

# Inexact Solves in Interpolatory Model Reduction<sup>☆☆ ☆</sup>

Christopher Beattie<sup>a</sup>, Serkan Gugercin<sup>a</sup>, Sarah Wyatt<sup>a</sup>

<sup>a</sup>*Department of Mathematics, Virginia Tech, Blacksburg, VA, 24061-0123*

---

## Abstract

We investigate the use of inexact solves for interpolatory model reduction and consider associated perturbation effects on the underlying model reduction problem. We give bounds on system perturbations induced by inexact solves and relate this to termination criteria for iterative solution methods. We show that when a Petrov-Galerkin framework is employed for the inexact solves, the associated reduced order model is an exact interpolatory model for a nearby full-order system; thus demonstrating backward stability. We also give evidence that for  $\mathcal{H}_2$ -optimal interpolation points, interpolatory model reduction is robust with respect to perturbations due to inexact solves. Finally, we demonstrate the effectiveness of direct use of inexact solves in optimal  $\mathcal{H}_2$  approximation. The result is an effective model reduction strategy that is applicable in realistically large-scale settings.

*Keywords:* Model reduction; system order reduction; tangential interpolation, iterative solves, Petrov-Galerkin

---

## 1. Introduction

The simulation of dynamical systems constitutes a basic framework for the modeling and control of many complex phenomena of interest in science and industry. The need for ever greater model fidelity often leads to computational tasks that make unmanageably large demands on resources. Efficient

---

<sup>☆</sup>Dedicated to Danny Sorensen on the occasion of his 65th birthday.

<sup>☆☆</sup>This work was supported in part by the NSF through Grants DMS-0505971 and DMS-0645347

*Email addresses:* `beattie@vt.edu` (Christopher Beattie), `gugercin@math.vt.edu` (Serkan Gugercin), `sawyatt@vt.edu` (Sarah Wyatt)

model utilization becomes a critical consideration in such large-scale problem settings and motivates the development of strategies for model reduction.

We consider here linear time invariant multi-input/multi-output (MIMO) systems that have a state space form (in the Laplace transform domain) as

$$\text{Find } \hat{\mathbf{v}}(s) \text{ such that } \mathbf{K}(s)\hat{\mathbf{v}}(s) = \mathbf{B}(s)\hat{\mathbf{u}}(s), \quad \text{then } \hat{\mathbf{y}}(s) \stackrel{\text{def}}{=} \mathbf{C}(s)\hat{\mathbf{v}}(s). \quad (1)$$

Here,  $\hat{\mathbf{u}}(s)$  and  $\hat{\mathbf{y}}(s)$  denote Laplace-transformed system inputs and outputs, respectively;  $\hat{\mathbf{v}}(s)$  represents the internal system state. We assume that  $\mathbf{C}(s) \in \mathbb{C}^{p \times n}$  and  $\mathbf{B}(s) \in \mathbb{C}^{n \times m}$  are analytic in the right half plane; and that  $\mathbf{K}(s) \in \mathbb{C}^{n \times n}$  is analytic and full rank throughout the right half plane. Solving for  $\hat{\mathbf{y}}(s)$  in terms of  $\hat{\mathbf{u}}(s)$ , we obtain

$$\hat{\mathbf{y}}(s) = \mathbf{C}(s)\mathbf{K}(s)^{-1}\mathbf{B}(s)\hat{\mathbf{u}}(s) = \mathbf{H}(s)\hat{\mathbf{u}}(s). \quad (2)$$

This representation of the *transfer function*,

$$\mathbf{H}(s) = \mathbf{C}(s)\mathbf{K}(s)^{-1}\mathbf{B}(s), \quad (3)$$

we refer to as a *generalized coprime realization*. Standard first-order descriptor system realizations, with  $\mathbf{H}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}$  for constant matrices  $\mathbf{E}, \mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times m}$ , and  $\mathbf{C} \in \mathbb{R}^{p \times n}$  evidently fit this pattern with  $\mathbf{C}(s) = \mathbf{C}$ ,  $\mathbf{B}(s) = \mathbf{B}$ , and  $\mathbf{K}(s) = s\mathbf{E} - \mathbf{A}$ . However, many dynamical systems can be described more naturally with generalized coprime realizations. For example, a system that includes internal system delays as well as transmission/propagation delays in its input and output could be described with a model

$$\mathbf{E}\dot{\mathbf{x}}(t) = \mathbf{A}_0\mathbf{x}(t) + \mathbf{A}_1\mathbf{x}(t - \tau_{sys}) + \mathbf{B}\mathbf{u}(t - \tau_{inp}), \quad \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t - \tau_{out}) \quad (4)$$

for  $\tau_{sys}, \tau_{inp}, \tau_{out} > 0$ , and  $\mathbf{E}, \mathbf{A}_0, \mathbf{A}_1 \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times m}$  and  $\mathbf{C} \in \mathbb{R}^{p \times n}$ . Taking the Laplace transformation of (4) yields the transfer function

$$\mathbf{H}(s) = \mathbf{C}(s)\mathbf{K}(s)^{-1}\mathbf{B}(s) = (e^{-s\tau_{out}}\mathbf{C})(s\mathbf{E} - \mathbf{A}_0 - e^{-s\tau_{sys}}\mathbf{A}_1)^{-1}(e^{-s\tau_{inp}}\mathbf{B}),$$

which has the form of (3). The form of (3) can accommodate greater generality than this, of course, including memory convolution involving higher derivatives, second and higher-order polynomial differential equations, systems described via integro-differential equations, and systems where state variables may be coupled through infinite dimensional subsystems (possibly

Table 1: Examples of Generalized Coprime System Realizations

Descriptor Systems	$\mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}$ ( $\mathbf{E}$ possibly singular)
Delay Systems	$(e^{-s\tau_{out}}\mathbf{C})(s\mathbf{I} - \mathbf{A}_0 - e^{-s\tau_{sys}}\mathbf{A}_1)^{-1}(e^{-s\tau_{inp}}\mathbf{B})$
Second Order Systems	$(s\mathbf{C}_1 + \mathbf{C}_0)(s^2\mathbf{M} + s\mathbf{G} + \mathbf{K})^{-1}\mathbf{B}$
Weighted Systems	$\mathbf{W}_o(s)\mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B}\mathbf{W}_i(s)$

modeling internal propagation or diffusion). See Table 1 for other examples and [6] for further discussion.

In many applications, the state space dimension,  $n$ , is too large for efficient system simulation and control computation, so the cases of interest for us here have state space dimension vastly larger than input and output dimensions:  $n \gg m, p$ . See [19] for a recent collection of such benchmark problems.

The goal is to produce a reduced system that will have approximately the same response (output) as the original system for any given input  $\mathbf{u}(t)$ . For a given reduced-order  $r \ll n$ , we construct reduced order models through a Petrov-Galerkin approximation of (1): Select full rank matrices  $\mathbf{V}_r \in \mathbb{R}^{n \times r}$  and  $\mathbf{W}_r \in \mathbb{R}^{n \times r}$ . For any input,  $\mathbf{u}(t)$ , the reduced system output,  $\mathbf{y}_r(t)$ , is then defined (in the Laplace transform domain) as:

$$\text{Find } \hat{\mathbf{v}}(s) \in \text{Ran}(\mathbf{V}_r) \text{ such that } \mathbf{W}_r^T(\mathcal{K}(s)\hat{\mathbf{v}}(s) - \mathcal{B}(s)\hat{\mathbf{u}}(s)) = 0 \quad (5)$$

$$\text{then } \hat{\mathbf{y}}_r(s) \stackrel{\text{def}}{=} \mathcal{C}(s)\hat{\mathbf{v}}(s) \quad (6)$$

which defines the reduced transfer function as,

$$\mathcal{H}_r(s) = \mathcal{C}_r(s)\mathcal{K}_r(s)^{-1}\mathcal{B}_r(s), \quad (7)$$

where

$$\begin{aligned} \mathcal{K}_r(s) &= \mathbf{W}_r^T\mathcal{K}(s)\mathbf{V}_r \in \mathbb{C}^{r \times r}, & \mathcal{B}_r(s) &= \mathbf{W}_r^T\mathcal{B}(s) \in \mathbb{C}^{r \times m}, \\ \text{and } \mathcal{C}_r(s) &= \mathcal{C}(s)\mathbf{V}_r \in \mathbb{C}^{p \times r}. \end{aligned} \quad (8)$$

## 2. Interpolatory Model Reduction

Interpolatory reduced order models are designed to exactly reproduce certain system response components that result from inputs having specified frequency content and growth. The approach has been described for standard

first-order system realizations in [13, 2, 11, 3] and extended to generalized coprime realizations in [6]. We summarize the basic elements of this approach below.

A set of points  $\{\mu_i\}_{i=1}^r \subset \mathbb{C}$  and (nontrivial) direction vectors  $\{\mathbf{c}_i\}_{i=1}^r \subset \mathbb{C}^p$  constitute left tangential interpolation data for the reduced model,  $\mathcal{H}_r(s)$ , if

$$\mathbf{c}_i^T \mathcal{H}(\mu_i) = \mathbf{c}_i^T \mathcal{H}_r(\mu_i) \quad \text{for each } i = 1, \dots, r. \quad (9)$$

Likewise,  $\{\sigma_j\}_{j=1}^r$ , and associated directions  $\{\mathbf{b}_j\}_{j=1}^r \subset \mathbb{C}^m$ , constitute right tangential interpolation data for the reduced model,  $\mathcal{H}_r(s)$ , if

$$\mathcal{H}(\sigma_j) \mathbf{b}_j = \mathcal{H}_r(\sigma_j) \mathbf{b}_j \quad \text{for each } j = 1, \dots, r. \quad (10)$$

Given left and right tangential interpolating data, interpolatory model reduction may be implemented by first solving the linear systems:

$$\text{Find } \mathbf{w}_i \text{ such that } \mathbf{w}_i^T \mathcal{K}(\mu_i) = \mathbf{c}_i^T \mathcal{C}(\mu_i) \quad \text{for } i = 1, \dots, r, \text{ and} \quad (11)$$

$$\text{find } \mathbf{v}_i \text{ such that } \mathcal{K}(\sigma_j) \mathbf{v}_j = \mathcal{B}(\sigma_j) \mathbf{b}_j \quad \text{for } j = 1, \dots, r. \quad (12)$$

We assume that the two point sets  $\{\mu_i\}_{i=1}^r$  and  $\{\sigma_j\}_{j=1}^r$  each consist of  $r$  distinct points and that the vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$  and  $\{\mathbf{w}_1, \dots, \mathbf{w}_r\}$  are linearly independent sets. These vectors constitute ‘‘primitive bases’’ for the subspaces  $\mathcal{V}_r = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$  and  $\mathcal{W}_r = \text{span}\{\mathbf{w}_1, \dots, \mathbf{w}_r\}$ . Define the associated matrices:

$$\mathbf{V}_r = [\mathbf{v}_1, \dots, \mathbf{v}_r] = [\mathcal{K}(\sigma_1)^{-1} \mathcal{B}(\sigma_1) \mathbf{b}_1, \dots, \mathcal{K}(\sigma_r)^{-1} \mathcal{B}(\sigma_r) \mathbf{b}_r], \quad (13)$$

$$\mathbf{W}_r^T = \begin{bmatrix} \mathbf{w}_1^T \\ \vdots \\ \mathbf{w}_r^T \end{bmatrix} = \begin{bmatrix} \mathbf{c}_1^T \mathcal{C}(\mu_1) \mathcal{K}(\mu_1)^{-1} \\ \vdots \\ \mathbf{c}_r^T \mathcal{C}(\mu_r) \mathcal{K}(\mu_r)^{-1} \end{bmatrix}. \quad (14)$$

The reduced model,  $\mathcal{H}_r(s)$ , as defined in (7) and (8) using  $\mathbf{V}_r$  and  $\mathbf{W}_r$  from (13) and (14), interpolates  $\mathcal{H}(s)$  at the  $2r$  points  $\{\mu_i\}_{i=1}^r$  and  $\{\sigma_j\}_{j=1}^r$ , in respective output directions  $\{\mathbf{c}_i\}_{i=1}^r$  and input directions  $\{\mathbf{b}_j\}_{j=1}^r$ ; that is, conditions (9) and (10) are satisfied. If  $\mu_k = \sigma_k$  for some  $k$  then first order bitangential moments match as well:

$$\mathbf{c}_k^T \mathcal{H}'(\mu_k) \mathbf{b}_k = \mathbf{c}_k^T \mathcal{H}'_r(\mu_k) \mathbf{b}_k$$

Interpolation of higher order derivatives of  $\mathcal{H}(s)$  can be accomplished with similar constructions as well; see [6, 3] and references therein.

For large-scale settings with millions of degrees of freedom, interpolatory model reduction has become the method of choice since it does not require dense matrix operations; the major computational cost lies in solving the (often sparse) linear systems in (11) and (12). This contrasts with Gramian-based model reduction approaches such as balanced truncation [25, 24], optimal Hankel norm approximation [12] and singular perturbation approximation [21] where large-scale Lyapunov equations need to be solved. Moreover, these computational advantages have been enhanced for standard first order state-space realizations by strategies for optimal selection of tangential interpolation data, see [16].

### 2.1. Inexact Interpolatory Model Reduction

The basic framework for interpolatory model reduction presumes that the key equations (11) and (12) may be solved exactly or nearly so, at least to an accuracy associated with machine precision. Direct solution methods, employing sparse factorization strategies, for example, are capable of handling systems of significantly large order. However since the need for ever greater modeling detail and fidelity can drive system order to the order of millions, the use of direct solvers for the linear systems (11) and (12) often becomes infeasible and iterative methods must be employed that terminate with possibly coarse approximate solutions to the linear systems. We consider and evaluate issues related to these approaches here.

Suppose  $\{\widehat{\mathbf{v}}_1, \dots, \widehat{\mathbf{v}}_r\}$  and  $\{\widehat{\mathbf{w}}_1, \dots, \widehat{\mathbf{w}}_r\}$  are linearly independent sets in  $\mathbb{C}^n$  and define

$$\widehat{\mathbf{V}}_r = [\widehat{\mathbf{v}}_1, \dots, \widehat{\mathbf{v}}_r] \quad \widehat{\mathbf{W}}_r^T = \begin{bmatrix} \widehat{\mathbf{w}}_1^T \\ \vdots \\ \widehat{\mathbf{w}}_r^T \end{bmatrix}. \quad (15)$$

$\widehat{\mathbf{w}}_i$  and  $\widehat{\mathbf{v}}_j$  will be viewed as approximate solutions to the linear systems (11) and (12) and accordingly we will refer to them as “inexact” solutions to (11) and (12). Nonetheless, unless otherwise stated, these vectors can be any arbitrarily chosen linearly independent vectors in  $\mathbb{C}^n$ .

Define residuals,  $\boldsymbol{\xi}_i$  and  $\boldsymbol{\eta}_j$ , corresponding to  $\widehat{\mathbf{w}}_i$  and  $\widehat{\mathbf{v}}_j$ , as

$$\boldsymbol{\xi}_i = \mathcal{K}(\mu_i)^T \widehat{\mathbf{w}}_i - \mathcal{C}(\mu_i)^T \mathbf{c}_i \quad \text{and} \quad \boldsymbol{\eta}_j = \mathcal{K}(\sigma_j) \widehat{\mathbf{v}}_j - \mathcal{B}(\sigma_j) \mathbf{b}_j \quad (16)$$

The deviations from the corresponding exact solutions are then

$$\delta \mathbf{w}_i = \widehat{\mathbf{w}}_i - \mathbf{w}_i = \mathcal{K}(\mu_i)^{-T} \boldsymbol{\xi}_i \quad \text{and} \quad \delta \mathbf{v}_j = \widehat{\mathbf{v}}_j - \mathbf{v}_j = \mathcal{K}(\sigma_j)^{-1} \boldsymbol{\eta}_j. \quad (17)$$

The resulting (inexact) basis matrices destined for use in a reduced order model are

$$\widehat{\mathbf{W}}_r = \mathbf{W}_r + [\delta\mathbf{w}_1, \dots, \delta\mathbf{w}_r] \quad (18)$$

$$\widehat{\mathbf{V}}_r = \mathbf{V}_r + [\delta\mathbf{v}_1, \dots, \delta\mathbf{v}_r]. \quad (19)$$

Define reduced order maps associated with these inexact bases:

$$\widehat{\mathcal{K}}_r(s) = \widehat{\mathbf{W}}_r^T \mathcal{K}(s) \widehat{\mathbf{V}}_r, \quad \widehat{\mathcal{B}}_r(s) = \widehat{\mathbf{W}}_r^T \mathcal{B}(s), \quad \text{and} \quad \widehat{\mathcal{C}}_r(s) = \mathcal{C}(s) \widehat{\mathbf{V}}_r, \quad (20)$$

together with the associated inexact reduced order transfer function

$$\widehat{\mathcal{H}}_r(s) = \widehat{\mathcal{C}}_r(s) \widehat{\mathcal{K}}_r(s)^{-1} \widehat{\mathcal{B}}_r(s).$$

Notice that we are free to make any choice for bases for the subspaces,  $\widehat{\mathcal{V}}_r$  and  $\widehat{\mathcal{W}}_r$ , in defining  $\widehat{\mathcal{H}}_r(s)$ ; no change in the definition of (20) is necessary. As a practical matter, it is generally prudent to choose well conditioned bases in computation.

### 3. Forward Error

#### 3.1. Interpolation Error

Inexactness in the solution of the key linear systems (11) and (12) produces a computed reduced order transfer function,  $\widehat{\mathcal{H}}_r(s)$  that no longer interpolates  $\mathcal{H}(s)$ ; typically, the reduced order system response will no longer match any component of the full order system response at any of the complex frequencies  $\{\mu_i\}_{i=1}^r$  and  $\{\sigma_i\}_{i=1}^r$  that have been specified. How much response error has been introduced at these points ?

