_1) = \text{im } E_2 \subset \mathcal{N}$, and so on. Hence $\text{im } E_\sigma \subset \mathcal{N}$.

We have that $\mathfrak{S}_{m-1}(\tilde{A}, \tilde{C}) = \ker R \supset \ker \tilde{R}$ and that $\mathfrak{S}_m(\tilde{A}, \tilde{C}) = \ker \tilde{R}$. From the statement of the theorem, $\mathfrak{S}_{m-1}(\tilde{A}, \tilde{C})$ is \tilde{W}_k -invariant for all $k \in 1, \dots, v$; hence $\text{im } E_\sigma \subset \ker R$, and thus $RE_\sigma = 0$. Using this and $\ell = 0$, we can write

$$\ell = \text{rank } \tilde{R}^* E_\sigma = \text{rank} \begin{bmatrix} R \\ \tilde{R}^* \\ 0 \end{bmatrix} E_\sigma = \text{rank } S \tilde{R} E_\sigma = \text{rank } \tilde{R} E_\sigma = 0.$$

Hence $\text{im } E_\sigma \subset \ker \tilde{R}$. Combined with $\ker \tilde{R} \subset \text{im } E_\sigma$, this implies that $\ker \tilde{R} = \text{im } E_\sigma$. Hence, $\mathfrak{S}_m(\tilde{A}, \tilde{C}) = \ker \tilde{R}$ is \tilde{W}_k -invariant for all $k \in 1, \dots, v$. \square

Using Lemma 1, we can state another lemma.

Lemma 2

After Step 2 of a given iteration, let $(A_e^*, C_e^*, W_{e1}^*, \dots, W_{ev}^*)$ denote the set of extended system matrices, n_e^* the corresponding dimension, and m^* the integer m . Then for each $i \in 1, \dots, m^* - 1$, $W_{ek}^* \mathfrak{S}_i(A_e^*, C_e^*) \subset \mathfrak{S}_i(A_e^*, C_e^*)$ for all $k \in 1, \dots, v$.

Proof

Since m is reset to 1 every time the execution of **extend** results in $\ell \neq 0$, we know that we must have had $\ell = 0$ during the previous $m^* - 1$ iterations. Moreover, since $\ell = 0$ leaves the filter quadruple (A_z, B_z, C_z, D_z) (and therefore the extended system in the following iteration) unchanged, we know that **extend** has been executed on $(A_e^*, C_e^*, W_{e1}^*, \dots, W_{ev}^*)$ for $m = 1, \dots, m^* - 1$, each time resulting in $\ell = 0$. We can now prove the lemma by noting that the premise of Lemma 1 is always satisfied for $m = 1$. Hence, Lemma 1 establishes an induction that implies that for each $i = 1, \dots, m^* - 1$, $W_{ek}^* \mathfrak{S}_i(A_e^*, C_e^*) \subset \mathfrak{S}_i(A_e^*, C_e^*)$ for all $k \in 1, \dots, v$. \square

From Lemma 2, we can now conclude that the termination criteria in Step 3 is satisfied only if the current set of extended matrices $(A_e, C_e, W_{e1}, \dots, W_{ev})$ is admissible. Moreover, (A_e, C_e) is built up by progressively adding invertible output filters returned by **extend** to the original pair (A, C) , and these filters contain no invariant zeros. It can be shown that this does not affect observability, and hence the filter solves the output shaping problem.

6.1. Termination in finite number of steps

Next, we shall prove that, if the output shaping problem is solvable, termination occurs after a finite number of steps. This part of the proof is quite involved; however, the idea is simple. By the premise of the theorem, the output shaping problem is solvable for the original set of matrices (A, C, W_1, \dots, W_v) by a filter of minimal order κ_1 . Letting $\iota \in 1, 2, \dots$, enumerate the iterations,

we shall establish an induction by focusing on the set of extended matrices defined in Step 2 of some arbitrary iteration l^* , denoted by $(A_e^*, C_e^*, W_{e1}^*, \dots, W_{ev}^*)$, and assume that the output shaping problem is solvable for this set of matrices by a filter of minimal order κ_{l^*} . Based on this assumption we shall show that for the set of extended matrices defined in Step 2 of iteration $l^* + 1$, the output shaping problem is solvable by a filter of minimal order κ_{l^*+1} . Moreover, we shall show that $\kappa_{l^*+1} = \kappa_{l^*} - \ell$, where ℓ is the order of the filter returned by **extend** during iteration l^* . Thus, as the algorithm progresses, the order of the filter needed to solve the output shaping problem for the current set of extended matrices is reduced, and we shall show that this implies eventual termination of the algorithm.

In the following, we denote by n_e^* the dimension of the system with associated set of matrices $(A_e^*, C_e^*, W_{e1}^*, \dots, W_{ev}^*)$. We denote by $(A_z^*, B_z^*, C_z^*, D_z^*)$ the κ_{l^*} -dimensional filter quadruple that solves the output shaping problem for the set of matrices $(A_e^*, C_e^*, W_{e1}^*, \dots, W_{ev}^*)$. For the purpose of achieving admissibility, extending the set of matrices $(A_e^*, C_e^*, W_{e1}^*, \dots, W_{ev}^*)$ with the filter quadruple $(A_z^*, B_z^*, C_z^*, D_z^*)$ is equivalent to extending the set of matrices $(A_e^*, \Lambda_u^{-1} C_e^*, W_{e1}^*, \dots, W_{ev}^*)$ with a filter quadruple $(A_z^\dagger, B_z^\dagger, C_z^\dagger, D_z^\dagger)$ of the same order, where $A_z^\dagger = \Lambda_z^{-1} A_z^* \Lambda_z$, $B_z^\dagger = \Lambda_z^{-1} B_z^* \Lambda_u$, $C_z^\dagger = \Lambda_y^{-1} C_z^* \Lambda_z$, and $D_z^\dagger = \Lambda_y^{-1} D_z^* \Lambda_u$ for some nonsingular matrices Λ_z , Λ_y , and Λ_u . To see this, note that the set of extended matrices