The particular realization taken for a transfer function can create innate sensitivities to perturbations associated with that representation. Define perturbed transfer functions,

$$\mathcal{H}_{\delta\mathcal{B}}(s) = \mathcal{C}(s) \mathcal{K}(s)^{-1} (\mathcal{B}(s) + \delta\mathcal{B}) \quad \text{and} \quad \mathcal{H}_{\delta\mathcal{C}}(s) = (\mathcal{C}(s) + \delta\mathcal{C}) \mathcal{K}(s)^{-1} \mathcal{B}(s).$$

In discussing perturbations in system response caused by  $\delta\mathcal{B}$  and  $\delta\mathcal{C}$  at  $s = \sigma$ , it is natural to introduce the following quantities:

$$\text{cond}_{\mathcal{B}}(\mathcal{H}(\sigma)) = \frac{\|\mathcal{C}(\sigma) \mathcal{K}(\sigma)^{-1}\| \|\mathcal{B}(\sigma)\|}{\|\mathcal{H}(\sigma)\|}$$

$$\text{cond}_{\mathcal{C}}(\mathcal{H}(\sigma)) = \frac{\|\mathcal{C}(\sigma)\| \|\mathcal{K}(\sigma)^{-1} \mathcal{B}(\sigma)\|}{\|\mathcal{H}(\sigma)\|}$$

to be *condition numbers of the transfer function response*, by way of analogy to the condition number of algebraic linear systems. (Unless otherwise noted, norms will always refer to the Euclidean 2-norm for vectors or the naturally induced spectral norm for matrices). It is straightforward to show that these quantities measure the relative sensitivity of the system with respect to perturbations in  $\mathbf{B}$  and  $\mathbf{C}$ , respectively:

$$\begin{aligned} \frac{\|\mathcal{H}_{\delta\mathbf{B}}(\sigma) - \mathcal{H}(\sigma)\|}{\|\mathcal{H}(\sigma)\|} &\leq \text{cond}_{\mathbf{B}}(\mathcal{H}(\sigma)) \frac{\|\delta\mathbf{B}\|}{\|\mathbf{B}(\sigma)\|} \quad \text{and} \\ \frac{\|\mathcal{H}_{\delta\mathbf{C}}(\sigma) - \mathcal{H}(\sigma)\|}{\|\mathcal{H}(\sigma)\|} &\leq \text{cond}_{\mathbf{C}}(\mathcal{H}(\sigma)) \frac{\|\delta\mathbf{C}\|}{\|\mathbf{C}(\sigma)\|}. \end{aligned}$$

For values of  $s$  such that  $\mathcal{K}_r(s)$  and  $\widehat{\mathcal{K}}_r(s)$  are nonsingular, define the matrix-valued functions,

$$\begin{aligned} \mathcal{P}_r(s) &= \mathcal{K}(s)\mathbf{V}_r\mathcal{K}_r(s)^{-1}\mathbf{W}_r^T, \quad \mathcal{Q}_r(s) = \mathbf{V}_r\mathcal{K}_r(s)^{-1}\mathbf{W}_r^T\mathcal{K}(s), \\ \widehat{\mathcal{P}}_r(s) &= \mathcal{K}(s)\widehat{\mathbf{V}}_r\widehat{\mathcal{K}}_r(s)^{-1}\widehat{\mathbf{W}}_r^T, \quad \text{and} \quad \widehat{\mathcal{Q}}_r(s) = \widehat{\mathbf{V}}_r\widehat{\mathcal{K}}_r(s)^{-1}\widehat{\mathbf{W}}_r^T\mathcal{K}(s) \end{aligned} \quad (21)$$

Where defined,  $\mathcal{P}_r(s)$ ,  $\mathcal{Q}_r(s)$ ,  $\widehat{\mathcal{P}}_r(s)$ , and  $\widehat{\mathcal{Q}}_r(s)$  are differentiable (indeed, analytic) with respect to  $s$ , having derivatives that satisfy:

$$\begin{aligned} \widehat{\mathcal{P}}_r'(s) &= (\mathbf{I} - \widehat{\mathcal{P}}_r) \mathcal{K}'(s)\mathcal{K}(s)^{-1}\widehat{\mathcal{P}}_r \\ \text{and} \quad \widehat{\mathcal{Q}}_r'(s) &= \widehat{\mathcal{Q}}_r\mathcal{K}(s)^{-1}\mathcal{K}'(s)(\mathbf{I} - \widehat{\mathcal{Q}}_r) \end{aligned} \quad (22)$$

with equivalent expressions for  $\mathcal{P}_r'(s)$  and  $\mathcal{Q}_r'(s)$ . We will make a series of observations about properties of  $\widehat{\mathcal{P}}_r(s)$  and  $\widehat{\mathcal{Q}}_r(s)$  which will have immediately apparent parallels to properties for  $\mathcal{P}_r(s)$  and  $\mathcal{Q}_r(s)$ .

Observe first that  $\widehat{\mathcal{P}}_r^2 = \widehat{\mathcal{P}}_r$  and  $\widehat{\mathcal{Q}}_r^2 = \widehat{\mathcal{Q}}_r$  so both  $\widehat{\mathcal{P}}_r(s)$  and  $\widehat{\mathcal{Q}}_r(s)$  are skew projectors. These projectors are of interest because the pointwise error in the transfer function can be expressed as

$$\begin{aligned} \mathcal{H}(s) - \widehat{\mathcal{H}}_r(s) &= \mathcal{C}(s) \left( \mathcal{K}(s)^{-1} - \widehat{\mathbf{V}}_r\widehat{\mathcal{K}}_r(s)^{-1}\widehat{\mathbf{W}}_r^T \right) \mathcal{B}(s) \\ &= \mathcal{C}(s)\mathcal{K}(s)^{-1} (\mathbf{I} - \widehat{\mathcal{P}}_r(s)) \mathcal{B}(s). \end{aligned}$$

Similarly,

$$\mathcal{H}(s) - \widehat{\mathcal{H}}_r(s) = \mathcal{C}(s) (\mathbf{I} - \widehat{\mathcal{Q}}_r(s)) \mathcal{K}(s)^{-1}\mathcal{B}(s)$$

and

$$\mathcal{H}(s) - \widehat{\mathcal{H}}_r(s) = \mathbf{C}(s) \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(s) \right) \mathcal{K}(s)^{-1} \left( \mathbf{I} - \widehat{\mathcal{P}}_r(s) \right) \mathcal{B}(s).$$

The derivative of this last expression can be computed with the aid of (22) and observing  $\mathcal{K}(s)^{-1} \widehat{\mathcal{P}}_r(s) = \widehat{\mathcal{Q}}_r(s) \mathcal{K}(s)^{-1}$ :

$$\begin{aligned} \mathcal{H}'(s) - \widehat{\mathcal{H}}_r'(s) &= \frac{d}{ds} [\mathbf{C}(s) \mathcal{K}(s)^{-1}] \left( \mathbf{I} - \widehat{\mathcal{P}}_r(s) \right) \mathcal{B}(s) \\ &\quad + \mathbf{C}(s) \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(s) \right) \frac{d}{ds} [\mathcal{K}(s)^{-1} \mathcal{B}(s)] \\ &\quad - \mathbf{C}(s) \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(s) \right) \frac{d}{ds} [\mathcal{K}(s)^{-1}] \left( \mathbf{I} - \widehat{\mathcal{P}}_r(s) \right) \mathcal{B}(s). \end{aligned} \quad (23)$$

We introduce the following ( $s$ -dependent) subspaces:

$$\begin{aligned} \mathfrak{P}_r(s) &= \text{Ran } \mathcal{P}_r(s) = \text{Ran } \mathcal{K}(s) \mathbf{V}_r, & \mathcal{Q}_r(s) &= \text{Ker } (\mathbf{W}_r^T \mathcal{K}(s))^\perp, \\ \widehat{\mathfrak{P}}_r(s) &= \text{Ran } \widehat{\mathcal{P}}_r(s) = \text{Ran } \mathcal{K}(s) \widehat{\mathbf{V}}_r, & \widehat{\mathcal{Q}}_r(s) &= \text{Ker } (\widehat{\mathbf{W}}_r^T \mathcal{K}(s))^\perp, \\ \mathfrak{B}_m(s) &= \text{Ran } \mathcal{K}(s)^{-1} \mathcal{B}(s), & \mathfrak{C}_p(s) &= \text{Ker } (\mathbf{C}(s) \mathcal{K}(s)^{-1})^\perp. \end{aligned}$$

$\widehat{\mathcal{P}}_r(s)$  maps vectors in  $\mathbb{C}^n$  onto  $\widehat{\mathfrak{P}}_r(s)$  along  $\widehat{\mathcal{W}}_r^\perp$  and  $\widehat{\mathcal{Q}}_r$  maps vectors in  $\mathbb{C}^n$  onto  $\widehat{\mathcal{V}}_r$  along  $\widehat{\mathcal{Q}}_r(s)^\perp$ .

Given two subspaces of  $\mathbb{C}^n$ , say  $\mathcal{M}$  and  $\mathcal{N}$ , we express the proximity of one to the other in terms of the angle between the subspaces,  $\Theta(\mathcal{M}, \mathcal{N}) \in [0, \frac{\pi}{2}]$  defined as

$$\sup_{\mathbf{x} \in \mathcal{M}} \inf_{\mathbf{y} \in \mathcal{N}} \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}\|} = \sin \Theta(\mathcal{M}, \mathcal{N}).$$

$\Theta(\mathcal{M}, \mathcal{N})$  is the largest canonical angle between  $\mathcal{M}$  and a ‘‘closest’’ subspace  $\widehat{\mathcal{N}}$  of  $\mathcal{N}$  having dimension equal to  $\dim \mathcal{M}$ . Notice that if  $\dim \mathcal{N} < \dim \mathcal{M}$  then  $\Theta(\mathcal{M}, \mathcal{N}) = \frac{\pi}{2}$  and  $\Theta(\mathcal{M}, \mathcal{N}) = 0$  if and only if  $\mathcal{M} \subset \mathcal{N}$ .  $\Theta(\mathcal{M}, \mathcal{N})$  is asymmetrically defined with respect to  $\mathcal{M}$  and  $\mathcal{N}$ , however if  $\dim \mathcal{N} = \dim \mathcal{M}$  then  $\Theta(\mathcal{M}, \mathcal{N}) = \Theta(\mathcal{N}, \mathcal{M})$ . If  $\mathbf{\Pi}_{\mathcal{M}}$  and  $\mathbf{\Pi}_{\mathcal{N}}$  denote orthogonal projectors onto  $\mathcal{M}$  and  $\mathcal{N}$ , respectively, then  $\sin \Theta(\mathcal{M}, \mathcal{N}) = \|(\mathbf{I} - \mathbf{\Pi}_{\mathcal{M}}) \mathbf{\Pi}_{\mathcal{N}}\|$ .

The spectral norm of a skew projector can be expressed in terms of the angle between its range and cokernel [27]. In particular,

$$\|\widehat{\mathcal{P}}_r(s)\| = \|\mathbf{I} - \widehat{\mathcal{P}}_r(s)\| = \frac{1}{\cos \Theta(\widehat{\mathfrak{P}}_r(s), \widehat{\mathcal{W}}_r)} \quad (24)$$

$$\|\widehat{\mathcal{Q}}_r(s)\| = \|\mathbf{I} - \widehat{\mathcal{Q}}_r(s)\| = \frac{1}{\cos \Theta(\widehat{\mathcal{Q}}_r(s), \widehat{\mathcal{V}}_r)} \quad (25)$$

**Theorem 3.1.** *Given the full-order model  $\mathcal{H}(s) = \mathcal{C}(s)\mathcal{K}(s)^{-1}\mathcal{B}(s)$ , interpolation points  $\{\sigma_j\} \subset \mathbb{C}$ ,  $\{\mu_i\} \subset \mathbb{C}$  and corresponding tangential directions,  $\{\mathbf{b}_j\} \subset \mathbb{C}^m$  and  $\{\mathbf{c}_i\} \subset \mathbb{C}^p$ , let the inexact interpolatory reduced model  $\widehat{\mathcal{H}}_r(s) = \widehat{\mathcal{C}}_r(s)\widehat{\mathcal{K}}_r(s)^{-1}\widehat{\mathcal{B}}_r(s)$  be constructed as defined in (15)-(20). The (tangential) interpolation error at  $\mu_i$  and  $\sigma_j$  is*

$$\frac{\|\widehat{\mathcal{H}}_r(\sigma_j)\mathbf{b}_j - \mathcal{H}(\sigma_j)\mathbf{b}_j\|}{\|\mathcal{H}(\sigma_j)\mathbf{b}_j\|} \leq \text{cond}_{\mathcal{B}}(\mathcal{H}(\sigma_j)\mathbf{b}_j) \frac{\sin \Theta(\mathcal{C}_p(\sigma_j), \widehat{\mathcal{W}}_r)}{\cos \Theta(\widehat{\mathcal{P}}_r(\sigma_j), \widehat{\mathcal{W}}_r)} \frac{\|\boldsymbol{\eta}_j\|}{\|\mathcal{B}(\sigma_j)\mathbf{b}_j\|} \quad (26)$$

$$\frac{\|\mathbf{c}_i^T \widehat{\mathcal{H}}_r(\mu_i) - \mathbf{c}_i^T \mathcal{H}(\mu_i)\|}{\|\mathbf{c}_i^T \mathcal{H}(\mu_i)\|} \leq \text{cond}_{\mathcal{C}}(\mathbf{c}_i^T \mathcal{H}(\mu_i)) \frac{\sin \Theta(\mathcal{B}_m(\mu_i), \widehat{\mathcal{V}}_r)}{\cos \Theta(\widehat{\mathcal{Q}}_r(\mu_i), \widehat{\mathcal{V}}_r)} \frac{\|\boldsymbol{\xi}_i\|}{\|\mathbf{c}_i^T \mathcal{C}(\mu_i)\|}. \quad (27)$$

If  $\mu_i = \sigma_i$  then,

$$|\mathbf{c}_i^T \widehat{\mathcal{H}}_r(\mu_i)\mathbf{b}_i - \mathbf{c}_i^T \mathcal{H}(\mu_i)\mathbf{b}_i| \leq \frac{\|\mathcal{K}(\mu_i)^{-1}\| \|\boldsymbol{\eta}_i\| \|\boldsymbol{\xi}_i\|}{\max\left(\cos \Theta(\widehat{\mathcal{P}}_r(\mu_i), \widehat{\mathcal{W}}_r), \cos \Theta(\widehat{\mathcal{Q}}_r(\mu_i), \widehat{\mathcal{V}}_r)\right)}. \quad (28)$$

and

$$|\mathbf{c}_i^T \mathcal{H}'(\mu_i)\mathbf{b}_i - \mathbf{c}_i^T \widehat{\mathcal{H}}_r'(\mu_i)\mathbf{b}_i| \leq M \left( \frac{\|\boldsymbol{\eta}_i\|}{\cos \Theta(\widehat{\mathcal{P}}_r(\mu_i), \widehat{\mathcal{W}}_r)} + \frac{\|\boldsymbol{\xi}_i\|}{\cos \Theta(\widehat{\mathcal{Q}}_r(\mu_i), \widehat{\mathcal{V}}_r)} + \frac{\|\boldsymbol{\eta}_i\|}{\cos \Theta(\widehat{\mathcal{P}}_r(\mu_i), \widehat{\mathcal{W}}_r)} \frac{\|\boldsymbol{\xi}_i\|}{\cos \Theta(\widehat{\mathcal{Q}}_r(\mu_i), \widehat{\mathcal{V}}_r)} \right) \quad (29)$$

with  $M = \max\left(\left\|\frac{d}{ds} [\mathbf{c}_i^T \mathcal{C} \mathcal{K}^{-1}]|_{\mu_i}\right\|, \left\|\frac{d}{ds} [\mathcal{K}^{-1} \mathcal{B} \mathbf{b}_i]|_{\mu_i}\right\|, \left\|\frac{d}{ds} [\mathcal{K}^{-1}]|_{\mu_i}\right\|\right)$ .

PROOF: From (17),  $\widehat{\mathbf{v}}_j = \mathcal{K}(\sigma_j)^{-1}(\mathcal{B}(\sigma_j)\mathbf{b}_j + \boldsymbol{\eta}_j)$ , which implies then that  $\mathcal{K}(\sigma_j)\widehat{\mathbf{v}}_j = \mathcal{B}(\sigma_j)\mathbf{b}_j + \boldsymbol{\eta}_j \in \widehat{\mathcal{P}}_r(\sigma_j)$  and  $(\mathbf{I} - \widehat{\mathcal{P}}_r(\sigma_j))(\mathcal{B}(\sigma_j)\mathbf{b}_j + \boldsymbol{\eta}_j) = 0$ , which may be rearranged to obtain

$$(\mathbf{I} - \widehat{\mathcal{P}}_r(\sigma_j)) \mathcal{B}(\sigma_j)\mathbf{b}_j = -(\mathbf{I} - \widehat{\mathcal{P}}_r(\sigma_j)) \boldsymbol{\eta}_j. \quad (30)$$

Let  $\widehat{\boldsymbol{\Pi}}$  be the orthogonal projector taking  $\mathbb{C}^n$  onto  $\widehat{\mathcal{W}}_r = \text{Ker}(\widehat{\mathcal{P}}_r(s))^\perp$ . One

may directly verify that  $\mathbf{I} - \widehat{\mathcal{P}}_r(s) = (\mathbf{I} - \widehat{\Pi}) (\mathbf{I} - \widehat{\mathcal{P}}_r(s))$ , and

$$\begin{aligned} \widehat{\mathcal{H}}_r(\sigma_j)\mathbf{b}_j - \mathcal{H}(\sigma_j)\mathbf{b}_j &= -\mathbf{C}(\sigma_j)\mathcal{K}(\sigma_j)^{-1} (\mathbf{I} - \widehat{\mathcal{P}}_r(\sigma_j)) \mathcal{B}(\sigma_j)\mathbf{b}_j \\ &= \mathbf{C}(\sigma_j)\mathcal{K}(\sigma_j)^{-1} (\mathbf{I} - \widehat{\mathcal{P}}_r(\sigma_j)) \boldsymbol{\eta}_j \\ &= \mathbf{C}(\sigma_j)\mathcal{K}(\sigma_j)^{-1} (\mathbf{I} - \widehat{\Pi}) (\mathbf{I} - \widehat{\mathcal{P}}_r(\sigma_j)) \boldsymbol{\eta}_j. \end{aligned} \quad (31)$$

Now suppose  $\Gamma$  is an orthogonal projector onto  $\mathfrak{C}_p(\sigma_j)$ . We have then that  $\text{Ran}(\mathbf{I} - \Gamma) = \text{Ker}(\mathbf{C}(\sigma_j)\mathcal{K}(\sigma_j)^{-1})$ , so that  $\mathbf{C}(\sigma_j)\mathcal{K}(\sigma_j)^{-1} = \mathbf{C}(\sigma_j)\mathcal{K}(\sigma_j)^{-1}\Gamma$  and

$$\widehat{\mathcal{H}}_r(\sigma_j)\mathbf{b}_j - \mathcal{H}(\sigma_j)\mathbf{b}_j = \mathbf{C}(\sigma_j)\mathcal{K}(\sigma_j)^{-1}\Gamma (\mathbf{I} - \widehat{\Pi}) (\mathbf{I} - \widehat{\mathcal{P}}_r(\sigma_j)) \boldsymbol{\eta}_j.$$

Taking norms, we obtain an estimate yielding (26):

$$\begin{aligned} \|\widehat{\mathcal{H}}_r(\sigma_j)\mathbf{b}_j - \mathcal{H}(\sigma_j)\mathbf{b}_j\| &\leq \|(\mathbf{I} - \widehat{\Pi}) \Gamma (\mathbf{C}(\sigma_j)\mathcal{K}(\sigma_j)^{-1})^T\| \cdot \|\mathbf{I} - \widehat{\mathcal{P}}_r(\sigma_j)\| \cdot \|\boldsymbol{\eta}_j\| \\ &\leq \|\mathbf{C}(\sigma_j)\mathcal{K}(\sigma_j)^{-1}\| \cdot \frac{\sin \Theta(\mathfrak{C}_p(\sigma_j), \widehat{\mathcal{W}}_r)}{\cos \Theta(\widehat{\mathcal{P}}_r(\sigma_j), \widehat{\mathcal{W}}_r)} \cdot \|\boldsymbol{\eta}_j\| \end{aligned}$$

(27) is shown similarly, noting first that

$$\mathbf{c}_i^T \mathbf{C}(\mu_i) (\mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i)) = -\boldsymbol{\xi}_i^T (\mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i)). \quad (32)$$

Defining the orthogonal projector,  $\widehat{\Xi}$ , that takes  $\mathbb{C}^n$  onto  $\widehat{\mathcal{V}}_r = \text{Ran}(\widehat{\mathcal{Q}}_r(s))$ , one observes next  $\mathbf{I} - \widehat{\mathcal{Q}}_r(s) = (\mathbf{I} - \widehat{\mathcal{Q}}_r(s)) (\mathbf{I} - \widehat{\Xi})$  so that

$$\begin{aligned} \|\mathbf{c}_i^T \widehat{\mathcal{H}}_r(\mu_i) - \mathbf{c}_i^T \mathcal{H}(\mu_i)\| &= \|\mathbf{c}_i^T \mathbf{C}(\mu_i) (\mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i)) \mathcal{K}(\mu_i)^{-1} \mathcal{B}(\mu_i)\| \\ &\leq \|\boldsymbol{\xi}_i^T (\mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i)) (\mathbf{I} - \widehat{\Xi}) \mathcal{K}(\mu_i)^{-1} \mathcal{B}(\mu_i)\| \\ &\leq \|\boldsymbol{\xi}_i\| \cdot \|\mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i)\| \cdot \|(\mathbf{I} - \widehat{\Xi}) \mathcal{K}(\mu_i)^{-1} \mathcal{B}(\mu_i)\| \\ &\leq \|\mathcal{K}(\mu_i)^{-1} \mathcal{B}(\mu_i)\| \cdot \frac{\sin \Theta(\mathfrak{B}_m(\mu_i), \widehat{\mathcal{V}}_r)}{\cos \Theta(\widehat{\mathcal{Q}}_r(\mu_i), \widehat{\mathcal{V}}_r)} \cdot \|\boldsymbol{\xi}_i\| \end{aligned}$$

When  $\mu_i = \sigma_i$ , we have

$$\begin{aligned} \mathbf{c}_i^T \mathcal{H}(\mu_i) \mathbf{b}_i - \mathbf{c}_i^T \widehat{\mathcal{H}}_r(\mu_i) \mathbf{b}_i &= \mathbf{c}_i^T \mathbf{e}(\mu_i) \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i) \right) \mathcal{K}(\mu_i)^{-1} \left( \mathbf{I} - \widehat{\mathcal{P}}_r(\mu_i) \right) \mathcal{B}(\mu_i) \mathbf{b}_i \\ &= \boldsymbol{\xi}_i^T \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i) \right) \mathcal{K}(\mu_i)^{-1} \left( \mathbf{I} - \widehat{\mathcal{P}}_r(\mu_i) \right) \boldsymbol{\eta}_i \\ &= \begin{cases} \boldsymbol{\xi}_i^T \mathcal{K}(\mu_i)^{-1} \left( \mathbf{I} - \widehat{\mathcal{P}}_r(\mu_i) \right) \boldsymbol{\eta}_i, & \text{or} \\ \boldsymbol{\xi}_i^T \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i) \right) \mathcal{K}(\mu_i)^{-1} \boldsymbol{\eta}_i, \end{cases} \end{aligned}$$

leading then to two estimates:

$$\left| \mathbf{c}_i^T \mathcal{H}(\mu_i) \mathbf{b}_i - \mathbf{c}_i^T \widehat{\mathcal{H}}_r(\mu_i) \mathbf{b}_i \right| \leq \|\boldsymbol{\xi}_i\| \cdot \|\boldsymbol{\eta}_i\| \cdot \|\mathcal{K}(\mu_i)^{-1}\| \cdot \|\mathbf{I} - \widehat{\mathcal{P}}_r(\mu_i)\|$$

and

$$\left| \mathbf{c}_i^T \mathcal{H}(\mu_i) \mathbf{b}_i - \mathbf{c}_i^T \widehat{\mathcal{H}}_r(\mu_i) \mathbf{b}_i \right| \leq \|\boldsymbol{\xi}_i\| \cdot \|\boldsymbol{\eta}_i\| \cdot \|\mathcal{K}(\mu_i)^{-1}\| \cdot \|\mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i)\|.$$

These can be combined to yield (28).

The last inequality comes from using (23) with  $s = \mu_i$ :

$$\begin{aligned} \mathbf{c}_i^T \mathcal{H}'(\mu_i) \mathbf{b}_i - \mathbf{c}_i^T \widehat{\mathcal{H}}_r'(\mu_i) \mathbf{b}_i &= \frac{d}{ds} \left[ \mathbf{c}_i^T \mathbf{e} \mathcal{K}^{-1} \right] \Big|_{\mu_i} \left( \mathbf{I} - \widehat{\mathcal{P}}_r(\mu_i) \right) \mathcal{B}(\mu_i) \mathbf{b}_i \\ &\quad + \mathbf{c}_i^T \mathbf{e}(\mu_i) \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i) \right) \frac{d}{ds} \left[ \mathcal{K}^{-1} \mathcal{B} \mathbf{b}_i \right] \Big|_{\mu_i} \\ &\quad - \mathbf{c}_i^T \mathbf{e}(\mu_i) \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i) \right) \frac{d}{ds} \left[ \mathcal{K}^{-1} \right] \Big|_{\mu_i} \left( \mathbf{I} - \widehat{\mathcal{P}}_r(\mu_i) \right) \mathcal{B}(\mu_i) \mathbf{b}_i. \end{aligned}$$

Then from (30), (32), and the Cauchy-Schwarz inequality

$$\begin{aligned} \left| \mathbf{c}_i^T \mathcal{H}'(\mu_i) \mathbf{b}_i - \mathbf{c}_i^T \widehat{\mathcal{H}}_r'(\mu_i) \mathbf{b}_i \right| &\leq \left| \frac{d}{ds} \left[ \mathbf{c}_i^T \mathbf{e} \mathcal{K}^{-1} \right] \Big|_{\mu_i} \left( \mathbf{I} - \widehat{\mathcal{P}}_r(\mu_i) \right) \boldsymbol{\eta}_i \right| \\ &\quad + \left| \boldsymbol{\xi}_i^T \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i) \right) \frac{d}{ds} \left[ \mathcal{K}^{-1} \mathcal{B} \mathbf{b}_i \right] \Big|_{\mu_i} \right| \\ &\quad + \left| \boldsymbol{\xi}_i^T \left( \mathbf{I} - \widehat{\mathcal{Q}}_r(\mu_i) \right) \frac{d}{ds} \left[ \mathcal{K}^{-1} \right] \Big|_{\mu_i} \left( \mathbf{I} - \widehat{\mathcal{P}}_r(\mu_i) \right) \boldsymbol{\eta}_i \right| \\ &\leq \left\| \frac{d}{ds} \left[ \mathbf{c}_i^T \mathbf{e} \mathcal{K}^{-1} \right] \Big|_{\mu_i} \right\| \cdot \frac{\|\boldsymbol{\eta}_i\|}{\cos \Theta(\widehat{\mathcal{P}}_r(\mu_i), \widehat{\mathcal{W}}_r)} \\ &\quad + \frac{\|\boldsymbol{\xi}_i\|}{\cos \Theta(\widehat{\mathcal{Q}}_r(\mu_i), \widehat{\mathcal{V}}_r)} \cdot \left\| \frac{d}{ds} \left[ \mathcal{K}^{-1} \mathcal{B} \mathbf{b}_i \right] \Big|_{\mu_i} \right\| \\ &\quad + \left\| \frac{d}{ds} \left[ \mathcal{K}^{-1} \right] \Big|_{\mu_i} \right\| \cdot \frac{\|\boldsymbol{\eta}_i\|}{\cos \Theta(\widehat{\mathcal{P}}_r(\mu_i), \widehat{\mathcal{W}}_r)} \frac{\|\boldsymbol{\xi}_i\|}{\cos \Theta(\widehat{\mathcal{Q}}_r(\mu_i), \widehat{\mathcal{V}}_r)} \end{aligned}$$

which yields the conclusion.  $\square$

Consider the effect of solving (11) and (12) approximately with successively increasing levels of accuracy that force the residual norms to zero,  $\|\boldsymbol{\eta}_j\| \rightarrow 0$  and  $\|\boldsymbol{\xi}_i\| \rightarrow 0$ . The multiplicative behavior of the error bound (28) with respect to  $\|\boldsymbol{\eta}_j\|$  and  $\|\boldsymbol{\xi}_i\|$  contrasts with the additive behavior seen in (26) and (27) and suggests some potential benefit in using the same interpolation points for both left and right interpolation, i.e., choosing  $\mu_i = \sigma_i$  for  $i = 1, \dots, r$ . Note that this choice also forces convergent (bitangential) derivative interpolation as shown in (29). Indeed, choosing  $\mu_i = \sigma_i$  for  $i = 1, \dots, r$  is a *necessary* condition for forming  $\mathcal{H}_2$ -optimal interpolatory reduced order models for first-order descriptor realizations, as we discuss in §5 (see also [16]). Beyond this, there can be notable computational advantages in choosing  $\mu_i = \sigma_i$ , since the linear systems to be solved in (11) and (12) then have the same coefficient matrix; allowing one potentially to reuse factorizations and preconditioners.

Certain applications require the retention of structural properties such as symmetry in passing from  $\mathcal{K}$  to  $\widehat{\mathcal{K}}_r$  and one is compelled to choose  $\widehat{\mathbf{W}}_r = \widehat{\mathbf{V}}_r$  (“one-sided” model reduction), so the vectors  $\{\widehat{\mathbf{w}}_1, \dots, \widehat{\mathbf{w}}_r\}$  might not be approximate solutions to (11) in the usual sense. Nonetheless, the behavior of the interpolation error is still governed by (26) and (27). We explore this in the following numerical example.

We illustrate the character of the results given in Theorem 3.1, bounding the response error at the nominal interpolation points caused by inexact solves in (11) and (12). To this end, we consider a delay differential equation of the form introduced in (4) taking  $n = 2000$ ,  $m = p = 1$  and  $\tau_{inp} = \tau_{out} = 0$ . The coefficient matrices for the full order model in (4) were taken from [6]. We construct multiple reduced models all of order  $r = 3$ , solving (11) and (12) with different levels of accuracy. We chose three logarithmically spaced values,  $\sigma_1 = 0.001$ ,  $\sigma_2 = 0.0316$ ,  $\sigma_3 = 1.0$ , and fixed them as interpolation points. We then obtained approximate solutions of varying accuracy to (11) and (12) in a manner described in more detail below, assembled the inexact interpolation basis matrices,  $\widehat{\mathbf{V}}_r$  and  $\widehat{\mathbf{W}}_r$ , and obtained reduced models of order  $r = 3$  having the same internal delay structure as the original system:

$$\begin{aligned} \widehat{\mathcal{H}}_r(s) &= \widehat{\mathbf{C}}_r(s) \widehat{\mathcal{K}}_r(s)^{-1} \widehat{\mathbf{B}}_r(s) \\ &= \mathbf{C} \widehat{\mathbf{V}}_r \left( s \widehat{\mathbf{W}}_r^T \mathbf{E} \widehat{\mathbf{V}}_r - \widehat{\mathbf{W}}_r^T \mathbf{A}_0 \widehat{\mathbf{V}}_r - e^{-s\tau_{sys}} \widehat{\mathbf{W}}_r^T \mathbf{A}_1 \widehat{\mathbf{V}}_r \right)^{-1} \widehat{\mathbf{W}}_r^T \mathbf{B} \end{aligned}$$

We considered both the usual “two-sided” model reduction process that in-

volves approximate solution of both (11) and (12) and the “one-sided” process that involves approximate solutions only to (12) to generate  $\widehat{\mathbf{V}}_r$  and then assigning  $\widehat{\mathbf{W}}_r = \widehat{\mathbf{V}}_r$ . Linear systems were solved with GMRES terminating with a final relative residual below a uniform tolerance denoted by  $\varepsilon$ .

We generated reduced order models in this way, varying the relative residual tolerance  $\varepsilon$  from  $10^{-1}$  down to  $10^{-8}$ . Figure 1 below shows the resulting interpolation errors  $|\mathfrak{H}(\sigma_1) - \widehat{\mathfrak{H}}_r(\sigma_1)|$  and bounds from equations (26) and (28) for one-sided and two-sided cases, respectively, as  $\varepsilon$  varies. Observe that the bounds in Theorem 3.1 predict the convergence behavior of the true error quite well; the rates (slopes) are matched almost exactly. Note also that the interpolation error decays much faster for two-sided reduction than for one-sided reduction. Indeed, the ratio of the two errors is close to  $\varepsilon$ , i.e., for a given tolerance  $\varepsilon$ , the interpolation error for two-sided reduction is approximately  $\varepsilon$  times smaller than the interpolation error for one-sided reduction.

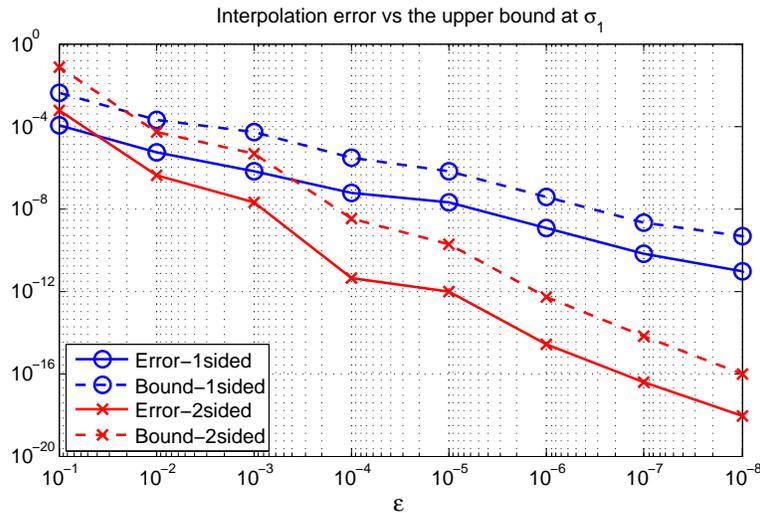


Figure 1: Behavior of interpolation error and upper bounds vs  $\varepsilon$

Analogous results regarding behavior of the bounds and interpolation error are observed at  $\sigma_2$  and  $\sigma_3$  and so are omitted for brevity.

### 3.2. Global Error Bounds

Thus far we have focussed on the extent to which interpolation properties are lost in the computed reduced models when inexact solves are introduced

into the process, considering in effect *local* error bounds. Clearly, it is important to understand the effect of inexact solves on the overall *global* quality of the reduced order model. There are two commonly used measures for closeness of two conforming dynamical systems (i.e., those with the same input and output dimensions):

$$\begin{aligned} \text{the } \mathcal{H}_2\text{-norm:} \quad & \|\mathcal{H} - \mathcal{G}\|_{\mathcal{H}_2} = \frac{1}{2\pi} \int_{-\infty}^{\infty} \|\mathcal{H}(i\omega) - \mathcal{G}(i\omega)\|_F^2 d\omega \\ \text{the } \mathcal{H}_\infty\text{-norm:} \quad & \|\mathcal{H} - \mathcal{G}\|_{\mathcal{H}_\infty} = \max_{\omega \in \mathbb{R}} \|\mathcal{H}(i\omega) - \mathcal{G}(i\omega)\|_2. \end{aligned}$$

Since reduced models are completely determined by the subspaces,  $\mathcal{V}_r$  and  $\mathcal{W}_r$ , as shown in (8), we first evaluate (in Theorem 3.2) how much inexact interpolatory subspaces,  $\widehat{\mathcal{V}}_r$  and  $\widehat{\mathcal{W}}_r$ , can deviate from the corresponding true subspaces,  $\mathcal{V}_r$  and  $\mathcal{W}_r$ , as a result of inexact solves. The effect of this deviation on the resulting model reduction (forward) error will be shown in Theorem 3.3. In this way, we are able to connect model reduction error to observable quantities that are associated with inexact solves, such as the relative stopping criterion  $\varepsilon$ .

**Theorem 3.2.** *Let the columns of  $\mathbf{V}_r$  and  $\widehat{\mathbf{V}}_r$  be exact and approximate solutions to (12) and the columns of  $\mathbf{W}_r$  and  $\widehat{\mathbf{W}}_r$  be exact and approximate solutions to (11). Suppose approximate solutions are computed to a relative residual tolerance of  $\varepsilon > 0$ , so that  $\|\boldsymbol{\eta}_i\| \leq \varepsilon \|\mathcal{B}(\sigma_i)\mathbf{b}_i\|$  and  $\|\boldsymbol{\xi}_i\| \leq \varepsilon \|\mathcal{C}(\mu_i)^T \mathbf{c}_i\|$ , where the residuals  $\boldsymbol{\eta}_i$  and  $\boldsymbol{\xi}_i$  are defined in (16).*

*Denoting the associated subspaces as  $\mathcal{V}_r$ ,  $\widehat{\mathcal{V}}_r$ ,  $\mathcal{W}_r$  and  $\widehat{\mathcal{W}}_r$  then*

$$\sin \Theta(\widehat{\mathcal{V}}_r, \mathcal{V}_r) \leq \frac{\varepsilon \sqrt{r}}{\varsigma_{\min}(\widehat{\mathbf{V}}_r \mathbf{D}_v)} \quad (33)$$

$$\sin \Theta(\widehat{\mathcal{W}}_r, \mathcal{W}_r) \leq \frac{\varepsilon \sqrt{r}}{\varsigma_{\min}(\widehat{\mathbf{W}}_r \mathbf{D}_w)} \quad (34)$$

where  $\mathbf{D}_v$  and  $\mathbf{D}_w$  are diagonal scaling matrices defined as

$$\begin{aligned} \mathbf{D}_v &= \text{diag}((\|\mathcal{K}(\sigma_1)^{-1}\| \|\mathcal{B}(\sigma_1)\mathbf{b}_1\|)^{-1}, \dots, (\|\mathcal{K}(\sigma_r)^{-1}\| \|\mathcal{B}(\sigma_r)\mathbf{b}_r\|)^{-1}) \text{ and} \\ \mathbf{D}_w &= \text{diag}((\|\mathcal{K}(\mu_1)^{-1}\| \|\mathcal{C}(\mu_1)^T \mathbf{c}_1\|)^{-1}, \dots, (\|\mathcal{K}(\mu_r)^{-1}\| \|\mathcal{C}(\mu_r)^T \mathbf{c}_r\|)^{-1}) \end{aligned}$$

and  $\varsigma_{\min}(\mathbf{M})$  denotes the smallest singular value of the matrix  $\mathbf{M}$ .

PROOF: We prove (33). The proof of (34) is similar.

Write  $\widehat{\mathbf{V}}_r = \mathbf{V}_r + \mathbf{E}$  with  $\mathbf{E} = [\mathbf{K}(\sigma_1)^{-1}\boldsymbol{\eta}_1, \dots, \mathbf{K}(\sigma_r)^{-1}\boldsymbol{\eta}_r]$ . Then

$$\begin{aligned} \sin \Theta(\widehat{\mathcal{V}}_r, \mathcal{V}_r) &= \max_{\widehat{\mathbf{v}} \in \widehat{\mathcal{V}}_r} \min_{\mathbf{v} \in \mathcal{V}_r} \frac{\|\mathbf{v} - \widehat{\mathbf{v}}\|}{\|\widehat{\mathbf{v}}\|} \\ &= \max_{x_i} \min_{z_i} \frac{\|\sum_{i=1}^r z_i \mathcal{K}(\sigma_i)^{-1} \mathcal{B}(\sigma_i) \mathbf{b}_i - \sum_{i=1}^r x_i \widehat{\mathbf{v}}_i\|}{\|\sum_{i=1}^r x_i \widehat{\mathbf{v}}_i\|} \\ &= \max_{x_i} \min_{z_i} \frac{\|\sum_{i=1}^r (z_i - x_i) \mathcal{K}(\sigma_i)^{-1} \mathcal{B}(\sigma_i) \mathbf{b}_i - x_i \mathcal{K}(\sigma_i)^{-1} \boldsymbol{\eta}_i\|}{\|\sum_{i=1}^r x_i \widehat{\mathbf{v}}_i\|} \\ &\leq \max_{x_i} \frac{\|\sum_{i=1}^r x_i \mathcal{K}(\sigma_i)^{-1} \boldsymbol{\eta}_i\|}{\|\sum_{i=1}^r x_i \widehat{\mathbf{v}}_i\|} = \max_{\mathbf{x}} \frac{\|\mathbf{E}\mathbf{x}\|}{\|\widehat{\mathbf{V}}_r \mathbf{x}\|} = \max_{\mathbf{x}} \frac{\|\mathbf{E}\mathbf{D}\mathbf{x}\|}{\|\widehat{\mathbf{V}}_r \mathbf{D}\mathbf{x}\|} \end{aligned}$$

where  $\mathbf{D} = \text{diag}(d_1, \dots, d_r)$  is a diagonal matrix with positive diagonal entries,  $d_i > 0$ , that are fixed but for the moment unspecified.