$$\left(\begin{bmatrix} A_z^* & B_z^* C_e^* \\ 0 & A_e^* \end{bmatrix}, [C_z^* \quad D_z^* C_e^*], \begin{bmatrix} 0 & 0 \\ 0 & W_{e1}^* \end{bmatrix}, \dots, \begin{bmatrix} 0 & 0 \\ 0 & W_{ev}^* \end{bmatrix} \right),$$

obtained by extending $(A_e^*, C_e^*, W_{e1}^*, \dots, W_{ev}^*)$ with the filter quadruple $(A_z^*, B_z^*, C_z^*, D_z^*)$, is admissible and that

$$\begin{aligned} A_e &:= \begin{bmatrix} A_z^\dagger & B_z^\dagger \Lambda_u^{-1} C_e^* \\ 0 & A_e^* \end{bmatrix} = \begin{bmatrix} \Lambda_z & 0 \\ 0 & I \end{bmatrix}^{-1} \begin{bmatrix} A_z^* & B_z^* C_e^* \\ 0 & A_e^* \end{bmatrix} \begin{bmatrix} \Lambda_z & 0 \\ 0 & I \end{bmatrix}, \\ C_e &:= \begin{bmatrix} C_z^\dagger & D_z^\dagger \Lambda_u^{-1} C_e^* \end{bmatrix} = \Lambda_y^{-1} [C_z^* \quad D_z^* C_e^*] \begin{bmatrix} \Lambda_z & 0 \\ 0 & I \end{bmatrix}, \\ W_{ek} &:= \begin{bmatrix} 0 & 0 \\ 0 & W_{ek}^* \end{bmatrix} = \begin{bmatrix} \Lambda_z & 0 \\ 0 & I \end{bmatrix}^{-1} \begin{bmatrix} 0 & 0 \\ 0 & W_{ek}^* \end{bmatrix} \begin{bmatrix} \Lambda_z & 0 \\ 0 & I \end{bmatrix}, \quad k \in 1, \dots, v. \end{aligned}$$

Hence, it follows from Remark 1 that $(A_e, C_e, W_{e1}, \dots, W_{ev})$, which corresponds to the set of extended matrices obtained by extending $(A_e^*, \Lambda_u^{-1} C_e^*, W_{e1}^*, \dots, W_{ev}^*)$ with the filter quadruple $(A_z^\dagger, B_z^\dagger, C_z^\dagger, D_z^\dagger)$, is admissible.

Let Λ_z , Λ_y , and Λ_u be defined such that $(A_z^\dagger, B_z^\dagger, C_z^\dagger, D_z^\dagger)$ is in the SCB [21, 24]. Then we can write the state of the filter as $z^\dagger = [z_a^\top, z_d^\top]^\top$, where $z_a \in \mathbb{R}^{n_a}$ and $z_d \in \mathbb{R}^{n_d}$ with $n_a + n_d = \kappa_{l^*}$. We can write the input as $u^\dagger = [u_0^\top, u_d^\top]^\top$, where $u_0 \in \mathbb{R}^{p_0}$ and $u_d \in \mathbb{R}^{p_d}$ with $p_0 + p_d = p$. We can write the output as $y_f^\dagger = [y_{f0}^\top, y_{fd}^\top]^\top$ where $y_{f0} \in \mathbb{R}^{p_0}$ and $y_{fd} \in \mathbb{R}^{p_d}$. Furthermore, we can write $z_d = [z_{d1}^\top, \dots, z_{dp_d}^\top]^\top$, where for each $i \in 1, \dots, p_d$, $z_{di} \in \mathbb{R}^{q_i}$ with $q_i \leq q_{i+1}$; $u_d = [u_1, \dots, u_{p_d}]^\top$; and $y_{fd} = [y_{fd1}, \dots, y_{fdp_d}]^\top$. The filter dynamics is then given by

$$\begin{aligned} \begin{bmatrix} \dot{z}_a \\ \dot{z}_d \end{bmatrix} &= \begin{bmatrix} A_a & 0 \\ B_d E_a & A_d + B_d E_d \end{bmatrix} \begin{bmatrix} z_a \\ z_d \end{bmatrix} + \begin{bmatrix} L_{a0} & L_{ad} \\ L_{d0} & L_{dd} \end{bmatrix} \begin{bmatrix} y_{f0} \\ y_{fd} \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & B_d \end{bmatrix} \begin{bmatrix} u_0 \\ u_d \end{bmatrix}, \\ \begin{bmatrix} y_{f0} \\ y_{fd} \end{bmatrix} &= \begin{bmatrix} C_{0a} & C_{0d} \\ 0 & C_d \end{bmatrix} \begin{bmatrix} z_a \\ z_d \end{bmatrix} + \begin{bmatrix} I_{p_0} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} u_0 \\ u_d \end{bmatrix}, \end{aligned}$$

where $A_d = \text{diag}(A_{q_1}, \dots, A_{q_{p_d}})$, $B_d = \text{diag}(B_{q_1}, \dots, B_{q_{p_d}})$, $C_d = \text{diag}(C_{q_1}, \dots, C_{q_{p_d}})$, and where the matrices A_{q_i} , B_{q_i} , and C_{q_i} have the special form

$$A_{q_i} = \begin{bmatrix} 0 & & \\ & I_{q_i-1} & \\ 0 & & 0 \end{bmatrix}, \quad B_{q_i} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}, \quad C_{q_i} = [1 \quad 0 \quad \dots \quad 0].$$

We can assume that columns $1, 1 + q_1, 1 + q_1 + q_2, \dots$ of E_d are zero, since any nonzero elements here can be absorbed in the term $L_{dd}y_{fd}$. We may also assume that the same columns of C_{0d} are zero, since any nonzero columns can be canceled through an output transformation by adding to y_{f0} a multiple of y_{fd} without changing the SCB structure.

From Lemma 5 in Appendix B, we can see that output injection terms can be removed from the filter without affecting the admissibility of $(A_e, C_e, W_{e1}, \dots, W_{ev})$. Moreover, removing output injection terms does not change observability of (A_e, C_e) or the invertibility properties of the filter. Thus, we can modify the filter by setting L_{a0}, L_{ad}, L_{d0} , and L_{dd} to zero. Next, we show that the z_a subsystem is of dimension zero.

Lemma 3

The dimension n_a of z_a is zero.

Proof

We shall show that, if we apply the filter with the state z_a removed, the corresponding set of extended matrices is also admissible. Moreover, the filter remains invertible, since it still has the form of an invertible filter in the SCB, and observability is preserved, since removing the z_a renders the filter zero-free. This will prove the lemma, since the filter is assumed to be of minimal order.

With z_a removed, the extended matrices after applying the filter to $(A_e^*, \Lambda_u^{-1}C_e^*, W_{e1}^*, \dots, W_{ev}^*)$ are

$$\bar{A}_e = \begin{bmatrix} \bar{A}_z & \bar{B}_z \Lambda_u^{-1} C_e^* \\ 0 & A_e^* \end{bmatrix}, \quad \bar{C}_e = [\bar{C}_z \quad \bar{D}_z \Lambda_u^{-1} C_e^*], \quad \bar{W}_{ek} = \begin{bmatrix} 0 & 0 \\ 0 & W_{ek}^* \end{bmatrix},$$

where

$$\bar{A}_z = A_d + B_d E_d, \quad \bar{B}_z = [0_{n_d \times p_0} \quad B_d], \quad \bar{C}_z = \begin{bmatrix} C_{0d} \\ C_d \end{bmatrix}, \quad \bar{D}_z = \begin{bmatrix} I_{p_0} & 0 \\ 0 & 0 \end{bmatrix}.$$

We have $\mathcal{S}_i(\bar{A}_e, \bar{C}_e) = \cap_{j=0}^{i-1} \ker \bar{X}_j$, where $\bar{X}_0 = \bar{C}_e$ and $\bar{X}_j = \bar{X}_{j-1} \bar{A}_e$ for each $j \in 1, \dots, i-1$.

Without z_a removed, the extended matrices can be written as

$$A_e = \begin{bmatrix} A_a & 0 \\ A_{21} & \bar{A}_e \end{bmatrix}, \quad W_{ek} = \begin{bmatrix} 0 & 0 \\ 0 & \bar{W}_{ek} \end{bmatrix}, \quad C_e = [C_1 \quad \bar{C}_e],$$

where

$$A_{21} = \begin{bmatrix} B_d E_a \\ 0_{n_e^* \times n_a} \end{bmatrix}, \quad C_1 = \begin{bmatrix} C_{0a} \\ 0_{p_d \times n_a} \end{bmatrix}.$$

We have $\mathcal{S}_i(A_e, C_e) = \cap_{j=0}^{i-1} \ker X_j$ where $X_0 = C_e$ and $X_j = X_{j-1} A_e$.

Let $X_j = [X_{j1}, X_{j2}]$, where $X_{j1} \in \mathbb{R}^{p \times n_a}$ and $X_{j2} \in \mathbb{R}^{p \times (n_e^* + n_d)}$. For a given $j \in 0, \dots, i-2$, suppose that $X_{j2} = \bar{X}_j$. This holds for $j=0$, because $X_0 = [C_1, \bar{C}_e] = [C_1, \bar{X}_0]$. Then we have $X_{j+1} = [X_{j1}, \bar{X}_j] A_e = [X_{j1} A_a + \bar{X}_j A_{21}, \bar{X}_j \bar{A}_e] = [X_{j1} A_a + \bar{X}_j A_{21}, \bar{X}_{j+1}]$. Hence, $X_{[j+1]2} = \bar{X}_{j+1}$ and by induction, we have $X_{j2} = \bar{X}_j$ for all $j \in 0, \dots, i-1$.

For some $i \in 1, \dots, n_e^* + n_d$ let $\bar{x}_e = [z_d^T, x_e^{*T}]^T \in \mathcal{S}_i(\bar{A}_e, \bar{C}_e)$ be chosen arbitrarily, where $z_d \in \mathbb{R}^{n_d}$ and $x_e^* \in \mathbb{R}^{n_e^*}$. For each $j \in 0, \dots, i-1$, we then have $\bar{X}_j \bar{x}_e = 0$. By the above derivation, this implies

$$X_j \begin{bmatrix} 0_{n_a \times 1} \\ \bar{x}_e \end{bmatrix} = 0 \implies \begin{bmatrix} 0_{n_a \times 1} \\ \bar{x}_e \end{bmatrix} \in \mathcal{S}_i(A_e, C_e).$$

Since $W_{ek} \mathcal{S}_i(A_e, C_e) \subset \mathcal{S}_i(A_e, C_e)$, we have

$$W_{ek} \begin{bmatrix} 0_{n_a \times 1} \\ \bar{x}_e \end{bmatrix} = \begin{bmatrix} 0_{(n_a+n_d) \times 1} \\ W_{ek}^* x_e^* \end{bmatrix} \in \mathcal{S}_i(A_e, C_e) \implies X_j \begin{bmatrix} 0_{(n_a+n_d) \times 1} \\ W_{ek}^* x_e^* \end{bmatrix} = 0, \quad j \in 0, \dots, i-1.$$

Noting that

$$X_j \begin{bmatrix} 0_{(n_a+n_d) \times 1} \\ W_{ek}^* x_e^* \end{bmatrix} = \bar{X}_j \begin{bmatrix} 0_{n_d \times 1} \\ W_{ek}^* x_e^* \end{bmatrix} = \bar{X}_j \bar{W}_{ek} \bar{x}_e,$$

we conclude that $\bar{W}_{ek} \bar{x}_e \in \mathcal{S}_i(\bar{A}_e, \bar{C}_e)$. Hence, admissibility is satisfied with z_a removed. \square