Note that

$$\begin{aligned} \|\mathbf{E}\mathbf{D}\mathbf{x}\| &\leq \|\mathbf{E}\mathbf{D}\| \|\mathbf{x}\| \leq \sqrt{r} \|\mathbf{x}\| \max_i (d_i \|\mathcal{K}(\sigma_i)^{-1} \boldsymbol{\eta}_i\|) \\ &\leq \sqrt{r} \|\mathbf{x}\| \max_i (d_i \|\mathcal{K}(\sigma_i)^{-1}\| \|\boldsymbol{\eta}_i\|) \end{aligned}$$

Thus we have,

$$\sin \Theta(\widehat{\mathcal{V}}_r, \mathcal{V}_r) \leq \sqrt{r} \frac{\max_i (d_i \|\mathcal{K}(\sigma_i)^{-1}\| \|\boldsymbol{\eta}_i\|)}{\min_{\mathbf{x}} (\|\widehat{\mathbf{V}}_r \mathbf{D}\mathbf{x}\| / \|\mathbf{x}\|)} = \sqrt{r} \frac{\max_i (d_i \|\mathcal{K}(\sigma_i)^{-1}\| \|\boldsymbol{\eta}_i\|)}{\varsigma_{\min}(\widehat{\mathbf{V}}_r \mathbf{D})} \quad (35)$$

This bound is valid for any choice of diagonal scalings,  $\mathbf{D}$ , so we can minimize the right hand side of (35) with respect to  $d_1, \dots, d_r$ . The *Column Equilibration Theorem* of van der Sluis [28] asserts that the optimal choice of  $d_1, \dots, d_r$  is such that  $d_i \|\mathcal{K}(\sigma_i)^{-1}\| \|\boldsymbol{\eta}_i\| = C$ , independent of  $i = 1, \dots, r$ . If inexact solves terminate with residuals satisfying  $\|\boldsymbol{\eta}_i\| \approx \varepsilon \|\mathcal{B}(\sigma_i) \mathbf{b}_i\|$  then we may take  $C = \varepsilon$  and  $d_i = (\|\mathcal{K}(\sigma_i)^{-1}\| \|\mathcal{B}(\sigma_i) \mathbf{b}_i\|)^{-1}$  to achieve the best bound possible with the information given. This leads to (33).  $\square$

As a practical matter, the column scalings used in (33) and (34) will not be computationally feasible in realistic settings. If instead we scale the columns of  $\widehat{\mathbf{V}}_r$  and  $\widehat{\mathbf{W}}_r$  to have unit norm (cheap!) — taking  $\widetilde{\mathbf{D}}_v = \text{diag}(1/\|\widehat{\mathbf{v}}_1\|, \dots, 1/\|\widehat{\mathbf{v}}_r\|)$  and  $\widetilde{\mathbf{D}}_w = \text{diag}(1/\|\widehat{\mathbf{w}}_1\|, \dots, 1/\|\widehat{\mathbf{w}}_r\|)$ , the bound for (33) degrades to

$$\sin \Theta(\widehat{\mathcal{V}}_r, \mathcal{V}_r) \leq \max_i \kappa_2(\mathcal{K}(\sigma_i), \widehat{\mathbf{v}}_i) \frac{\varepsilon \sqrt{r}}{\varsigma_{\min}(\widehat{\mathbf{V}}_r \widetilde{\mathbf{D}}_v)}$$

where  $\kappa_2(\mathcal{K}(\sigma_i), \widehat{\mathbf{v}}_i) = \frac{\|\mathcal{K}(\sigma_i)^{-1}\| \|\mathcal{B}(\sigma_i)\mathbf{b}_i\|}{\|\widehat{\mathbf{v}}_i\|} > 1$  is the *condition number* of the linear system (12). A similar expression holds for  $\sin \Theta(\widehat{\mathcal{W}}_r, \mathcal{W}_r)$ . In many cases, these condition numbers have only modest magnitude and the bounds (33) and (34) remain descriptive.

**Theorem 3.3.** *Let the columns of  $\mathbf{V}_r$  and  $\widehat{\mathbf{V}}_r$  be exact and approximate solutions to (12) and the columns of  $\mathbf{W}_r$  and  $\widehat{\mathbf{W}}_r$  be exact and approximate solutions to (11). Let the associated subspaces be denoted as  $\mathcal{V}_r$ ,  $\widehat{\mathcal{V}}_r$ ,  $\mathcal{W}_r$  and  $\widehat{\mathcal{W}}_r$  and the associated reduced order systems be denoted as  $\mathcal{H}_r(s)$  (exact) and  $\widehat{\mathcal{H}}_r(s)$  (inexact). Then*

$$\frac{\|\mathcal{H}_r - \widehat{\mathcal{H}}_r\|_{\mathcal{H}_\infty}}{\frac{1}{2}(\|\mathcal{H}_r\|_{\mathcal{H}_\infty} + \|\widehat{\mathcal{H}}_r\|_{\mathcal{H}_\infty})} \leq M \max\left(\sin \Theta(\widehat{\mathcal{V}}_r, \mathcal{V}_r), \sin \Theta(\widehat{\mathcal{W}}_r, \mathcal{W}_r)\right),$$

where

$$M = 2 \max\left(\frac{\max_{\omega \in \mathbb{R}} \text{cond}_{\mathcal{E}}(\mathcal{H}_r(\omega))}{\min_{\omega \in \mathbb{R}} \cos \Theta(\widehat{\mathcal{Q}}_r(\omega), \widehat{\mathcal{V}}_r)}, \frac{\max_{\omega \in \mathbb{R}} \text{cond}_{\mathcal{B}}(\widehat{\mathcal{H}}_r(\omega))}{\min_{\omega \in \mathbb{R}} \cos \Theta(\widehat{\mathcal{P}}_r(\omega), \mathcal{W}_r)}\right)$$

and

$$\begin{aligned} \text{cond}_{\mathcal{B}}(\widehat{\mathcal{H}}_r(s)) &= \frac{\|\widehat{\mathcal{C}}_r(s)\widehat{\mathcal{K}}_r(s)^{-1}\widehat{\mathbf{W}}_r^T\| \|\mathcal{B}(s)\|}{\|\widehat{\mathcal{H}}_r(s)\|} \\ \text{cond}_{\mathcal{E}}(\mathcal{H}_r(s)) &= \frac{\|\mathcal{C}(s)\| \|\mathbf{V}_r\mathcal{K}_r(s)^{-1}\mathcal{B}_r(s)\|}{\|\mathcal{H}_r(s)\|} \end{aligned}$$

PROOF: Note that for all  $s \in \mathbb{C}$  for which  $\mathcal{H}_r$  and  $\widehat{\mathcal{H}}_r$  are both analytic,

$$\begin{aligned} \|\mathcal{H}_r(s) - \widehat{\mathcal{H}}_r(s)\| &= \|\mathcal{C}(s) \left( \mathbf{V}_r\mathcal{K}_r(s)^{-1}\mathbf{W}_r^T - \widehat{\mathbf{V}}_r\widehat{\mathcal{K}}_r(s)^{-1}\widehat{\mathbf{W}}_r^T \right) \mathcal{B}(s)\| \\ &= \|\mathcal{C}(s) \left( \mathcal{Q}_r(s) - \widehat{\mathcal{Q}}_r(s) \right) \mathcal{K}(s)^{-1}\mathcal{B}(s)\| \\ &= \|\mathcal{C}(s) \left( (\mathbf{I} - \widehat{\mathcal{Q}}_r(s)) \mathcal{Q}_r(s) - \widehat{\mathcal{Q}}_r(s) (\mathbf{I} - \mathcal{Q}_r(s)) \right) \mathcal{K}(s)^{-1}\mathcal{B}(s)\| \end{aligned}$$

So,

$$\begin{aligned}
\|\mathcal{H}_r(s) - \widehat{\mathcal{H}}_r(s)\| &\leq \|\mathbf{C}(s) \left( \mathbf{I} - \widehat{\mathbf{Q}}_r(s) \right) \mathbf{Q}_r(s) \mathcal{K}(s)^{-1} \mathbf{B}(s)\| \\
&\quad + \|\mathbf{C}(s) \widehat{\mathbf{Q}}_r(s) (\mathbf{I} - \mathbf{Q}_r(s)) \mathcal{K}(s)^{-1} \mathbf{B}(s)\| \\
&\leq \|\mathbf{C}(s) \left( \mathbf{I} - \widehat{\mathbf{Q}}_r(s) \right) \mathbf{Q}_r(s) \mathcal{K}(s)^{-1} \mathbf{B}(s)\| \\
&\quad + \|\mathbf{C}(s) \mathcal{K}(s)^{-1} \widehat{\mathcal{P}}_r(s) (\mathbf{I} - \mathcal{P}_r(s)) \mathbf{B}(s)\| \\
&\leq \|\mathbf{C}(s) \left( \mathbf{I} - \widehat{\mathbf{Q}}_r(s) \right) \left( \mathbf{I} - \widehat{\mathbf{\Xi}} \right) \mathbf{\Xi} \mathbf{Q}_r(s) \mathcal{K}(s)^{-1} \mathbf{B}(s)\| \\
&\quad + \|\mathbf{C}(s) \mathcal{K}(s)^{-1} \widehat{\mathcal{P}}_r(s) \widehat{\mathbf{\Pi}} (\mathbf{I} - \mathbf{\Pi}) (\mathbf{I} - \mathcal{P}_r(s)) \mathbf{B}(s)\| \\
&\leq \|\mathbf{C}(s)\| \left\| \mathbf{I} - \widehat{\mathbf{Q}}_r(s) \right\| \left\| \left( \mathbf{I} - \widehat{\mathbf{\Xi}} \right) \mathbf{\Xi} \right\| \|\mathbf{Q}_r(s) \mathcal{K}(s)^{-1} \mathbf{B}(s)\| \\
&\quad + \|\mathbf{C}(s) \mathcal{K}(s)^{-1} \widehat{\mathcal{P}}_r(s)\| \|\widehat{\mathbf{\Pi}} (\mathbf{I} - \mathbf{\Pi})\| \|\mathbf{I} - \mathcal{P}_r(s)\| \|\mathbf{B}(s)\| \\
&\leq \|\mathbf{C}(s)\| \frac{\sin \Theta(\widehat{\mathcal{V}}_r, \mathcal{V}_r)}{\cos \Theta(\widehat{\mathcal{Q}}_r(s), \widehat{\mathcal{V}}_r)} \|\mathbf{Q}_r(s) \mathcal{K}(s)^{-1} \mathbf{B}(s)\| \\
&\quad + \|\mathbf{C}(s) \mathcal{K}(s)^{-1} \widehat{\mathcal{P}}_r(s)\| \frac{\sin \Theta(\widehat{\mathcal{W}}_r, \mathcal{W}_r)}{\cos \Theta(\widehat{\mathcal{P}}_r(s), \mathcal{W}_r)} \|\mathbf{B}(s)\| \\
&\leq \text{cond}_{\mathcal{E}}(\mathcal{H}_r(s)) \frac{\sin \Theta(\widehat{\mathcal{V}}_r, \mathcal{V}_r)}{\cos \Theta(\widehat{\mathcal{Q}}_r(s), \widehat{\mathcal{V}}_r)} \|\mathcal{H}_r(s)\| \\
&\quad + \|\widehat{\mathcal{H}}_r(s)\| \frac{\sin \Theta(\widehat{\mathcal{W}}_r, \mathcal{W}_r)}{\cos \Theta(\widehat{\mathcal{P}}_r(s), \mathcal{W}_r)} \text{cond}_{\mathcal{B}}(\widehat{\mathcal{H}}_r(s))
\end{aligned}$$

Maximizing over  $s = \imath\omega$  with  $\omega \in \mathbb{R}$  gives

$$\begin{aligned}
\|\mathcal{H}_r - \widehat{\mathcal{H}}_r\|_{\mathcal{H}_\infty} &\leq \max_{\omega \in \mathbb{R}} \text{cond}_{\mathcal{E}}(\mathcal{H}_r(\imath\omega)) \frac{\sin \Theta(\widehat{\mathcal{V}}_r, \mathcal{V}_r)}{\min_{\omega \in \mathbb{R}} \cos \Theta(\widehat{\mathcal{Q}}_r(s), \widehat{\mathcal{V}}_r)} \|\mathcal{H}_r\|_{\mathcal{H}_\infty} \\
&\quad + \max_{\omega \in \mathbb{R}} \text{cond}_{\mathcal{B}}(\widehat{\mathcal{H}}_r(\imath\omega)) \frac{\sin \Theta(\widehat{\mathcal{W}}_r, \mathcal{W}_r)}{\min_{\omega \in \mathbb{R}} \cos \Theta(\widehat{\mathcal{P}}_r(s), \mathcal{W}_r)} \|\widehat{\mathcal{H}}_r\|_{\mathcal{H}_\infty}
\end{aligned}$$

which leads immediately to the conclusion.  $\square$

### 3.3. Illustrative examples

The process to be modeled arises in cooling within a rolling mill and is modeled as boundary control of a two dimensional heat equation. A finite element discretization results in a descriptor system of the form

$$\mathbf{E}\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \quad y(t) = \mathbf{C}\mathbf{x}(t).$$

where  $\mathbf{A}, \mathbf{E} \in \mathbb{R}^{5177 \times 5177}$ ,  $\mathbf{B} \in \mathbb{R}^{5177 \times 7}$ ,  $\mathbf{C} \in \mathbb{R}^{6 \times 5177}$ . For simplicity, we focus on a SISO full-order subsystem that relates the sixth input to the second output. For details regarding the modeling, discretization, optimal control design, and model reduction, see [8, 9].

We show the results of interpolatory model reduction using an *ad hoc* choice of interpolation points: 6 logarithmically spaced points between  $10^{0.5}$  and 10; and an  $\mathcal{H}_2$ -optimal choice of interpolation points obtained by the method of [16]. For each case, we reduce the system order to  $r = 6$  using first *exact* interpolatory model reduction (i.e., the linear systems are solved directly) and then with inexact model reduction with varying choices of termination criteria. The resulting reduced-order models are denoted by  $\mathcal{H}_r(s)$  and  $\widehat{\mathcal{H}}_r(s)$ , respectively. To see the effect of the choice of interpolation points on the underlying model reduction problem, we vary the relative residual termination tolerance,  $\varepsilon$  between  $10^{-1}$  and  $10^{-10}$  and show how quickly  $\widehat{\mathcal{H}}_r(s)$  converges to  $\mathcal{H}_r(s)$  for both the *ad hoc* selection and the  $\mathcal{H}_2$ -optimal selection of interpolation points. Table 2 shows the relative  $\mathcal{H}_\infty$  error between  $\widehat{\mathcal{H}}_r(s)$  and  $\mathcal{H}_r(s)$  as  $\varepsilon$  decreases. For the  $\mathcal{H}_2$ -optimal choice of interpolation points,  $\widehat{\mathcal{H}}_r(s)$  converges to  $\mathcal{H}_r(s)$  as  $\varepsilon$  decreases, for the *ad hoc* choice of points, there is almost no improvement in accuracy until  $\varepsilon = 1 \times 10^{-6}$ .

$\varepsilon$	$\mathcal{H}_2$ -optimal $\{\sigma_i\}$	<i>ad hoc</i> $\{\sigma_i\}$
$10^{-1}$	$7.22 \times 10^{-1}$	$5.05 \times 10^{-1}$
$10^{-2}$	$2.00 \times 10^{-1}$	$1.64 \times 10^{-1}$
$10^{-3}$	$4.27 \times 10^{-2}$	$4.11 \times 10^{-1}$
$10^{-4}$	$1.07 \times 10^{-2}$	$2.38 \times 10^{-1}$
$10^{-5}$	$2.76 \times 10^{-4}$	$5.62 \times 10^{-1}$
$10^{-6}$	$2.56 \times 10^{-5}$	$2.13 \times 10^{-2}$
$10^{-7}$	$2.91 \times 10^{-6}$	$3.52 \times 10^{-3}$
$10^{-8}$	$1.51 \times 10^{-7}$	$6.18 \times 10^{-5}$
$10^{-9}$	$2.07 \times 10^{-8}$	$1.76 \times 10^{-5}$
$10^{-10}$	$2.17 \times 10^{-9}$	$5.15 \times 10^{-6}$

Table 2: The relative error  $\frac{\|\mathcal{H}_r - \widehat{\mathcal{H}}_r\|_{\mathcal{H}_\infty}}{\|\mathcal{H}_r\|_{\mathcal{H}_\infty}}$  as  $\varepsilon$  varies

The behavior exhibited in Table 2 becomes clearer once we inspect the subspace angles between the exact interpolatory subspaces  $\mathcal{V}_r$ ,  $\mathcal{W}_r$  and the in-

exact ones  $\widehat{\mathcal{V}}_r$  and  $\widehat{\mathcal{W}}_r$ . Table 3 shows the sine of the angle between the exact and inexact interpolatory subspaces as  $\varepsilon$  varies. While the gap decreases significantly as  $\varepsilon$  decreases for an  $\mathcal{H}_2$ -optimal selection of interpolation points, there is a much smaller improvement in the gap with respect to  $\varepsilon$  for an *ad hoc* choice of points. This behavior will be re-visited in more detail in §4.2 revealing that the  $\mathcal{H}_2$ -optimal (or good) interpolation points are expected to produce reduced order models that are more robust with respect to perturbations due to inexact solves.