Using Lemma 3, we can write the filter equations as

$$\begin{aligned}\dot{z}_d &= (A_d + B_d E_d)z_d + B_d u_d, \\ y_{f0} &= C_{0d}z_d + u_0, \\ y_{fd} &= C_d z_d.\end{aligned}$$

Let \mathcal{B} , \mathcal{D} , and ℓ refer to the internal values produced during the execution of **extend** in iteration ι^* , where A_e^* , C_e^* , W_{e1}^* , \dots , W_{ev}^* , and an integer $m \geq 1$ are parameters. Define a transformation of the filter input as

$$\begin{bmatrix} \bar{u}_1 \\ \bar{u}_2 \end{bmatrix} = \begin{bmatrix} \mathcal{B} \\ \mathcal{D} \end{bmatrix} \Lambda_u \begin{bmatrix} u_0 \\ u_d \end{bmatrix},$$

where $\bar{u}_1 \in \mathbb{R}^\ell$ and $\bar{u}_2 \in \mathbb{R}^{p-\ell}$.

Lemma 4

We can write

$$\begin{bmatrix} u_0 \\ u_d \end{bmatrix} = \begin{bmatrix} 0 & \bar{C}_{02} \\ \bar{C}_{d1} & \bar{C}_{d2} \end{bmatrix} \begin{bmatrix} \bar{u}_1 \\ \bar{u}_2 \end{bmatrix},$$

where $\bar{C}_{02} \in \mathbb{R}^{p_0 \times (p-\ell)}$ and $\bar{C}_{d1} \in \mathbb{R}^{p_d \times \ell}$ are of rank p_0 and ℓ , respectively.

Proof

Let m refer to the integer parameter to **extend** at iteration ι^* , and let R , \tilde{R} , \tilde{R}^* , r , \tilde{r} , S , E_σ , ρ , T , and U refer to the internal values during the execution of **extend**. We start by showing that $\ker \mathcal{D} \subset \ker D_z^\dagger \Lambda_u^{-1}$. We can write

$$\begin{aligned}\mathfrak{S}_i(A_e, C_e) &= \ker [Y_{i1} \quad Y_{i2}], \\ Y_{i1} &:= \begin{bmatrix} C_z^\dagger \\ \vdots \\ C_z^\dagger A_z^{\dagger i-1} \end{bmatrix}, \quad Y_{i2} := \begin{bmatrix} D_z^\dagger \Lambda_u^{-1} C_e^* \\ \vdots \\ D_z^\dagger \Lambda_u^{-1} C_e^* A_e^{*i-1} + \sum_{j=0}^{i-2} C_z^\dagger A_z^{\dagger i-j-2} B_z^\dagger \Lambda_u^{-1} C_e^* A_e^{*j} \end{bmatrix}.\end{aligned}$$

Let $\mathcal{X} \subset \mathbb{R}^{n_e^*}$ consist of all x_e^* such that $[0_{1 \times n_d}, x_e^{*\top}]^\top \in \mathfrak{S}_m(A_e, C_e)$, which is equivalent to $Y_{m2} x_e^* = 0$. For any $x_e^* \in \mathcal{X}$, we then have

$$W_{ek} \begin{bmatrix} 0_{n_d \times 1} \\ x_e^* \end{bmatrix} = \begin{bmatrix} 0_{n_d \times 1} \\ W_{ek}^* x_e^* \end{bmatrix} \in \mathfrak{S}_m(A_e, C_e) \implies W_{ek}^* x_e^* \in \mathcal{X}.$$

Hence \mathcal{X} is W_{ek}^* -invariant. For any $\chi \in \ker \tilde{R}$, we have $C_e^* A_e^{*i-1} \chi = 0$ for all $i \in 1, \dots, m$. It follows that $Y_{m2} \chi = 0$, which implies $\chi \in \mathcal{X}$. Thus, $\ker \tilde{R} \subset \mathcal{X}$. Following the proof of Lemma 1, we know that $\text{im } E_\sigma$ is the smallest subspace containing $\ker \tilde{R}$ that is W_{ek}^* -invariant for all $k \in 1, \dots, v$. Hence, we must have $\text{im } E_\sigma \subset \mathcal{X}$, which implies $Y_{m2} E_\sigma = 0$.

From Lemma 2, we know that if $m > 1$, then $\mathfrak{S}_{m-1}(A_e^*, C_e^*)$ is W_{ek}^* -invariant for all $k \in 1, \dots, v$. Since $\mathfrak{S}_{m-1}(A_e^*, C_e^*) = \ker R$ and $\ker \tilde{R} \subset \ker R$, we therefore have $\text{im } E_\sigma \subset \ker R$, and hence $RE_\sigma = 0$. This implies that $C_e^* A_e^{*i-1} E_\sigma = 0$ for all $i \in 1, \dots, m-1$. It therefore follows from $Y_{m2} E_\sigma = 0$ that $D_z^\dagger \Lambda_u^{-1} C_e^* A_e^{*m-1} E_\sigma = 0$. This can be rewritten as

$$\begin{aligned}D_z^\dagger \Lambda_u^{-1} C_e^* A_e^{*m-1} E_\sigma &= D_z^\dagger \Lambda_u^{-1} \begin{bmatrix} S_{22} \\ S_{32} \end{bmatrix}^{-1} \begin{bmatrix} S_{22} \\ S_{32} \end{bmatrix} C_e^* A_e^{*m-1} E_\sigma \\ &= D_z^\dagger \Lambda_u^{-1} \begin{bmatrix} S_{22} \\ S_{32} \end{bmatrix}^{-1} \begin{bmatrix} \tilde{R}^* \\ -S_{31} R \end{bmatrix} E_\sigma \\ &= D_z^\dagger \Lambda_u^{-1} \begin{bmatrix} S_{22} \\ S_{32} \end{bmatrix}^{-1} \begin{bmatrix} \tilde{R}^* E_\sigma \\ 0 \end{bmatrix}\end{aligned}$$

$$\begin{aligned}
&= D_z^\dagger \Lambda_u^{-1} \begin{bmatrix} S_{22} \\ S_{32} \end{bmatrix}^{-1} \begin{bmatrix} T & 0 \\ 0 & I_{p+r-\bar{r}} \end{bmatrix}^{-1} \begin{bmatrix} T & 0 \\ 0 & I_{p+r-\bar{r}} \end{bmatrix} \begin{bmatrix} \tilde{R}^* E_\sigma \\ 0 \end{bmatrix} \\
&= D_z^\dagger \Lambda_u^{-1} \begin{bmatrix} S_{22} \\ S_{32} \end{bmatrix}^{-1} \begin{bmatrix} T & 0 \\ 0 & I_{p+r-\bar{r}} \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ U \\ 0 \end{bmatrix} = 0.
\end{aligned}$$

Define

$$Q = \begin{bmatrix} S_{22} \\ S_{32} \end{bmatrix}^{-1} \begin{bmatrix} T & 0 \\ 0 & I_{p+r-\bar{r}} \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ U \\ 0 \end{bmatrix}.$$

Then $D_z^\dagger \Lambda_u^{-1} Q = 0$, so $\text{im } Q \subset \ker D_z^\dagger \Lambda_u^{-1}$. Since $\text{rank } U = \ell$, it follows that $\dim \text{im } Q = \ell$. On the other hand, we can write

$$\mathcal{D} = \begin{bmatrix} 0 & 0 & I_{p+r-\bar{r}} \\ I_{\bar{r}-r-\ell} & 0 & 0 \end{bmatrix} \begin{bmatrix} T_1 & 0 \\ T_2 & 0 \\ 0 & I_{p+r-\bar{r}} \end{bmatrix} \begin{bmatrix} S_{22} \\ S_{32} \end{bmatrix} \implies \mathcal{D}Q = 0.$$

Hence, $\text{im } Q \subset \ker \mathcal{D}$. Since \mathcal{D} has full row rank $p - \ell$, we have $\dim \ker \mathcal{D} = \ell$, and we can therefore conclude that $\text{im } Q = \ker \mathcal{D}$. It now follows that $\ker \mathcal{D} \subset \ker D_z^\dagger \Lambda_u^{-1}$.

Since $\ker \mathcal{D} \subset \ker D_z^\dagger \Lambda_u^{-1} = \ker [I_{p_0}, 0_{p_0 \times p_d}] \Lambda_u^{-1}$, we have $\text{im } \mathcal{D}^\top \supset \text{im}([I_{p_0}, 0_{p_0 \times p_d}] \Lambda_u^{-1})^\top$, and hence there is a matrix \bar{C}_{02} such that $\bar{C}_{02} \mathcal{D} = [I_{p_0}, 0_{p_0 \times p_d}] \Lambda_u^{-1}$. It follows that

$$[0_{p_0 \times \ell} \quad \bar{C}_{02}] \begin{bmatrix} \mathcal{B} \\ \mathcal{D} \end{bmatrix} = [I_{p_0} \quad 0_{p_0 \times p_d}] \Lambda_u^{-1} \implies [0_{p_0 \times \ell} \quad \bar{C}_{02}] = [I_{p_0} \quad 0_{p_0 \times p_d}] \Lambda_u^{-1} \begin{bmatrix} \mathcal{B} \\ \mathcal{D} \end{bmatrix}^{-1}.$$

We therefore have

$$u_0 = [I_{p_0} \quad 0_{p_0 \times p_d}] \begin{bmatrix} u_0 \\ u_d \end{bmatrix} = [I_{p_0} \quad 0_{p_0 \times p_d}] \Lambda_u^{-1} \begin{bmatrix} \mathcal{B} \\ \mathcal{D} \end{bmatrix}^{-1} \begin{bmatrix} \bar{u}_1 \\ \bar{u}_2 \end{bmatrix} = [0_{p_0 \times \ell} \quad \bar{C}_{02}] \begin{bmatrix} \bar{u}_1 \\ \bar{u}_2 \end{bmatrix}.$$

Now writing

$$u_d = [\bar{C}_{d1} \quad \bar{C}_{d2}] \begin{bmatrix} \bar{u}_1 \\ \bar{u}_2 \end{bmatrix},$$

we get the desired expression. Since the transformation from $[u_0^\top, u_d^\top]^\top$ to $[\bar{u}_1^\top, \bar{u}_2^\top]^\top$ is nonsingular, it immediately follows that \bar{C}_{02} must have full row rank p_0 and \bar{C}_{d1} must have full column rank ℓ . \square

We can now write the filter equations as

$$\begin{aligned}
\dot{z}_d &= (A_d + B_d E_d) z_d + B_d (\bar{C}_{d1} \bar{u}_1 + \bar{C}_{d2} \bar{u}_2), \\
y_{f0} &= C_{0d} z_d + \bar{C}_{02} \bar{u}_2, \\
y_{fd} &= C_d z_d.
\end{aligned}$$

We can assume that $\bar{C}_{d2} \in \mathbb{R}^{p_d \times (p-\ell)}$ has rank $p_d - \ell$; otherwise, we can select an L such that $\bar{C}_{d2} + L \bar{C}_{02}$ has rank $p_d - \ell$ and add the output injection term $B_d L y_{f0}$ to the dynamics. The term $B_d L C_{0d} z_d$ thus generated in the dynamics can be absorbed by $B_d E_d z_d$.