$\varepsilon$	$\sin \Theta(\mathcal{V}_r, \widehat{\mathcal{V}}_r)$		$\sin \Theta(\mathcal{W}_r, \widehat{\mathcal{W}}_r)$	
	$\mathcal{H}_2$ -optimal $\{\sigma_i\}$	<i>ad hoc</i> $\{\sigma_i\}$	$\mathcal{H}_2$ -optimal $\{\sigma_i\}$	<i>ad hoc</i> $\{\sigma_i\}$
$10^{-1}$	$9.85 \times 10^{-1}$	$9.99 \times 10^{-1}$	$9.99 \times 10^{-1}$	$9.99 \times 10^{-1}$
$10^{-2}$	$1.99 \times 10^{-1}$	$9.99 \times 10^{-1}$	$9.97 \times 10^{-1}$	$9.93 \times 10^{-1}$
$10^{-3}$	$2.36 \times 10^{-2}$	$9.99 \times 10^{-1}$	$4.87 \times 10^{-1}$	$9.83 \times 10^{-1}$
$10^{-4}$	$4.39 \times 10^{-3}$	$9.60 \times 10^{-1}$	$6.38 \times 10^{-2}$	$9.99 \times 10^{-1}$
$10^{-5}$	$2.72 \times 10^{-4}$	$5.80 \times 10^{-1}$	$7.09 \times 10^{-3}$	$7.20 \times 10^{-1}$
$10^{-6}$	$2.90 \times 10^{-5}$	$4.57 \times 10^{-2}$	$9.88 \times 10^{-4}$	$1.19 \times 10^{-1}$
$10^{-7}$	$3.46 \times 10^{-6}$	$6.90 \times 10^{-3}$	$6.87 \times 10^{-5}$	$2.00 \times 10^{-2}$
$10^{-8}$	$3.85 \times 10^{-7}$	$7.92 \times 10^{-4}$	$6.71 \times 10^{-6}$	$2.26 \times 10^{-3}$
$10^{-9}$	$3.63 \times 10^{-8}$	$1.01 \times 10^{-4}$	$9.16 \times 10^{-7}$	$2.60 \times 10^{-4}$
$10^{-10}$	$2.71 \times 10^{-9}$	$1.28 \times 10^{-5}$	$6.35 \times 10^{-8}$	$3.10 \times 10^{-5}$

Table 3:  $r = 6$ ;  $\sin \Theta(\mathcal{V}_r, \widehat{\mathcal{V}}_r)$  and  $\sin \Theta(\mathcal{W}_r, \widehat{\mathcal{W}}_r)$  as  $\varepsilon$  varies

#### 4. Backward error

Instead of seeking bounds on how much an inexactly computed reduced model differs from an exactly computed counterpart, one may view an inexactly computed reduced order model as an exactly computed reduced order model of a perturbed full order system. That is, we wish to find a full order system

$$\widetilde{\mathcal{H}}(s) = \widetilde{\mathcal{C}}(s)\widetilde{\mathcal{K}}(s)^{-1}\widetilde{\mathcal{B}}(s) \quad (36)$$

so that the inexactly computed reduced model for  $\mathcal{H}(s) = \mathcal{C}(s)\mathcal{K}(s)^{-1}\mathcal{B}(s)$  would be an *exactly* computed interpolatory reduced model for  $\widetilde{\mathcal{H}}(s)$ . Given left and right tangential interpolation data as in (9) and (10) that has contributed toward producing the inexactly computed interpolatory reduced

model  $\widehat{\mathcal{H}}_r(s)$ , find  $\widetilde{\mathcal{H}}(s)$  as in (36) so that

$$\begin{aligned} \mathbf{c}_i^T \widetilde{\mathcal{H}}(\mu_i) &= \mathbf{c}_i^T \widehat{\mathcal{H}}_r(\mu_i) \quad \text{for } i = 1, \dots, r, \text{ and} \\ \widetilde{\mathcal{H}}(\sigma_j) \mathbf{b}_j &= \widehat{\mathcal{H}}_r(\sigma_j) \mathbf{b}_j \quad \text{for } j = 1, \dots, r. \end{aligned}$$

and so that  $\widehat{\mathcal{H}}_r$  could have been computed from the perturbed system  $\widetilde{\mathcal{H}}$  from the given tangential interpolation data via an exact computation. Specifically, given computed (inexact) projecting bases

$$\widehat{\mathbf{V}}_r = [\widehat{\mathbf{v}}_1, \dots, \widehat{\mathbf{v}}_r] \quad \widehat{\mathbf{W}}_r^T = \begin{bmatrix} \widehat{\mathbf{w}}_1^T \\ \vdots \\ \widehat{\mathbf{w}}_r^T \end{bmatrix}.$$

as in (15), and a resulting (inexact) reduced order coprime realization

$$\widehat{\mathcal{H}}_r(s) = \widehat{\mathbf{C}}_r(s) \widehat{\mathcal{K}}_r(s)^{-1} \widehat{\mathbf{B}}_r(s),$$

find a full-order system  $\widetilde{\mathcal{H}}(s) = \widetilde{\mathbf{C}}(s) \widetilde{\mathcal{K}}(s)^{-1} \widetilde{\mathbf{B}}(s)$  so that left and right interpolation conditions hold:

$$\mathbf{c}_i^T \widetilde{\mathbf{C}}(\mu_i) = \widehat{\mathbf{w}}_i^T \widetilde{\mathcal{K}}(\mu_i) \quad \text{for } i = 1, \dots, r, \quad (37)$$

$$\widetilde{\mathcal{K}}(\sigma_j) \widehat{\mathbf{v}}_j = \widetilde{\mathbf{B}}(\sigma_j) \mathbf{b}_j \quad \text{for } j = 1, \dots, r, \quad (38)$$

and so that

$$\widehat{\mathcal{K}}_r(s) = \widehat{\mathbf{W}}_r^T \widetilde{\mathcal{K}}(s) \widehat{\mathbf{V}}_r, \quad \widehat{\mathbf{B}}_r(s) = \widehat{\mathbf{W}}_r^T \widetilde{\mathbf{B}}(s), \quad \text{and} \quad \widehat{\mathbf{C}}_r(s) = \widetilde{\mathbf{C}}(s) \widehat{\mathbf{V}}_r, \quad (39)$$

There (typically) will be an infinite number of possible systems,  $\widetilde{\mathcal{H}}$ , that are consistent with the computed reduced system  $\widehat{\mathcal{H}}_r$  in this sense — we are interested in those that are *close* to the original system  $\mathcal{H}$  with respect to a convenient system norm such as  $\mathcal{H}_\infty$  or  $\mathcal{H}_2$ . In order to proceed, it is convenient to restrict the class of backwardly compatible systems,  $\widetilde{\mathcal{H}}$ . We consider those that have realizations that are *constant* perturbations from the corresponding original system factors:

$$\widetilde{\mathcal{K}}(s) = \mathcal{K}(s) + \mathbf{F}, \quad \widetilde{\mathbf{B}}(s) = \mathbf{B}(s) + \mathbf{E}, \quad \text{and} \quad \widetilde{\mathbf{C}}(s) = \mathbf{C}(s) + \mathbf{G}. \quad (40)$$

where  $\mathbf{E}$ ,  $\mathbf{F}$ , and  $\mathbf{G}$  are *constant* matrices. The conditions (37), (38), and (39) impose constraints on  $\mathbf{E}$ ,  $\mathbf{F}$ , and  $\mathbf{G}$ . Indeed, (37) and (38) imply that

$$\begin{aligned} \widehat{\mathbf{w}}_i^T \mathbf{F} + \boldsymbol{\xi}_i^T &= \mathbf{c}_i^T \mathbf{E} \quad \text{for } i = 1, \dots, r, \text{ and} \\ \mathbf{F} \widehat{\mathbf{v}}_j + \boldsymbol{\eta}_j &= \mathbf{G} \mathbf{b}_j \quad \text{for } j = 1, \dots, r. \end{aligned}$$

(39) implies that

$$\widehat{\mathbf{W}}_r^T \mathbf{F} \widehat{\mathbf{V}}_r = \mathbf{0}, \quad \widehat{\mathbf{W}}_r^T \mathbf{G} = \mathbf{0}, \quad \text{and} \quad \mathbf{E} \widehat{\mathbf{V}}_r = \mathbf{0}.$$

Taken together, we find that backward perturbations of the form (40) can exist only if

$$\boldsymbol{\xi}_i^T \widehat{\mathbf{V}}_r = \mathbf{0} \quad \text{for } i = 1, \dots, r, \quad \text{and} \quad \widehat{\mathbf{W}}_r^T \boldsymbol{\eta}_j = \mathbf{0} \quad \text{for } j = 1, \dots, r. \quad (41)$$

Thus, we find constraints on the inexact interpolation residuals  $\boldsymbol{\xi}_i$  and  $\boldsymbol{\eta}_j$  in order for a backwardly compatible system of the form (40) to exist. More complicated perturbation classes than (40) may be considered that would allow us to remove the conditions (41), of course, but instead we choose to focus on a computational framework that guarantees (41). The Biconjugate Gradient Algorithm (BiCG) will be an example of an iterative solution strategy that fits this framework [1, 4]; others can be constructed without difficulty, although many standard strategies such as GMRES, do not fit this framework.

#### 4.1. The Petrov-Galerkin Framework for Inexact Solves

We have observed above that (41) is necessary for there to be a well-defined backward error of the form (40) to exist. The simplest framework within which one may generate reduced order models that are guaranteed to satisfy this condition involves a Petrov-Galerkin formalism for producing approximate solutions to (11) and (12). For simplicity, we restrict our discussion to the case that  $\mu_i = \sigma_i$  (identical left and right interpolation points).

Let  $\mathcal{P}_N$  and  $\mathcal{Q}_N$  be  $N$ -dimensional subspaces of  $\mathbb{C}^n$  satisfying a nondegeneracy condition:  $(\mathcal{K}(\sigma_i)\mathcal{P}_N)^\perp \cap \mathcal{Q}_N = \{0\}$  for all shifts,  $\sigma_i$  to be considered. The *Petrov-Galerkin framework* for generating approximate solutions to the interpolation conditions (11) and (12) proceeds as follows:

$$\begin{aligned} \text{Find } \tilde{\mathbf{v}}_j \in \mathcal{P}_N \text{ so that } \mathcal{K}(\sigma_j)\tilde{\mathbf{v}}_j - \mathcal{B}(\sigma_j)\mathbf{b}_j &\perp \mathcal{Q}_N \quad \text{and} \\ \text{find } \tilde{\mathbf{w}}_j \in \mathcal{Q}_N \text{ so that } \mathcal{K}(\sigma_j)^T \tilde{\mathbf{w}}_j - \mathcal{C}(\sigma_j)^T \mathbf{c}_j &\perp \mathcal{P}_N \end{aligned} \quad (42)$$

Computed quantities generated within a Petrov-Galerkin framework will be denoted with a “tilde” to distinguish them from earlier “hat” quantities where no structure was assumed in the inexact solves. The following theorem asserts

that if a reduced order model is computed within a Petrov-Galerkin framework (42), then one can obtain a structured backward error that throws the effect of inexact solves back onto a perturbation on the original dynamical system.

**Theorem 4.1.** *Given a full order model  $\mathcal{H}(s) = \mathcal{C}(s)\mathcal{K}(s)^{-1}\mathcal{B}(s)$ , interpolation points  $\{\sigma_j\}_{j=1}^r$ , and tangent directions  $\{\mathbf{b}_i\}_{i=1}^r$  and  $\{\mathbf{c}_i\}_{i=1}^r$ , let the inexact solutions  $\tilde{\mathbf{v}}_j$  for  $\mathcal{K}(\sigma_j)^{-1}\mathcal{B}(\sigma_j)\mathbf{b}_j$  and  $\tilde{\mathbf{w}}_j$  for  $\mathcal{K}(\sigma_j)^{-T}\mathcal{C}(\sigma_j)^T\mathbf{c}_j$  be obtained in a Petrov-Galerkin framework as in (42). Let  $\tilde{\mathbf{V}}_r$  and  $\tilde{\mathbf{W}}_r$  denote the corresponding inexact interpolatory bases; i.e.*

$$\tilde{\mathbf{V}}_r = [\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_r] \quad \text{and} \quad \tilde{\mathbf{W}}_r = [\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_r]. \quad (43)$$

Define residuals

$$\boldsymbol{\eta}_j = \mathcal{K}(\sigma_j)\tilde{\mathbf{v}}_j - \mathcal{B}(\sigma_j)\mathbf{b}_j \quad \text{and} \quad \boldsymbol{\xi}_j = \mathcal{K}(\sigma_j)^T\tilde{\mathbf{w}}_j - \mathcal{C}(\sigma_j)^T\mathbf{c}_j,$$

residual matrices

$$\mathbf{R}_b = [\boldsymbol{\eta}_1, \boldsymbol{\eta}_2, \dots, \boldsymbol{\eta}_r], \quad \mathbf{R}_c = [\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_r], \quad (44)$$

and the rank  $2r$  matrix

$$\mathbf{F}_{2r} = \mathbf{R}_b(\tilde{\mathbf{W}}_r^T\tilde{\mathbf{V}}_r)^{-1}\tilde{\mathbf{W}}_r^T + \tilde{\mathbf{V}}_r(\tilde{\mathbf{W}}_r^T\tilde{\mathbf{V}}_r)^{-1}\mathbf{R}_c^T. \quad (45)$$

Let  $\tilde{\mathcal{H}}_r(s) = \tilde{\mathcal{C}}_r(s)\tilde{\mathcal{K}}_r(s)^{-1}\tilde{\mathcal{B}}_r(s)$  denote the computed inexact reduced model via the Petrov-Galerkin process where

$$\tilde{\mathcal{K}}_r(s) = \tilde{\mathbf{W}}_r^T\mathcal{K}(s)\tilde{\mathbf{V}}_r, \quad \tilde{\mathcal{B}}_r(s) = \tilde{\mathbf{W}}_r^T\mathcal{B}(s), \quad \text{and} \quad \tilde{\mathcal{C}}_r(s) = \mathcal{C}(s)\tilde{\mathbf{V}}_r. \quad (46)$$

Then,  $\tilde{\mathcal{H}}_r(s)$  exactly tangentially interpolates the perturbed full-order model

$$\tilde{\mathcal{H}}(s) = \mathcal{C}(s)(\mathcal{K}(s) + \mathbf{F}_{2r})^{-1}\mathcal{B}(s), \quad (47)$$

at each  $\sigma_i$ :

$$\begin{aligned} \tilde{\mathcal{H}}(\sigma_i)\mathbf{b}_i &= \tilde{\mathcal{H}}_r(\sigma_i)\mathbf{b}_i, \quad \mathbf{c}_i^T\tilde{\mathcal{H}}(\sigma_i) = \mathbf{c}_i^T\tilde{\mathcal{H}}_r(\sigma_i), \\ \text{and } \mathbf{c}_i^T\tilde{\mathcal{H}}'(\sigma_i)\mathbf{b}_i &= \mathbf{c}_i^T\tilde{\mathcal{H}}_r'(\sigma_i)\mathbf{b}_i \quad \text{for each } i = 1, \dots, r. \end{aligned}$$

PROOF: The computed model,  $\tilde{\mathcal{H}}_r(s)$ , will (exactly) tangentially interpolate a perturbed model  $\tilde{\mathcal{H}}(s) = \mathbf{C}(s)(\mathcal{K}(s) + \mathbf{F})^{-1}\mathbf{B}(s)$  provided the following interpolation conditions hold:

$$(\mathcal{K}(\sigma_i) + \mathbf{F})\tilde{\mathbf{v}}_i = \mathbf{B}(\sigma_i)\mathbf{b}_i \text{ and } \tilde{\mathbf{w}}_i^T(\mathcal{K}(\sigma_i) + \mathbf{F}) = \mathbf{c}_i^T\mathbf{C}(\sigma_i) \text{ for } i = 1, \dots, r.$$

Equivalently, these can be interpreted as conditions on the perturbation  $\mathbf{F}$ . Rewriting this using notation defined above,  $\mathbf{F}$  must satisfy

$$\mathbf{F}\tilde{\mathbf{V}}_r = \mathbf{R}_b \text{ and } \tilde{\mathbf{W}}_r^T\mathbf{F} = \mathbf{R}_c^T. \quad (48)$$

The Petrov-Galerkin framework guarantees  $\tilde{\mathbf{W}}_r^T\mathbf{R}_b = \mathbf{0}$  and  $\mathbf{R}_c^T\tilde{\mathbf{V}}_r = \mathbf{0}$ . Substitution of  $\mathbf{F}_{2r}$  from (45) into (48) verifies that  $\mathbf{F}_{2r}$  is a perturbation to  $\mathcal{K}(s)$  for which the computed (inexact) vectors become (exact) interpolation vectors.

Note that since  $\tilde{\mathbf{W}}_r^T\mathbf{F}_{2r}\tilde{\mathbf{V}}_r = \mathbf{0}$ ,

$$\tilde{\mathcal{K}}_r(s) = \tilde{\mathbf{W}}_r^T\mathcal{K}(s)\tilde{\mathbf{V}}_r = \tilde{\mathbf{W}}_r^T(\mathcal{K}(s) + \mathbf{F}_{2r})\tilde{\mathbf{V}}_r.$$

Consequently, the reduced model  $\tilde{\mathcal{H}}_r(s)$  obtained by inexact solves in (46) is what one would have obtained by exact interpolatory model reduction of  $\tilde{\mathcal{H}}(s)$ .  $\square$

**Theorem 4.2.** *Assume the hypotheses of Theorem 4.1 and that  $\tilde{\mathbf{W}}_r^T\tilde{\mathbf{V}}_r$  is nonsingular. Define an oblique projector,  $\tilde{\Phi}_r = \tilde{\mathbf{V}}_r(\tilde{\mathbf{W}}_r^T\tilde{\mathbf{V}}_r)^{-1}\tilde{\mathbf{W}}_r^T$ . The backward perturbation  $\mathbf{F}_{2r}$  given in Theorem 4.1 satisfies*

$$\|\mathbf{F}_{2r}\|_F \leq \sqrt{r} \|\tilde{\Phi}_r\| \cdot \left( \max_i \frac{\|\boldsymbol{\eta}_i\|}{\|\tilde{\mathbf{v}}_i\|} \varsigma_{\min}(\tilde{\mathbf{V}}_r\mathbf{D})^{-1} + \max_i \frac{\|\boldsymbol{\xi}_i\|}{\|\tilde{\mathbf{w}}_i\|} \varsigma_{\min}(\tilde{\mathbf{W}}_r\mathbf{D})^{-1} \right)$$

where  $\varsigma_{\min}$  denotes the smallest singular value and  $\|\mathbf{M}\|_F = \sqrt{\text{trace}(\mathbf{M}^T\mathbf{M})}$  denotes the Frobenius norm of a matrix,  $\mathbf{M}$ .

PROOF: Note that

$$\|\mathbf{F}_{2r}\|_F \leq \|\mathbf{R}_b(\tilde{\mathbf{W}}_r^T\tilde{\mathbf{V}}_r)^{-1}\tilde{\mathbf{W}}_r^T\|_F + \|\tilde{\mathbf{V}}_r(\tilde{\mathbf{W}}_r^T\tilde{\mathbf{V}}_r)^{-1}\mathbf{R}_c^T\|_F.$$

Let  $\tilde{\mathbf{V}}_r$  have an orthogonal factorization as  $\tilde{\mathbf{V}}_r = \mathbf{Q}_v \mathbf{L}_v$  with  $\mathbf{Q}_v^* \mathbf{Q}_v = \mathbf{I}$ . Then

$$\begin{aligned}
\|\mathbf{R}_b (\tilde{\mathbf{W}}_r^T \tilde{\mathbf{V}}_r)^{-1} \tilde{\mathbf{W}}_r^T\|_F &= \|\mathbf{R}_b \mathbf{L}_v^{-1} \mathbf{L}_v (\tilde{\mathbf{W}}_r^T \tilde{\mathbf{V}}_r)^{-1} \tilde{\mathbf{W}}_r^T\|_F \\
&\leq \|\mathbf{R}_b \mathbf{L}_v^{-1}\|_F \cdot \|\mathbf{L}_v (\tilde{\mathbf{W}}_r^T \tilde{\mathbf{V}}_r)^{-1} \tilde{\mathbf{W}}_r^T\| \\
&\leq \|\mathbf{R}_b \mathbf{L}_v^{-1}\|_F \cdot \|\tilde{\Phi}_r\| \\
&\leq \|\mathbf{R}_b \tilde{\mathbf{D}}_v (\mathbf{L}_v \tilde{\mathbf{D}}_v)^{-1}\|_F \cdot \|\tilde{\Phi}_r\| \\
&\leq \|\mathbf{R}_b \tilde{\mathbf{D}}_v\|_F \cdot \|(\mathbf{L}_v \tilde{\mathbf{D}}_v)^{-1}\| \cdot \|\tilde{\Phi}_r\|
\end{aligned}$$

where we have introduced a diagonal scaling matrix

$$\tilde{\mathbf{D}}_v = \text{diag}(1/\|\tilde{\mathbf{v}}_1\|, 1/\|\tilde{\mathbf{v}}_2\|, \dots, 1/\|\tilde{\mathbf{v}}_r\|).$$

Easily one sees  $\|\mathbf{R}_b \tilde{\mathbf{D}}_v\|_F \leq \sqrt{r} \max_i \frac{\|\boldsymbol{\eta}_i\|}{\|\tilde{\mathbf{v}}_i\|}$ . For the remaining term, note that

$$\|(\mathbf{L}_v \tilde{\mathbf{D}}_v)^{-1}\| = \left( \min_{\mathbf{x}} \frac{\|\tilde{\mathbf{V}}_r \tilde{\mathbf{D}}_v \mathbf{x}\|}{\|\mathbf{x}\|} \right)^{-1} = \varsigma_{\min}(\tilde{\mathbf{V}}_r \tilde{\mathbf{D}}_v)^{-1}$$

A similar bound for  $\|\tilde{\mathbf{V}}_r (\tilde{\mathbf{W}}_r^T \tilde{\mathbf{V}}_r)^{-1} \mathbf{R}_c^T\|_F$  is produced by an analogous process, which leads then to the final estimate for  $\|\mathbf{F}_{2r}\|_F$ .  $\square$

Note that the perturbation  $\mathbf{F}_{2r}$  is completely determined by accessible, computed quantities. Hence, one can use  $\mathbf{F}_{2r}$  to determine how accurately one must solve the underlying linear systems in order to assure system fidelity of a given order.

**Theorem 4.3.** *If  $\|\mathbf{F}_{2r}\| < 1/\|\mathcal{K}(s)^{-1}\|_{\mathcal{H}_\infty}$  then*

$$\|\mathcal{H}(s) - \tilde{\mathcal{H}}(s)\|_{\mathcal{H}_2} \leq \frac{\|\mathcal{C}(s) \mathcal{K}(s)^{-1}\|_{\mathcal{H}_2} \|\mathcal{K}(s)^{-1} \mathcal{B}(s)\|_{\mathcal{H}_\infty}}{1 - \|\mathcal{K}(s)^{-1}\|_{\mathcal{H}_\infty} \|\mathbf{F}_{2r}\|} \|\mathbf{F}_{2r}\|$$

**PROOF:** The system-wise backward error associated with inexact solves may be written as

$$\begin{aligned}
\mathcal{H}(s) - \tilde{\mathcal{H}}(s) &= \mathcal{C}(s) \mathcal{K}(s)^{-1} \mathcal{B}(s) - \mathcal{C}(s) (\mathcal{K}(s) + \mathbf{F}_{2r})^{-1} \mathcal{B}(s) \\
&= \mathcal{C}(s) \mathcal{K}(s)^{-1} \mathbf{F}_{2r} (\mathcal{K}(s) + \mathbf{F}_{2r})^{-1} \mathcal{B}(s) \\
&= \mathcal{C}(s) \mathcal{K}(s)^{-1} \mathbf{F}_{2r} (\mathbf{I} + \mathcal{K}(s)^{-1} \mathbf{F}_{2r})^{-1} \mathcal{K}(s)^{-1} \mathcal{B}(s)
\end{aligned}$$

Define  $\mathbf{M}(s) = \mathbf{F}_{2r} (\mathbf{I} + \mathcal{K}(s)^{-1} \mathbf{F}_{2r})^{-1}$  and observe that

$$\begin{aligned}
\|\mathcal{H}(s) - \tilde{\mathcal{H}}(s)\|_{\mathcal{H}_2}^2 &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \|\mathbf{C}(i\omega) \mathcal{K}(i\omega)^{-1} \mathbf{M}(i\omega) \mathcal{K}(i\omega)^{-1} \mathbf{B}(i\omega)\|_F^2 d\omega \\
&\leq \frac{1}{2\pi} \int_{-\infty}^{\infty} \|\mathbf{C}(i\omega) \mathcal{K}(i\omega)^{-1}\|_F^2 \cdot \|\mathbf{M}(i\omega)\|^2 \cdot \|\mathcal{K}(i\omega)^{-1} \mathbf{B}(i\omega)\|^2 d\omega \\
&\leq \left( \frac{1}{2\pi} \int_{-\infty}^{\infty} \|\mathbf{C}(i\omega) \mathcal{K}(i\omega)^{-1}\|_F^2 d\omega \right) \cdot \max_{\omega} \|\mathbf{M}(i\omega)\|^2 \cdot \max_{\omega} \|\mathcal{K}(i\omega)^{-1} \mathbf{B}(i\omega)\|^2 \\
&\leq \|\mathbf{C}(s) \mathcal{K}(s)^{-1}\|_{\mathcal{H}_2}^2 \cdot \|\mathcal{K}(s)^{-1} \mathbf{B}(s)\|_{\mathcal{H}_\infty}^2 \cdot \|\mathbf{M}(s)\|_{\mathcal{H}_\infty}^2.
\end{aligned}$$

To estimate  $\|\mathbf{M}(s)\|_{\mathcal{H}_\infty}$ , a rearrangement of the definition of  $\mathbf{M}(s)$  provides

$$\mathbf{M}(s) = (\mathbf{I} - \mathbf{M}(s) \mathcal{K}(s)^{-1}) \mathbf{F}_{2r}.$$

So we have immediately,

$$\begin{aligned}
\|\mathbf{M}(s)\|_{\mathcal{H}_\infty} &= \max_{\omega \in \mathbb{R}} \|\mathbf{M}(i\omega)\| \leq \max_{\omega \in \mathbb{R}} \|\mathbf{I} - \mathbf{M}(i\omega) \mathcal{K}(i\omega)^{-1}\| \cdot \|\mathbf{F}_{2r}\| \\
&\leq \left( 1 + \max_{\omega \in \mathbb{R}} \|\mathbf{M}(i\omega) \mathcal{K}(i\omega)^{-1}\| \right) \|\mathbf{F}_{2r}\| \\
&\leq (1 + \|\mathbf{M}(s)\|_{\mathcal{H}_\infty} \|\mathcal{K}(s)^{-1}\|_{\mathcal{H}_\infty}) \|\mathbf{F}_{2r}\|
\end{aligned}$$

Since  $\|\mathcal{K}(s)^{-1}\|_{\mathcal{H}_\infty} \|\mathbf{F}_{2r}\| < 1$ , this last expression can be rearranged to obtain

$$\|\mathbf{M}(s)\|_{\mathcal{H}_\infty} \leq \frac{\|\mathbf{F}_{2r}\|}{1 - \|\mathcal{K}(s)^{-1}\|_{\mathcal{H}_\infty} \|\mathbf{F}_{2r}\|}$$

which implies the conclusion.  $\square$

By combining Theorem 3.3 with Theorem 3.2 or combining Theorem 4.2 with Theorem 4.3, we approach our goal of connecting quantities that we have control over, such as the termination threshold,  $\varepsilon$ , to relevant system theoretic errors,  $\|\mathcal{H}_r - \hat{\mathcal{H}}_r\|$  and  $\|\mathcal{H} - \tilde{\mathcal{H}}\|$ , which are quantities we would like to control.

One may use these expressions as a basis to devise and investigate different, effective stopping criteria in large-scale numerical settings. For example, while  $\varepsilon$  appears explicitly in Theorem 3.2 in a way that suggests its use as a relative residual norm threshold; while Theorem 4.2 suggests a scaling of the residual norm by the norm of the solution vector as another possible stopping criterion. These and related ideas are the focus of on-going work.

#### 4.2. Quantities of interest in derived bounds

By combining Theorem 4.2 with Theorem 4.3, one observes that perturbation effects of the inexact solves on the system theoretical (model reduction related) measures critically depend on the four quantities: The norm of the oblique projector  $\tilde{\Phi}_r = \tilde{\mathbf{V}}_r(\tilde{\mathbf{W}}_r^T \tilde{\mathbf{V}}_r)^{-1} \tilde{\mathbf{W}}_r^T$  of the underlying model reduction problem, reciprocals of the minimum singular values of the scaled primitive bases  $\tilde{\mathbf{V}}_r \mathbf{D}$  and  $\tilde{\mathbf{W}}_r \mathbf{D}$ ; and the stopping criterion  $\varepsilon$  for the inexact solves, (which affects  $\max_i \frac{\|\boldsymbol{\eta}_i\|}{\|\tilde{\mathbf{v}}_i\|}$  and  $\max_i \frac{\|\boldsymbol{\xi}_i\|}{\|\tilde{\mathbf{w}}_i\|}$ .)

The  $\varepsilon$  term is associated directly with inexact solves and is under the control of the user. The remaining quantities  $\varsigma_{\min}(\tilde{\mathbf{V}}_r \mathbf{D})^{-1}$ ,  $\varsigma_{\min}(\tilde{\mathbf{W}}_r \mathbf{D})^{-1}$  and  $\|\tilde{\Phi}_r\|$ , depend largely on the selection of interpolation points  $\{\sigma_i\}$  and tangent directions, but the influence of interpolation data on the magnitude of these quantities is difficult to anticipate.

In this section, we will investigate experimentally the effects of the interpolation point selection on the three quantities of interest,  $\varsigma_{\min}(\tilde{\mathbf{V}}_r \mathbf{D})^{-1}$ ,  $\varsigma_{\min}(\tilde{\mathbf{W}}_r \mathbf{D})^{-1}$  and  $\|\tilde{\Phi}_r\|$ , appearing in the derived bounds. These quantities are continuous with respect to the primitive basis vectors,  $\{\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_r\}$  and  $\{\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_r\}$  in neighborhoods where  $\tilde{\mathbf{W}}_r^T \tilde{\mathbf{V}}_r$  is nonsingular (i.e., where the projector  $\tilde{\Phi}_r$  is well defined). Thus it will be sufficient to examine how the magnitudes of the quantities of interest depend on interpolation data presuming that the necessary linear solves are done *exactly*; for modest convergence thresholds, the effect of inexact solves on these magnitudes is secondary to the effect of interpolation point location.

For our numerical study, we use the International Space Station 12A Module as the full-order model. The model has order  $n = 1412$ . We examine a single-input single-output subsystem,  $\mathcal{H}(s)$ , reducing the order from 1412 to order  $r$  with  $r$  varying from 2 to 70 in increments of two. For each reduced order, we chose 2000 random shift selections and computed  $\varsigma_{\min}(\mathbf{V}_r \mathbf{D})^{-1}$ ,  $\varsigma_{\min}(\mathbf{W}_r \mathbf{D})^{-1}$  and  $\|\Phi_r\|$ . For each  $r$ ,  $r/2$  shifts were sampled from a uniform distribution on a rectangular region in the positive half-plane:  $\left\{ z \in \mathbb{C} \mid \begin{array}{l} \min_{\lambda} |\operatorname{Re}(\lambda)| \leq \operatorname{Re}(z) \leq \max_{\lambda} |\operatorname{Re}(\lambda)| \\ |\operatorname{Im}(z)| \leq \max_{\lambda} |\operatorname{Im}(\lambda)| \end{array} \right\}$ , where the max and min are taken over all the poles of the system. The remaining  $r/2$  shifts were taken to be the complex conjugates of this random sample, so as to produce a shift configuration that was closed under conjugation. Additionally for each  $r$ , we applied model reduction using the  $\mathcal{H}_2$ -optimal interpolation points gen-

erated by the method of [16]. Then, for each  $r$ , out of the 2000 randomly generated shift selections, we counted the number of cases where the random shift selection yielded smaller values of  $\varsigma_{\min}(\mathbf{V}_r\mathbf{D})^{-1}$ ,  $\varsigma_{\min}(\mathbf{W}_r\mathbf{D})^{-1}$  and  $\|\Phi_r\|$ . The results are shown in Figure 2. Figure 2-(a) and -(b) show that for most of the cases, the  $\mathcal{H}_2$ -optimal interpolation points yield smaller values for  $\varsigma_{\min}(\mathbf{V}_r\mathbf{D})^{-1}$ ,  $\varsigma_{\min}(\mathbf{W}_r\mathbf{D})^{-1}$ . Indeed, for  $r \geq 48$ , the  $\mathcal{H}_2$ -optimal points produced smaller values in more than 99% of the cases. Also, for the last three cases:  $r = 66$ ,  $r = 68$ , and  $r = 70$ , the  $\mathcal{H}_2$ -optimal interpolation points always yielded smaller quantities. The results are even more dramatic for the projector norm, which is important in scaling the perturbation effects caused by inexact solves, see Theorem 4.2: Out of 70,000 cases (2000 selections for each  $r$  value), the  $\mathcal{H}_2$ -optimal interpolation point selection produced smaller condition numbers in all except 7 instances: 5 instances for  $r = 2$ , and 2 instances for  $r = 8$ . These numerical results illustrate that  $\mathcal{H}_2$ -optimal interpolation points can be expected to yield smaller values for  $\varsigma_{\min}(\mathbf{V}_r\mathbf{D})^{-1}$ ,  $\varsigma_{\min}(\mathbf{W}_r\mathbf{D})^{-1}$  and  $\|\Phi_r\|$ , and hence should produce reduced order models that are more robust with respect to perturbations.

Figure 2 also shows that for  $r = 14$ , 48% of the randomly selected shifts yielded smaller values of  $\varsigma_{\min}(\mathbf{V}_r\mathbf{D})^{-1}$ . However, when we inspected the 2000 randomly selected shift sets for  $r = 14$  in more detail, we observed some interesting additional features. We computed the three quantities  $\varsigma_{\min}(\mathbf{V}_r\mathbf{D})^{-1}$ ,  $\varsigma_{\min}(\mathbf{W}_r\mathbf{D})^{-1}$  and  $\|\Phi_r\|$  for each of the 2000 randomly selected shift sets, and compared them with the corresponding value derived from an  $\mathcal{H}_2$ -optimal shift selection. The results are shown in Figure 3. The top plot shows  $\varsigma_{\min}(\mathbf{V}_r\mathbf{D})/\varsigma_{\min}(\mathbf{V}_r^{\text{opt}}\mathbf{D})$  where  $\mathbf{V}_r^{\text{opt}}$  stands for the primitive interpolatory basis for the  $\mathcal{H}_2$ -optimal points. The bigger this ratio, the better the random shift selection. Even though for 48% of the cases, the random selection was better, the highest this ratio becomes is 2.20, i.e., the random shifts were never much better than a factor of 2 better than what  $\mathcal{H}_2$ -optimal shifts provided. For the remaining 52% of the cases, the randomly selected shifts were worse, and often worse by a factor of 100 or more. The situation for  $\mathbf{W}_r$  is shown in the middle plot. Once more, the situation is much more drastically in the favor of the  $\mathcal{H}_2$ -optimal interpolation points when the projector norm is inspected; the bottom plot in Figure 3 which depicts the ratio  $\|\Phi_r\|/\|\Phi_r^{\text{opt}}\|$  where  $\Phi_r^{\text{opt}}$  denotes the projector for the  $\mathcal{H}_2$ -optimal points. As illustrated in Figure 2, there are no random shift cases yielding a smaller projector norm. Furthermore, in many cases the projector norm for the random shift selection is almost 4 order of magnitudes higher than that of the

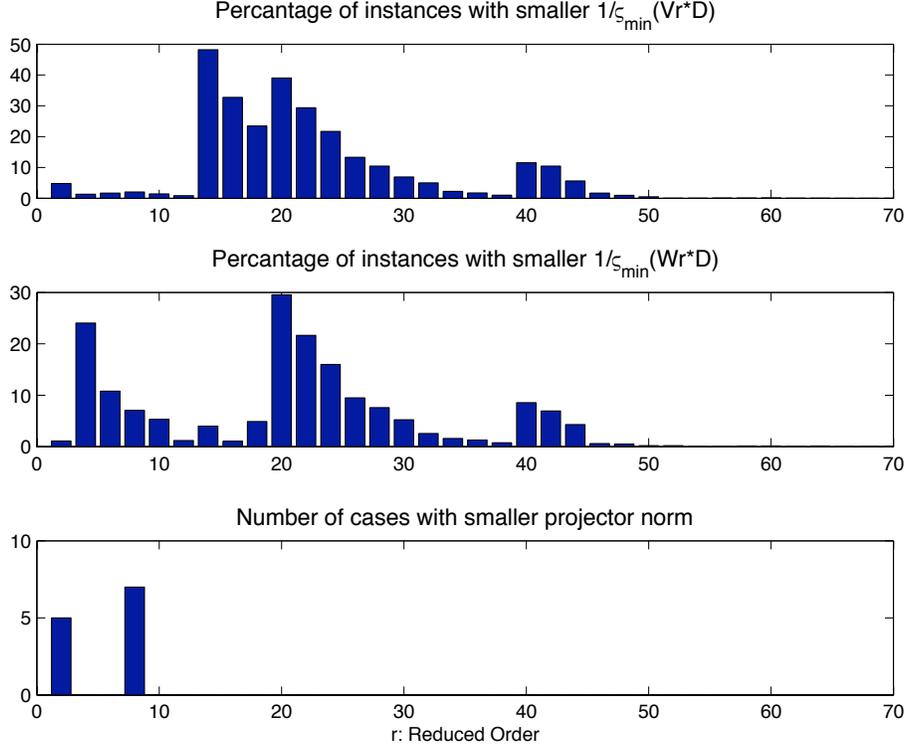


Figure 2: Comparison of  $\varsigma_{\min}(\mathbf{V}_r \mathbf{D})^{-1}$ ,  $\varsigma_{\min}(\mathbf{W}_r \mathbf{D})^{-1}$  and  $\|\Phi_r\|$  for random shift selections relative to values for  $\mathcal{H}_2$ -optimal shifts

$\mathcal{H}_2$ -optimal points. Indeed, on average the projector norm for the random points is  $8.19 \times 10^1$  times higher. These numbers change more in the favor of the  $\mathcal{H}_2$ -optimal points as  $r$  increases. For example, for  $r = 50$ , while the ratio  $\varsigma_{\min}(\mathbf{V}_r \mathbf{D})/\varsigma_{\min}(\mathbf{V}_r^{\text{opt}} \mathbf{D})$  becomes only as high as 1.48, it becomes as low as  $2.89 \times 10^{-4}$  for some random selections; Also, the ratio 1 can reach as high as  $2.91 \times 10^5$ . For  $r = 70$ ,  $\|\Phi_r\|$  for random selection is  $1.73 \times 10^2$  times higher than  $\|\Phi_r^{\text{opt}}\|$  on average. The three quantities we have been investigating appear to be extremely well conditioned for  $\mathcal{H}_2$ -optimal interpolation points. Even for  $r = 70$ , both  $\varsigma_{\min}(\mathbf{V}_r^{\text{opt}} \mathbf{D})^{-1}$ ,  $\varsigma_{\min}(\mathbf{W}_r^{\text{opt}} \mathbf{D})^{-1}$  remain smaller than 10 and  $\|\Phi_r^{\text{opt}}\|$  is smaller than 7.