Next, perform a state transformation by defining $\tilde{z}_d = (I - B_d \tilde{E}_d) z_d$, where \tilde{E}_d is defined as E_d , but with each column shifted once to the left and with the last column zero, which implies $E_d = \tilde{E}_d A_d$ and $\tilde{E}_d B_d = 0$. Then

$$\begin{aligned}
\dot{\tilde{z}}_d &= (I - B_d \tilde{E}_d) ((A_d + B_d E_d) z_d + B_d (\bar{C}_{d1} \bar{u}_1 + \bar{C}_{d2} \bar{u}_2)) \\
&= (A_d + B_d E_d - B_d \tilde{E}_d A_d - B_d \tilde{E}_d B_d E_d) z_d + B_d (\bar{C}_{d1} \bar{u}_1 + \bar{C}_{d2} \bar{u}_2)
\end{aligned}$$

$$\begin{aligned}
& -B_d \tilde{E}_d B_d (\bar{C}_{d1} \bar{u}_1 + \bar{C}_{d2} \bar{u}_2) \\
& = (A_d + B_d E_d - B_d E_d) z_d + B_d (\bar{C}_{d1} \bar{u}_1 + \bar{C}_{d2} \bar{u}_2) \\
& = A_d z_d + B_d (\bar{C}_{d1} \bar{u}_1 + \bar{C}_{d2} \bar{u}_2).
\end{aligned}$$

We furthermore have

$$A_d z_d = A_d (I - B_d \tilde{E}_d)^{-1} \tilde{z}_d = (A_d + \tilde{B}_d \tilde{E}_d) \tilde{z}_d, \quad (25)$$

where $\tilde{B}_d := A_d B_d$. The last equality in (25) can be confirmed by noting that

$$\begin{aligned}
A_d & = A_d + \tilde{B}_d \tilde{E}_d - \tilde{B}_d \tilde{E}_d - \tilde{B}_d \tilde{E}_d B_d \tilde{E}_d \\
& = A_d + \tilde{B}_d \tilde{E}_d - A_d B_d \tilde{E}_d - \tilde{B}_d \tilde{E}_d B_d \tilde{E}_d \\
& = (A_d + \tilde{B}_d \tilde{E}_d) (I - B_d \tilde{E}_d),
\end{aligned}$$

which implies (by post-multiplication with $(I - B_d \tilde{E}_d)^{-1}$) that $A_d (I - B_d \tilde{E}_d)^{-1} = A_d + \tilde{B}_d \tilde{E}_d$. Hence, the filter can be written as

$$\begin{aligned}
\dot{\tilde{z}}_d & = (A_d + \tilde{B}_d \tilde{E}_d) \tilde{z}_d + B_d (\bar{C}_{d1} \bar{u}_1 + \bar{C}_{d2} \bar{u}_2), \\
y_{f0} & = \tilde{C}_{0d} \tilde{z}_d + \tilde{C}_{02} \bar{u}_2, \\
y_{fd} & = \tilde{C}_d \tilde{z}_d,
\end{aligned}$$

where $\tilde{C}_{0d} = C_{0d} (I - B_d \tilde{E}_d)^{-1}$ and $\tilde{C}_d = C_d (I - B_d \tilde{E}_d)^{-1}$.

Let \tilde{z}_d be partitioned in the same way as z_d . Considering the dynamics of $\tilde{z}_{di} = [\tilde{z}_{di1}, \dots, \tilde{z}_{diq_i}]^T$ for $i \in 1, \dots, p_d$, we have

$$\begin{aligned}
\dot{\tilde{z}}_{dij} & = \tilde{z}_{di[j+1]}, \quad j \in 1, \dots, q_i - 2, \\
\dot{\tilde{z}}_{di[q_i-1]} & = \tilde{z}_{q_i} + \tilde{E}_{di} \tilde{z}_d, \\
\dot{\tilde{z}}_{diq_i} & = \tilde{C}_{d1i} \bar{u}_1 + \tilde{C}_{d2i} \bar{u}_2,
\end{aligned}$$

where \tilde{E}_{di} , \tilde{C}_{d1i} , and \tilde{C}_{d2i} represent the i 'th row of \tilde{E}_d , \tilde{C}_{d1} , and \tilde{C}_{d2} , respectively. Gathering the states \tilde{z}_{diq_i} , $i \in 1, \dots, p_d$, together in a vector \tilde{z}_{dq} , we get

$$\dot{\tilde{z}}_{dq} = [\tilde{C}_{d1} \quad \tilde{C}_{d2}] \begin{bmatrix} \bar{u}_1 \\ \bar{u}_2 \end{bmatrix}.$$

Let $\Gamma = [\Gamma_1^T, \Gamma_2^T]^T$ be a nonsingular matrix, where $\Gamma_1 \in \mathbb{R}^{(p_d - \ell) \times p_d}$ and $\Gamma_2 \in \mathbb{R}^{\ell \times p_d}$, such that

$$\Gamma \tilde{C}_{d2} = \begin{bmatrix} \Gamma_1 \tilde{C}_{d2} \\ 0 \end{bmatrix},$$

which is possible because \tilde{C}_{d2} has rank $p_d - \ell$. Using the state transformation $z'_{dq} = \Gamma \tilde{z}_{dq}$ then yields

$$\dot{z}'_{dq} = \begin{bmatrix} \Gamma_1 \tilde{C}_{d1} & \Gamma_1 \tilde{C}_{d2} \\ \Gamma_2 \tilde{C}_{d1} & 0 \end{bmatrix} \begin{bmatrix} \bar{u}_1 \\ \bar{u}_2 \end{bmatrix},$$

where $\Gamma_2 \tilde{C}_{d1} \in \mathbb{R}^{\ell \times \ell}$ is invertible, which follows from the fact that $[\tilde{C}_{d1}, \tilde{C}_{d2}]$ has full row rank. Furthermore defining

$$\bar{z}_{dq} = \begin{bmatrix} I_{p_d - \ell} & -\Gamma_1 \tilde{C}_{d1} (\Gamma_2 \tilde{C}_{d1})^{-1} \\ 0 & (\Gamma_2 \tilde{C}_{d1})^{-1} \end{bmatrix} z'_{dq}$$

we obtain

$$\dot{\bar{z}}_{dq} = \begin{bmatrix} 0 & \Gamma_1 \tilde{C}_{d2} \\ I_\ell & 0 \end{bmatrix} \begin{bmatrix} \bar{u}_1 \\ \bar{u}_2 \end{bmatrix}.$$