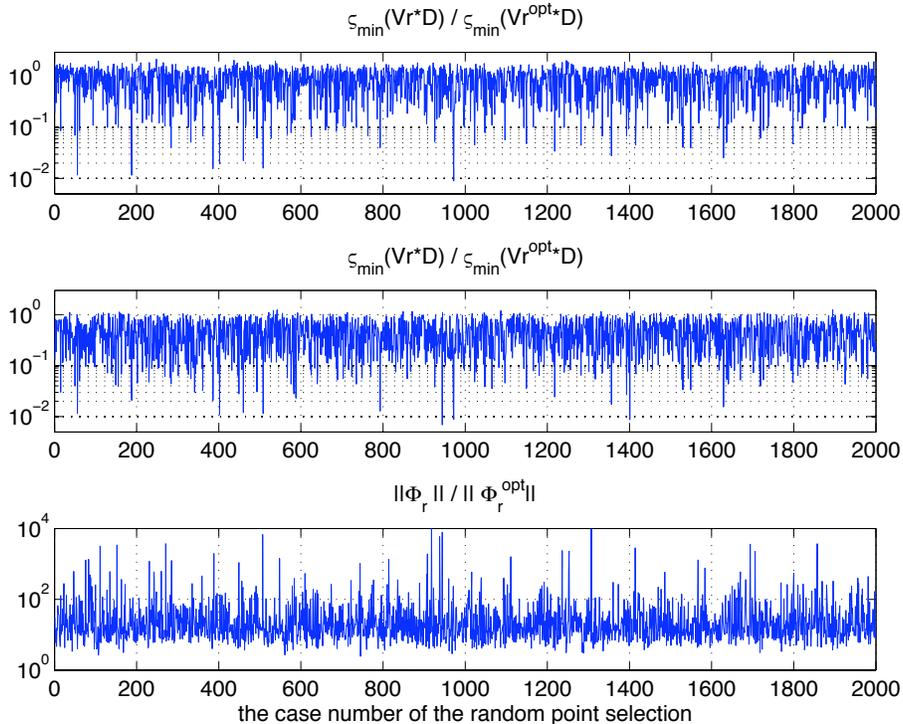


Figure 3: Detailed comparison for  $r = 14$

## 5. Inexact Solves in Optimal Interpolatory Approximation

The quality of the reduced-order model in interpolatory model reduction clearly depends on the selection of interpolation points and tangent directions. Until recently, this selection process was mostly *ad hoc*, and this factor had been the principal disadvantage of interpolatory model reduction. For systems in standard first-order state-space form, Gugercin *et al.* [16] have produced that an  $\mathcal{H}_2$ -optimal interpolation point / tangent direction selection strategy and proposed an Iterative Rational Krylov Algorithm (IRKA) to generate interpolatory reduced-order models that are (locally) optimal with respect to the  $\mathcal{H}_2$  norm. (An  $\mathcal{H}_2$ -optimal interpolation point selection strategy is still unknown for the general coprime factorization framework.) In this section, we investigate the behavior of inexact solves within the  $\mathcal{H}_2$ -optimal interpolatory approximation setting, specifically examining the behavior when inexact solves are employed in IRKA. In the rest of this section, we briefly review the optimal  $\mathcal{H}_2$  approximation problem and the method of

[16]. We then show how inexact solves can be employed effectively in this setting and discuss observed effects on optimality of the final reduced model. Our discussion focuses on systems in first-order descriptor form:

$$\mathcal{H}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B} \quad (49)$$

where  $\mathbf{E}, \mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times m}$  and  $\mathbf{C} \in \mathbb{R}^{p \times n}$ .

### 5.1. Optimal $\mathcal{H}_2$ approximation problem

Given the full-order system as in (49), the goal of the optimal  $\mathcal{H}_2$  model reduction problem is to find a reduced-order model  $\mathcal{H}_r(s)$  that minimizes the  $\mathcal{H}_2$  error; i.e.

$$\|\mathcal{H} - \mathcal{H}_r\|_{\mathcal{H}_2} = \min_{\substack{G_r \text{ stable} \\ \dim(G_r)=r}} \|\mathcal{H} - G_r\|_{\mathcal{H}_2}. \quad (50)$$

Many researchers have worked on this problem. These efforts can be grouped into two categories: Lyapunov-based optimal  $\mathcal{H}_2$  methods such as [31, 26, 17, 18, 30, 32]; and interpolation-based optimal  $\mathcal{H}_2$  methods such as [23, 16, 15, 29, 10, 14, 20, 5, 7]. Here, we will focus on the interpolation-based approach. However we note that Gugercin *et al.* [16] has shown that these two frameworks are theoretically equivalent; hence motivating the use of interpolatory approaches to optimal  $\mathcal{H}_2$  approximation since they are numerically superior to the Lyapunov-based approaches.

Since the optimization problem (50) is nonconvex, obtaining a global minimizer is a hard task and can be intractable. The usual approach is to find reduced order models that satisfy first-order necessary optimality conditions. Meier and Luenberger [23] introduced interpolation-based  $\mathcal{H}_2$ -optimality conditions for SISO systems. Analogous  $\mathcal{H}_2$ -optimality conditions for MIMO systems have recently been developed by [16, 10, 29] which in turn have led to analogous algorithms for the MIMO case; see [16, 10] for more details.

**Theorem 5.1.** *Given  $\mathcal{H}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}$ , let  $\mathcal{H}_r(s) = \sum_{i=1}^r \frac{1}{s-\hat{\lambda}_i} \mathbf{c}_i \mathbf{b}_i^T$  be the best  $r^{\text{th}}$  order approximation of  $\mathcal{H}$  with respect to the  $\mathcal{H}_2$  norm. Then*

$$\begin{aligned} (a) \quad & \mathcal{H}(-\hat{\lambda}_k) \mathbf{b}_k = \mathcal{H}_r(-\hat{\lambda}_k) \mathbf{b}_k, & (b) \quad & \mathbf{c}_k^T \mathcal{H}(-\hat{\lambda}_k) = \mathbf{c}_k^T \mathcal{H}_r(-\hat{\lambda}_k), & (51) \\ \text{and} \quad (c) \quad & \mathbf{c}_k^T \mathcal{H}'(-\hat{\lambda}_k) \mathbf{b}_k = \mathbf{c}_k^T \mathcal{H}'_r(-\hat{\lambda}_k) \mathbf{b}_k & & \text{for } k = 1, 2, \dots, r. \end{aligned}$$

### 5.1.1. An algorithm for interpolatory optimal $\mathcal{H}_2$ model reduction

Theorem 5.1 reveals that any  $\mathcal{H}_2$  optimal reduced-order model  $\mathfrak{H}_r(s)$  is a bi-tangential Hermite interpolant to  $\mathfrak{H}(s)$  at mirror images of the reduced-order poles. However, since the interpolation points and the tangent directions (and consequently,  $\mathbf{V}_r$  and  $\mathbf{W}_r$ ), depend on the final reduced-model to be computed, they are not known *a priori*. The *Iterative Rational Krylov Algorithm* (IRKA) of [16] resolves this problem by iteratively correcting the interpolation points and the directions as outlined in Algorithm 1: The reduced-order order poles are reflected across the imaginary axis to become the next set of interpolation points; the tangent directions are corrected using residue directions from the current reduced model. Upon convergence, the resulting interpolatory reduced-order model satisfies the necessary conditions of Theorem 5.1. For further details on IRKA, see [16].

#### Algorithm 1. IRKA for MIMO $\mathcal{H}_2$ Optimal Tangential Interpolation

1. Make an initial  $r$ -fold shift selection:  $\{\sigma_1, \dots, \sigma_r\}$  and initial tangent directions  $\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_r$  and  $\hat{\mathbf{c}}_1, \dots, \hat{\mathbf{c}}_r$ .
2.  $\mathbf{V}_r = \left[ (\sigma_1 \mathbf{E} - \mathbf{A})^{-1} \mathbf{B} \hat{\mathbf{b}}_1 \ \cdots \ (\sigma_r \mathbf{E} - \mathbf{A})^{-1} \mathbf{B} \hat{\mathbf{b}}_r \right]$   
 $\mathbf{W}_r = \left[ (\sigma_1 \mathbf{E} - \mathbf{A}^T)^{-1} \mathbf{C}^T \hat{\mathbf{c}}_1 \ \cdots \ (\sigma_r \mathbf{E} - \mathbf{A}^T)^{-1} \mathbf{C}^T \hat{\mathbf{c}}_1 \right]$ .
3. while (not converged)
  - (a)  $\mathbf{A}_r = \mathbf{W}_r^T \mathbf{A} \mathbf{V}_r$ ,  $\mathbf{E}_r = \mathbf{W}_r^T \mathbf{E} \mathbf{V}_r$ ,  $\mathbf{B}_r = \mathbf{W}_r^T \mathbf{B}$ , and  $\mathbf{C}_r = \mathbf{C} \mathbf{V}_r$
  - (b) Compute  $\mathbf{Y}^* \mathbf{A}_r \mathbf{X} = \text{diag}(\tilde{\lambda}_i)$  and  $\mathbf{Y}^* \mathbf{E}_r \mathbf{X} = \mathbf{I}_r$  where  $\mathbf{Y}^*$  and  $\mathbf{X}$  are the left and right eigenvector matrices for  $\lambda \mathbf{E}_r - \mathbf{A}_r$ .
  - (c)  $\sigma_i \leftarrow -\lambda_i(\mathbf{A}_r, \mathbf{E}_r)$  for  $i = 1, \dots, r$ ,  $\hat{\mathbf{b}}_i^* \leftarrow \mathbf{e}_i^T \mathbf{Y}^* \mathbf{B}_r$  and  $\hat{\mathbf{c}}_i \leftarrow \mathbf{C}_r \mathbf{X} \mathbf{e}_i$ .
  - (d)  $\mathbf{V}_r = \left[ (\sigma_1 \mathbf{E} - \mathbf{A})^{-1} \mathbf{B} \hat{\mathbf{b}}_1 \ \cdots \ (\sigma_r \mathbf{E} - \mathbf{A})^{-1} \mathbf{B} \hat{\mathbf{b}}_r \right]$
  - (e)  $\mathbf{W}_r = \left[ (\sigma_1 \mathbf{E} - \mathbf{A}^T)^{-1} \mathbf{C}^T \hat{\mathbf{c}}_1 \ \cdots \ (\sigma_r \mathbf{E} - \mathbf{A}^T)^{-1} \mathbf{C}^T \hat{\mathbf{c}}_1 \right]$ .
4.  $\mathbf{A}_r = \mathbf{W}_r^T \mathbf{A} \mathbf{V}_r$ ,  $\mathbf{E}_r = \mathbf{W}_r^T \mathbf{E} \mathbf{V}_r$ ,  $\mathbf{B}_r = \mathbf{W}_r^T \mathbf{B}$ ,  $\mathbf{C}_r = \mathbf{C} \mathbf{V}_r$

### 5.2. Inexact Iterative Rational Krylov Algorithm (InxIRKA)

For large system order, one may see from Algorithm 1, that the main cost of IRKA will generally be solving  $2r$  large linear systems at each step. If the IRKA iteration converges in  $k$  steps, a total of  $2rk$  linear systems will need to be solved. In settings where system dimension reaches into the millions, iterative linear system solvers become necessary and inexact linear system solves must be incorporated into IRKA. We refer to the modified algorithm as the *Inexact Iterative Rational Krylov Algorithm* (InxIRKA) and describe it

in Algorithm 2 below. We employ the Petrov-Galerkin framework for the inexact solves. In Algorithm 2, the function  $\mathbf{F}_{\text{PG}}$  in

$$[\tilde{\mathbf{v}}_i, \tilde{\mathbf{w}}_i] = \mathbf{F}_{\text{PG}}(\mathbf{A}, \mathbf{E}, \mathbf{B}, \sigma_i, \mathbf{b}_i, \mathbf{c}_i, \mathbf{v}^{(0)}, \mathbf{w}^{(0)}, \varepsilon)$$

denotes an inexact solve using a Petrov-Galerkin framework to approximately solve the linear systems  $(\sigma_i \mathbf{E} - \mathbf{A})\mathbf{v}_i = \mathbf{B}\mathbf{b}_i$  and  $(\sigma_i \mathbf{E} - \mathbf{A})^T \mathbf{w}_i = \mathbf{C}^T \mathbf{c}_i$  with initial guesses  $\mathbf{v}^{(0)}$  and  $\mathbf{w}^{(0)}$ , respectively, and a relative residual termination tolerance  $\varepsilon$ , i.e., at the end,

$$\frac{\|(\sigma_i \mathbf{E} - \mathbf{A})\tilde{\mathbf{v}}_i - \mathbf{B}\mathbf{b}_i\|}{\|\mathbf{B}\mathbf{b}_i\|} \leq \varepsilon \quad \text{and} \quad \frac{\|(\sigma_i \mathbf{E} - \mathbf{A})^T \tilde{\mathbf{w}}_i - \mathbf{C}^T \mathbf{c}_i\|}{\|\mathbf{C}^T \mathbf{c}_i\|} \leq \varepsilon$$

**Algorithm 2. InxIRKA for MIMO  $\mathcal{H}_2$  Optimal Tangential Interpolation**

1. *Make an initial  $r$ -fold shift selection:  $\{\sigma_1, \dots, \sigma_r\}$  and initial tangent directions  $\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_r$  and  $\hat{\mathbf{c}}_1, \dots, \hat{\mathbf{c}}_r$ .*
2.  $[\tilde{\mathbf{v}}_i, \tilde{\mathbf{w}}_i] = \mathbf{F}_{\text{PG}}(\mathbf{A}, \mathbf{E}, \mathbf{B}, \sigma_i, \mathbf{b}_i, \mathbf{c}_i, \mathbf{0}, \mathbf{0}, \varepsilon)$  for  $i = 1, \dots, r$
3.  $\tilde{\mathbf{V}}_r = [\tilde{\mathbf{v}}_1, \tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_r]$  and  $\tilde{\mathbf{W}}_r = [\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2, \dots, \tilde{\mathbf{w}}_r]$ .
4. *while (not converged)*
  - (a)  $\tilde{\mathbf{A}}_r = \tilde{\mathbf{W}}_r^T \mathbf{A} \tilde{\mathbf{V}}_r$ ,  $\tilde{\mathbf{E}}_r = \tilde{\mathbf{W}}_r^T \mathbf{E} \tilde{\mathbf{V}}_r$ ,  $\tilde{\mathbf{B}}_r = \tilde{\mathbf{W}}_r^T \mathbf{B}$ , and  $\tilde{\mathbf{C}}_r = \mathbf{C} \tilde{\mathbf{V}}_r$
  - (b) *Compute  $\mathbf{Y}^* \tilde{\mathbf{A}}_r \mathbf{X} = \text{diag}(\lambda_i)$  and  $\mathbf{Y}^* \tilde{\mathbf{E}}_r \mathbf{X} = \mathbf{I}_r$  where  $\mathbf{Y}^*$  and  $\mathbf{X}$  are the left and right eigenvector matrices of  $\lambda \tilde{\mathbf{E}}_r - \tilde{\mathbf{A}}_r$ .*
  - (c)  $\sigma_i \leftarrow -\lambda_i(\tilde{\mathbf{A}}_r, \tilde{\mathbf{E}}_r)$  for  $i = 1, \dots, r$ ,  $\hat{\mathbf{b}}_i^* \leftarrow \mathbf{e}_i^T \mathbf{Y}^* \tilde{\mathbf{B}}_r$  and  $\hat{\mathbf{c}}_i \leftarrow \tilde{\mathbf{C}}_r \mathbf{X} \mathbf{e}_i$ .
  - (d)  $[\tilde{\mathbf{v}}_i, \tilde{\mathbf{w}}_i] = \mathbf{F}_{\text{PG}}(\mathbf{A}, \mathbf{E}, \mathbf{B}, \sigma_i, \mathbf{b}_i, \mathbf{c}_i, \tilde{\mathbf{v}}_i, \tilde{\mathbf{w}}_i, \varepsilon)$  for  $i = 1, \dots, r$
  - (e)  $\tilde{\mathbf{V}}_r = [\tilde{\mathbf{v}}_1, \tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_r]$  and  $\tilde{\mathbf{W}}_r = [\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2, \dots, \tilde{\mathbf{w}}_r]$ .
5.  $\tilde{\mathbf{A}}_r = \tilde{\mathbf{W}}_r^T \mathbf{A} \tilde{\mathbf{V}}_r$ ,  $\tilde{\mathbf{E}}_r = \tilde{\mathbf{W}}_r^T \mathbf{E} \tilde{\mathbf{V}}_r$ ,  $\tilde{\mathbf{B}}_r = \tilde{\mathbf{W}}_r^T \mathbf{B}$ , and  $\tilde{\mathbf{C}}_r = \mathbf{C} \tilde{\mathbf{V}}_r$

As discussed and illustrated in [16, 3], in most cases IRKA converges rapidly; that is, the interpolation points and directions at the  $k^{\text{th}}$  step of IRKA stagnate rapidly with respect to  $k$ . Let  $\sigma_i^{(k)}$  and  $\mathbf{b}_i^{(k)}$  denote the  $i^{\text{th}}$  interpolation point and right-tangential direction, respectively, at the  $k^{\text{th}}$  step. Then we expect that as  $k$  increases, the solution  $\mathbf{v}_i^{(k)}$  of the linear system  $(\sigma_i^{(k)} \mathbf{E} - \mathbf{A})\mathbf{v}_i^{(k)} = \mathbf{B}\mathbf{b}_i^{(k)}$  from the  $k^{\text{th}}$  step approaches to the solution  $\mathbf{v}_i^{(k+1)}$  of the linear system  $(\sigma_i^{(k+1)} \mathbf{E} - \mathbf{A})\mathbf{v}_i^{(k+1)} = \mathbf{B}\mathbf{b}_i^{(k+1)}$  at the  $(k+1)^{\text{st}}$  step. This is precisely the reason that in Step 4.(d) of Algorithm 2, we use  $\mathbf{v}_i^{(k)}$  as an initial guess in solving  $(\sigma_i^{(k+1)} \mathbf{E} - \mathbf{A})\mathbf{v}_i^{(k+1)} = \mathbf{B}\mathbf{b}_i^{(k+1)}$  at the  $(k+1)^{\text{st}}$ . We

expect that this initialization strategy will speed-up the convergence of the iterative solves.

The development of effective stopping criteria based rationally on system theoretic error measures as we have introduced them here is the focus of on-going work. Similar approaches toward the design of effective preconditioning techniques and reuse of preconditioners tailored for interpolatory model reduction and especially for optimal  $\mathcal{H}_2$  approximation are also under investigation.

### 5.3. Effect of Inexact Solves in the InxIRKA Setting

The first question to answer in InxIRKA is whether a statement can be made about the optimality as in the exact IRKA case. Employing the Petrov-Galerkin framework makes this possible:

**Corollary 5.1.** *Let  $\tilde{\mathcal{H}}_r(s)$  be obtained by Algorithm 2. Then  $\tilde{\mathcal{H}}_r(s)$  satisfies the necessary conditions for optimal  $\mathcal{H}_2$  approximation of a near-by full-order model  $\tilde{\mathcal{H}}(s) = \mathbf{C}(s\mathbf{E} - (\mathbf{A} + \mathbf{F}_{2r}))^{-1}\mathbf{B}$  where  $\mathbf{F}_{2r}$  is the rank- $2r$  perturbation matrix defined in (45).*

Corollary 5.1 shows that with the help of the underlying Petrov-Galerkin framework, we state that the final reduced model of InxIRKA is an optimal  $\mathcal{H}_2$  approximation to a nearby full-order model.