Notice now that, because the dynamics of \bar{z}_{dq} are not affected by any other states of the filter, and because \bar{u}_1 does not affect the filter except in the last ℓ states of \bar{z}_{dq} , we can split the filter into two cascaded filters. Specifically, we can separate out the last ℓ states of \bar{z}_{dq} as \bar{z}_{dq2} to create a separate filter Σ_1 that integrates \bar{u}_1 and feeds \bar{u}_2 directly through:

$$\dot{\bar{z}}_{dq2} = \bar{u}_1, \quad (26a)$$

$$\bar{y}_1 = \bar{z}_{dq2}, \quad (26b)$$

$$\bar{y}_2 = \bar{u}_2. \quad (26c)$$

The overall filter can now be viewed as a cascade of Σ_1 and an $(n_d - \ell)$ -dimensional filter Σ_2 with input $[\bar{y}_1^\top, \bar{y}_2^\top]^\top$. Since Σ_1 and Σ_2 are square and the cascade is invertible, both Σ_1 and Σ_2 are invertible. Notice that we can write Σ_1 as

$$\begin{aligned} \dot{\bar{z}}_{dq2} &= \mathcal{B}\Lambda_u \begin{bmatrix} u_0 \\ u_d \end{bmatrix} = \tilde{A}_z \bar{z}_{dq2} + \tilde{B}_z \Lambda_u \begin{bmatrix} u_0 \\ u_d \end{bmatrix}, \\ \begin{bmatrix} \bar{y}_1 \\ \bar{y}_2 \end{bmatrix} &= \begin{bmatrix} I_\ell \\ 0 \end{bmatrix} \bar{z}_{dq2} + \begin{bmatrix} 0 \\ \mathcal{D} \end{bmatrix} \Lambda_u \begin{bmatrix} u_0 \\ u_d \end{bmatrix} = \tilde{C}_z \bar{z}_{dq2} + \tilde{D}_z \Lambda_u \begin{bmatrix} u_0 \\ u_d \end{bmatrix}, \end{aligned}$$

where \tilde{A}_z , \tilde{B}_z , \tilde{C}_z , and \tilde{D}_z are the return values produced by **extend** during iteration i^* ; that is, Σ_1 is represented by the quadruple $(\tilde{A}_z, \tilde{B}_z \Lambda_u, \tilde{C}_z, \tilde{D}_z \Lambda_u)$. Using Σ_1 to extend $(A_e^*, \Lambda_u^{-1} C_e^*, W_{e1}^*, \dots, W_{ev}^*)$, we obtain

$$\left(\begin{bmatrix} \tilde{A}_z & \tilde{B}_z C_e^* \\ 0 & A_e^* \end{bmatrix}, [\tilde{C}_z \quad \tilde{D}_z C_e^*], \begin{bmatrix} 0 & 0 \\ 0 & W_{e1}^* \end{bmatrix}, \dots, \begin{bmatrix} 0 & 0 \\ 0 & W_{ev}^* \end{bmatrix} \right), \quad (27)$$

and we know that, if we use Σ_2 to further extend (27), we obtain an admissible set of matrices. It is easy to confirm that (27) is precisely the set of extended matrices defined in Step 2 of iteration $i^* + 1$. Hence, Σ_2 represents an invertible filter of order $\kappa_{i^*+1} := n_d - \ell = \kappa_{i^*} - \ell$ that solves the output shaping problem for the set of extended matrices defined in Step 2 of iteration $i^* + 1$. This completes the induction, showing that at each iteration of the algorithm, the output shaping problem for the current extended system is solvable by a filter of minimal order $\kappa_i = \kappa_{i-1} - \ell$, which is always proper.

Since κ_i cannot become non-negative, $\ell > 0$ can occur for a maximum of κ_0 iterations. Hence, there must be a point at which $\ell = 0$ occurs a sufficient number of times in a row that the condition $m = n_e + 1$ in Step 3 is satisfied, thus causing termination. Clearly, the filter returned by the algorithm is of minimum order κ_0 .

7. CONCLUDING REMARKS

In this paper, we have investigated the observer design problem for a class of linear systems perturbed by nonlinear, time-varying terms. We have shown that there exist linear, nonsingular transformations to a canonical form suitable for high-gain observer design if a certain admissibility condition on the system data is satisfied. We have furthermore introduced an algorithm that solves the problem of making the system data admissible through the addition of an invertible output filter, whenever such a solution exists.

The results presented here are not motivated by a particular observability property of the nonlinear system in question, but by the goal of providing a constructive design methodology for a large class of systems. No direct connection to observability is currently known, and the exploration of such a connection is a topic of interest for future work.

A. REWRITING THE NONLINEAR ERROR TERM

To see why (6) holds, note that we can use Taylor's theorem [25, Theorem 11.1] to write

$$\begin{aligned}
\phi(t, x) - \phi(t, \hat{x}) &= \int_0^1 \frac{\partial \phi}{\partial x}(t, \hat{x} + p(x - \hat{x})) dp(x - \hat{x}) \\
&= \int_0^1 \sum_{k=1}^v \zeta_k(t, \hat{x} + p(x - \hat{x})) W_k dp(x - \hat{x}) \\
&= \sum_{k=1}^v \int_0^1 \zeta_k(t, \hat{x} + p(x - \hat{x})) dp W_k(x - \hat{x}) \\
&= \sum_{k=1}^v \mu_k(t, x, \hat{x}) W_k(x - \hat{x}),
\end{aligned} \tag{28}$$

where

$$\mu_k(t, x, \hat{x}) := \int_0^1 \zeta_k(t, \hat{x} + p(x - \hat{x})) dp$$

is uniformly bounded due to the boundedness of $\zeta_k(t, x)$.

B. A USEFUL LEMMA

Lemma 5

For all $i \in 1, \dots, n$, we have that $\mathfrak{S}_i(A + LC, C) = \mathfrak{S}_i(A, C)$, where $A \in \mathbb{R}^{n \times n}$, $L \in \mathbb{R}^{n \times p}$, and $C \in \mathbb{R}^{p \times n}$.

Proof

Suppose that for some $i \in 1, \dots, n-1$, $\mathfrak{S}_i(A + LC, C) = \mathfrak{S}_i(A, C)$. This holds for $i = 1$, since $\mathfrak{S}_1(A + LC, C) = \mathfrak{S}_1(A, C) = \ker C$. Using the fact that for all $j \in 1, \dots, i$, $\ker CA^{j-1} \supset \mathfrak{S}_i(A, C)$, we can write

$$\begin{aligned}
\mathfrak{S}_{i+1}(A + LC, C) &= \mathfrak{S}_i(A, C) \cap \ker C(A + LC)^i \\
&= \mathfrak{S}_i(A, C) \cap \ker C(A + LC)^{i-1} A \\
&= \dots \\
&= \mathfrak{S}_i(A, C) \cap \ker C(A + LC)A^{i-1} \\
&= \mathfrak{S}_i(A, C) \cap \ker CA^i = \mathfrak{S}_{i+1}(A, C).
\end{aligned}$$

Hence, the lemma holds by induction. \square

REFERENCES

1. Saberi A, Sannuti P. Time-scale structure assignment in linear multivariable systems using high-gain feedback. *Int. J. Contr.* 1989; **49**(6):2191–2213.
2. Doyle J, Stein G. Robustness with observers. *IEEE Trans. Automat. Contr.* 1979; **24**(4):607–611.
3. Esfandiari F, Khalil HK. Observer-based design of uncertain systems: Recovering state feedback robustness under matching conditions. *Proc. Allerton Conf.*, Monticello, IL, 1987; 97–106.
4. Saberi A, Sannuti P. Observer design for loop transfer recovery and for uncertain dynamical systems. *IEEE Trans. Automat. Contr.* 1990; **35**(8):878–897.
5. Khalil HK. High-gain observers in nonlinear feedback control. *Proc. Int. Conf. Contr. Autom. Syst.*, Seoul, South Korea, 2008.
6. Gauthier JP, Hammouri H, Othman S. A simple observer for nonlinear systems: Applications to bioreactors. *IEEE Trans. Automat. Contr.* 1992; **37**(6):875–880.

7. Williamson D. Observation of bilinear systems with application to biological control. *Automatica* 1977; **13**(3):243–254.
8. Bornard G, Hammouri H. A high gain observer for a class of uniformly observable systems. *Proc. IEEE Conf. Dec. Contr.*, Brighton, England, 1991; 1494–1496.
9. Deza F, Bossanne D, Busvelle E, Gauthier JP, Rakotopara D. Exponential observers for nonlinear systems. *IEEE Trans. Automat. Contr.* 1993; **38**(3):482–484.
10. Busawon K, Farza M, Hammouri H. Observer design for a special class of nonlinear systems. *Int. J. Contr.* 1998; **71**(3):405–418.
11. Rudolph J, Zeitz M. A block triangular nonlinear observer normal form. *Syst. Contr. Lett.* 1994; **23**(1):1–8.
12. Shim H, Son YI, Seo JH. Semi-global observer for multi-output nonlinear systems. *Syst. Contr. Lett.* 2001; **42**:233–244.
13. Bornard G, Hammouri H. A graph approach to uniform observability of nonlinear multi output systems. *Proc. IEEE Conf. Dec. Contr.*, Las Vegas, NV, 2002; 701–706.
14. Hammouri H, Targui B, Armanet F. High gain observer based on a triangular structure. *Int. J. Robust Nonlin. Contr.* 2002; **12**(6):497–518.
15. Hammouri H, Farza M. Nonlinear observers for locally uniformly observable systems. *ESAIM: Contr. Opt. Calc. Var.* 2003; **9**:353–370.
16. Liu FL, Farza M, M'Saad M, Hammouri H. Observer design for a class of uniformly observable MIMO nonlinear systems with coupled structure. *Proc. IFAC World Congr.*, Seoul, South Korea, 2008; 7630–7635.
17. Hammouri H, Bornard G, Busawon K. High gain observer for structured multi-output nonlinear systems. *IEEE Trans. Automat. Contr.* 2010; **55**(4):987–992.
18. Farza M, M'Saad M, Triki M, Maatoug T. High gain observer for a class of non-triangular systems. *Syst. Contr. Lett.* 2011; **60**(1):27–35.
19. Hou M, Busawon K, Saif M. Observer design based on triangular form generated by injective map. *IEEE Trans. Automat. Contr.* 2000; **45**(7):1350–1355.
20. Grip HF, Saberi A. High-gain observer design for domination of nonlinear perturbations: Transformation to a canonical form by dynamic output shaping. *Proc. IEEE Conf. Dec. Contr.*, Atlanta, GA, 2010; 5043–5049.
21. Sannuti P, Saberi A. Special coordinate basis for multivariable linear systems—Finite and infinite zero structure, squaring down and decoupling. *Int. J. Contr.* 1987; **45**(5):1655–1704.
22. Liu X, Chen BM, Lin Z. Linear systems toolkit in Matlab: Structural decompositions and their applications. *J. Contr. Theor. Appl.* 2005; **3**(3):287–294.
23. Grip HF, Saberi A. Structural decomposition of linear multivariable systems using symbolic computations. *Int. J. Contr.* 2010; **83**(7):1414–1426.
24. Saberi A, Sannuti P. Squaring down of non-strictly proper systems. *Int. J. Contr.* 1990; **51**(3):621–629.
25. Nocedal J, Wright SJ. *Numerical Optimization*. Springer-Verlag: New York, 1999.