As we discussed in Section 4.2, for a good selection of interpolation points, interpolatory model reduction is expected to be robust with respect to perturbations due to inexact solves. Hence, if one feeds the optimal interpolation points from IRKA into an inexact interpolation framework, we expect that the resulting reduced model will be close to the optimal reduced model of IRKA. However, the optimal interpolation points are not known initially and InxIRKA will be initiated with a nonoptimal initial shift selection. If the initial interpolation points and directions are poorly selected, at the early stages of the iteration, perturbations due to inexact solves might be magnified by this poor selection. One can avoid this scenario by using a small termination threshold  $\varepsilon$  in the early steps of InxIRKA, and then gradually increase  $\varepsilon$  as the iteration starts to converge. However, we note that in our numerical experiments using random initialization strategies, InxIRKA performed robustly and yielded high fidelity reduced models that are also close to the true optimal reduced model. This is illustrated in §5.4 below. Effective initialization strategies are discussed in [16] as well.

#### 5.4. Numerical results for $\text{InxIRKA}$

Here we illustrate the usage of inexact solves in the optimal  $\mathcal{H}_2$  approximation setting by comparing IRKA with  $\text{InxIRKA}$ . We use the example of §3.3, but with a finer discretization leading to a state-space dimension of  $n = 20209$ . We focus on a MIMO version using 2-inputs and 2-outputs.

We reduce the order to  $r = 6$  using both IRKA and  $\text{InxIRKA}$ . In  $\text{InxIRKA}$ , the dual linear systems are solved in a Petrov-Galerkin framework using BiCG [4] where we use three different values for the relative residual termination threshold of  $\varepsilon$ :  $10^{-5}$ ,  $10^{-3}$ , and  $10^{-1}$ . In all cases, the behavior of  $\text{InxIRKA}$  is virtually indistinguishable from that of IRKA. Starting with the same initial conditions, both IRKA and  $\text{InxIRKA}$  converge within 10 iteration steps in all 5 cases. The evolution of the  $\mathcal{H}_2$  errors  $\|\mathcal{H} - \mathcal{H}_r\|_{\mathcal{H}_2}$  and  $\|\mathcal{H} - \tilde{\mathcal{H}}_r\|_{\mathcal{H}_2}$  during the course of IRKA and  $\text{InxIRKA}$ , respectively, are depicted in the top plot of Figure 4. The figure shows that  $\text{InxIRKA}$  behavior is almost an exact replica of that of IRKA. The deviation from the exact IRKA is noticeable in the graph only for  $\varepsilon = 10^{-1}$ . To illustrate how much  $\mathcal{H}_r$  deviates from  $\tilde{\mathcal{H}}_r$  as IRKA and  $\text{InxIRKA}$  evolve, we show the progress of  $\|\mathcal{H}_r - \tilde{\mathcal{H}}_r\|_{\mathcal{H}_2}$  in the bottom plot of Figure 4. For this example, we initialized both IRKA and  $\text{InxIRKA}$  with an initial reduced-order model (as opposed to specifying initial interpolation points and tangent directions). Thus,  $\mathcal{H}_r = \tilde{\mathcal{H}}_r$  initially and no linear solvers are involved in the first ( $k = 0$ ) step. One could expect that perturbation errors due to inexact solves might accumulate over the course of the  $\text{InxIRKA}$  iteration, but this does not appear to be the case as this figure illustrates. The magnitude of  $\|\mathcal{H}_r - \tilde{\mathcal{H}}_r\|_{\mathcal{H}_2}$  remains relatively constant throughout the iteration at a magnitude proportional to the termination criterion.

The resulting  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  model reduction errors,  $\|\mathcal{H} - \mathcal{H}_r\|_{\mathcal{H}_2}$  and  $\|\mathcal{H} - \mathcal{H}_r\|_{\mathcal{H}_\infty}$  (with  $\mathcal{H}_r$  obtained from IRKA), versus  $\|\mathcal{H} - \tilde{\mathcal{H}}_r\|_{\mathcal{H}_2}$  and  $\|\mathcal{H} - \tilde{\mathcal{H}}_r\|_{\mathcal{H}_\infty}$  (with  $\tilde{\mathcal{H}}_r$  obtained from  $\text{InxIRKA}$ ) are given as  $\varepsilon$  varies in Table 4 below. The row corresponding to  $\varepsilon = 0$  represents the errors due to exact IRKA. These numbers demonstrate that employing inexact solves in  $\text{InxIRKA}$  does not degrade the model reduction performance. We also measure the difference between  $\mathcal{H}_r$  and  $\tilde{\mathcal{H}}_r$  in both  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  norms as  $\varepsilon$  varies. These results are tabulated in Table 5: Note that while  $\|\mathcal{H} - \mathcal{H}_r\|_{\mathcal{H}_2}$  and  $\|\mathcal{H} - \mathcal{H}_r\|_{\mathcal{H}_\infty}$  are respectively  $\mathcal{O}(10^{-4})$  and  $\mathcal{O}(10^{-2})$ , the contributions attributable to  $\mathcal{H}_r - \tilde{\mathcal{H}}_r$  are much smaller in magnitude and do not alter the resulting (optimal) model reduction performance in any significant way. If one were to convert the perturbation errors in Table 5 to relative

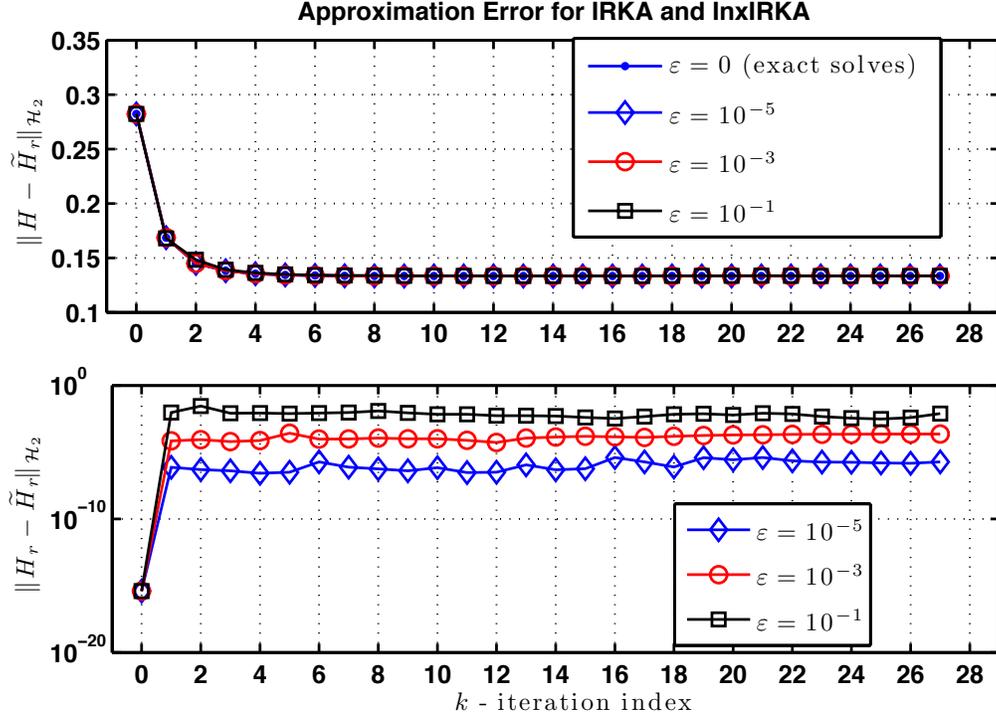


Figure 4: Evolution of the  $\mathcal{H}_2$  error during IRKA and lnxIRKA

error (as opposed to the displayed absolute error), both  $\|\mathcal{H}_r - \tilde{\mathcal{H}}_r\|_{\mathcal{H}_2}$  and  $\|\mathcal{H}_r - \tilde{\mathcal{H}}_r\|_{\mathcal{H}_\infty}$  starts at  $\mathcal{O}(10^{-6})$  for  $\varepsilon = 10^{-5}$ , and increases linearly by one order as  $\varepsilon$  increases by the same amount.

We finally list, in Table 6, the final exact and inexact optimal interpolation points due to IRKA, and lnxIRKA for  $\varepsilon = 10^{-3}$  and  $\varepsilon = 10^{-1}$ : Not surprisingly, the resulting interpolation points are very close to each other (though not the same). This can be viewed as another illustration of the fact that  $\mathcal{H}_r$  is an  $\mathcal{H}_2$  optimal approximation to a nearby full-order system.

As discussed above, in the implementation of lnxIRKA, we used the solution vectors from the previous step as the initial guess for the linear system in the next step taking advantage of the convergence in the interpolation points and tangent directions. To illustrate the effectiveness of this simple approach, throughout lnxIRKA we monitor the number of BiCG steps required to solve each linear system. We illustrate the behavior only for one of the interpolation points. We choose the interpolation points closest to the imaginary

$\varepsilon$	$\mathcal{H}_2$ error	$\mathcal{H}_\infty$ error
0	$3.708415753 \times 10^{-4}$	$1.084442854 \times 10^{-2}$
$10^{-5}$	$3.708415754 \times 10^{-4}$	$1.084425703 \times 10^{-2}$
$10^{-4}$	$3.708415778 \times 10^{-4}$	$1.084282001 \times 10^{-2}$
$10^{-3}$	$3.708418102 \times 10^{-4}$	$1.082437228 \times 10^{-2}$
$10^{-2}$	$3.708621743 \times 10^{-4}$	$1.064836300 \times 10^{-2}$
$10^{-1}$	$3.716780975 \times 10^{-4}$	$1.055441476 \times 10^{-2}$

Table 4: Evolution of the model reduction errors as  $\varepsilon$  varies

$\varepsilon$	$\ \mathcal{H}_r - \tilde{\mathcal{H}}_r\ _{\mathcal{H}_2}$	$\ \mathcal{H}_r - \tilde{\mathcal{H}}_r\ _{\mathcal{H}_\infty}$
$10^{-5}$	$5.1921 \times 10^{-9}$	$2.7776 \times 10^{-7}$
$10^{-4}$	$5.7156 \times 10^{-8}$	$2.4611 \times 10^{-6}$
$10^{-3}$	$6.3982 \times 10^{-7}$	$2.1043 \times 10^{-5}$
$10^{-2}$	$5.9277 \times 10^{-6}$	$2.0910 \times 10^{-4}$
$10^{-1}$	$2.2056 \times 10^{-5}$	$2.9228 \times 10^{-3}$

Table 5: Evolution of the perturbation error as  $\varepsilon$  varies

axis since these produce the hardest linear systems to solve and invariably contribute most to the cost of inexact solves. Figure 5 depicts the the number of BiCG steps required as InxIRKA proceeds for these interpolation points using three different stopping criteria  $\varepsilon = 10^{-5}$ ,  $\varepsilon = 10^{-3}$  and  $\varepsilon = 10^{-1}$ . The figure clearly illustrates that re-using the solutions from the previous steps works very effectively in reducing the overall cost of the BiCG. The number of BiCG steps goes from 1200 down to 200 in 3 to 4 steps.

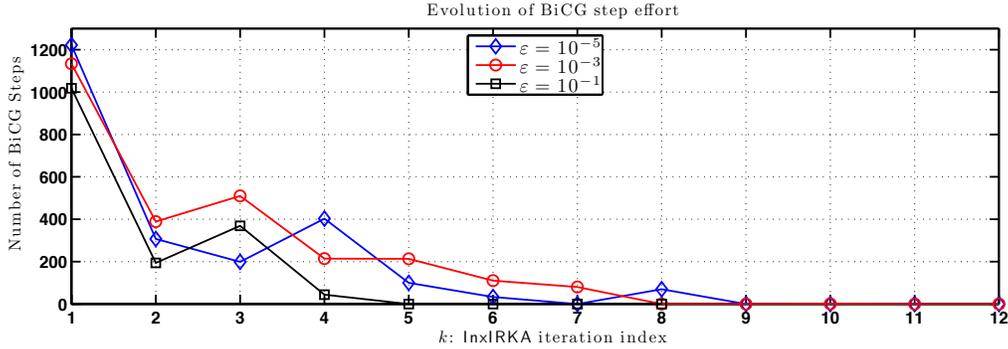


Figure 5: Evolution of BiCG effort during InxIRKA for shift closest to the imaginary axis

$\sigma_i(\text{IRKA})$	$\sigma_i(\text{lnxIRKA}), \varepsilon = 1 \times 10^{-3}$	$\sigma_i(\text{lnxIRKA}), \varepsilon = 1 \times 10^{-1}$
$1.0802 \times 10^{-5}$	$1.0800 \times 10^{-5}$	$1.2396 \times 10^{-5}$
$9.7164 \times 10^{-4}$	$9.7080 \times 10^{-4}$	$9.5860 \times 10^{-4}$
$6.6310 \times 10^{-3}$	$6.6246 \times 10^{-3}$	$6.5923 \times 10^{-3}$
$5.7925 \times 10^{-2}$	$5.7938 \times 10^{-2}$	$5.7929 \times 10^{-2}$
$9.0460 \times 10^{-1}$	$9.0419 \times 10^{-1}$	$8.9877 \times 10^{-1}$
$1.4127 \times 10^0$	$1.4126 \times 10^0$	$1.4104 \times 10^0$

Table 6: Optimal interpolations points as  $\varepsilon$  varies

## 6. Structure-preserving interpolation for descriptor systems

The backward error analysis of §4 has been presented for the transfer functions in the generalized coprime factorization form as in (2). In this section, we show that stronger conclusions on the structure of the reduced system can be drawn in the case the system has a realization as a descriptor system, that is,

$$\mathcal{H}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B} \quad (52)$$

where  $\mathbf{E}, \mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times m}$ , and  $\mathbf{C} \in \mathbb{R}^{p \times n}$  are constant matrices. In this case, for the interpolation points  $\{\sigma_j\}_{j=1}^r$ , and the tangent directions  $\{\mathbf{b}_j\}_{j=1}^r$  and  $\{\mathbf{c}_j\}_{j=1}^r$ , the associated primitive interpolatory bases  $\mathbf{V}_r$  and  $\mathbf{W}_r$  can be obtained from (13) and (14) using  $\mathcal{K}(s) = s\mathbf{E} - \mathbf{A}$ ,  $\mathcal{B}(s) = \mathbf{B}$  (constant matrix) and  $\mathcal{C}(s) = \mathbf{C}$  (constant matrix). Then, the resulting reduced-order model is given by

$$\mathcal{H}_r(s) = \mathbf{C}_r(s\mathbf{E}_r - \mathbf{A}_r)^{-1}\mathbf{B}_r \quad (53)$$

where

$$\mathbf{E}_r = \mathbf{W}_r^T \mathbf{E} \mathbf{V}_r, \quad \mathbf{A}_r = \mathbf{W}_r^T \mathbf{A} \mathbf{V}_r, \quad \mathbf{B}_r = \mathbf{W}_r^T \mathbf{B}, \quad \text{and} \quad \mathbf{C}_r = \mathbf{C} \mathbf{V}_r. \quad (54)$$

Let the set  $\mathcal{S} = \{\sigma_i, \mathbf{b}_i, \mathbf{c}_i\}$  denote given tangential interpolation data. Define the matrices  $\mathbb{L}[\mathcal{H}, \mathcal{S}] \in \mathbb{C}^{r \times r}$  and  $\mathbb{M}[\mathcal{H}, \mathcal{S}] \in \mathbb{C}^{r \times r}$  corresponding to the dynamical system  $\mathcal{H}(s)$  and interpolation data  $\mathcal{S}$ :

$$(\mathbb{L}[\mathcal{H}, \mathcal{S}])_{i,j} := \begin{cases} \frac{\mathbf{c}_i^T (\mathcal{H}(\sigma_i) - \mathcal{H}(\sigma_j)) \mathbf{b}_j}{\sigma_i - \sigma_j} & \text{if } i \neq j \\ \mathbf{c}_i^T \mathcal{H}'(\sigma_i) \mathbf{b}_i & \text{if } i = j \end{cases} \quad (55)$$

$$(\mathbb{M}[\mathcal{H}, \mathcal{S}])_{i,j} := \begin{cases} \frac{\mathbf{c}_i^T (\sigma_i \mathcal{H}(\sigma_i) - \sigma_j \mathcal{H}(\sigma_j)) \mathbf{b}_j}{\sigma_i - \sigma_j} & \text{if } i \neq j \\ \mathbf{c}_i^T [s\mathcal{H}(s)]'|_{s=\sigma_i} \mathbf{b}_i & \text{if } i = j \end{cases} \quad (56)$$

$\mathbb{L}[\mathcal{H}, \mathcal{S}]$  is the *Loewner matrix* associated with the interpolation data  $\mathcal{S}$  and the dynamical system  $\mathcal{H}(s)$ ,  $\mathbb{M}[\mathcal{H}, \mathcal{S}]$  is the *shifted Loewner matrix* associated with the interpolation data  $\mathcal{S}$  and the system  $s\mathcal{H}(s)$ , see [3, 22]. The next theorem presents a canonical structure for the exact interpolatory reduced-order model (53)-(54).

**Theorem 6.1.** [22] *Given a full-order model  $\mathcal{H}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}$  and tangential interpolation data  $\mathcal{S} = \{\sigma_i, \mathbf{b}_i, \mathbf{c}_i\}$ , then the reduced-order quantities in (54) satisfy*

$$\begin{aligned} \mathbf{E}_r &= -\mathbb{L}[\mathcal{H}, \mathcal{S}], \\ \mathbf{A}_r &= -\mathbb{M}[\mathcal{H}, \mathcal{S}], \end{aligned} \quad \mathbf{B}_r = \begin{bmatrix} \mathbf{c}_1^T \mathcal{H}(\sigma_1) \\ \vdots \\ \mathbf{c}_r^T \mathcal{H}(\sigma_r) \end{bmatrix}, \quad (57)$$

$$\mathbf{C}_r = [ \mathcal{H}(\sigma_1)\mathbf{b}_1, \dots, \mathcal{H}(\sigma_r)\mathbf{b}_r ].$$

### 6.1. The Petrov-Galerkin framework and structure preservation

Theorem 6.1 presents a canonical form for the exact bitangential Hermite interpolant in the case of standard state-space model. Next we show that if a Petrov-Galerkin framework is employed in the solution of the linear systems, the inexact reduced-model will have exactly the same form as the exact one. The result is a direct consequence of Theorems 4.1 and 6.1.

**Corollary 6.1.** *Given the standard full-order model  $\mathcal{H}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}$  together with the the interpolation data  $\mathcal{S} = \{\sigma_i, \mathbf{b}_i, \mathbf{c}_i\}$ , let the inexact solutions  $\tilde{\mathbf{v}}_j$  for  $(\sigma_j\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}\mathbf{b}_j$  and  $\tilde{\mathbf{w}}_j$  for  $(\sigma_j\mathbf{E} - \mathbf{A})^{-T}\mathbf{C}^T\mathbf{c}_j$  be obtained in a Petrov-Galerkin framework as in (42). Let  $\tilde{\mathbf{V}}_r$  and  $\tilde{\mathbf{W}}_r$  denote the corresponding inexact Krylov bases as in (43). Define the residuals*

$$\tilde{\boldsymbol{\eta}}_j = (\sigma_j\mathbf{E} - \mathbf{A})\tilde{\mathbf{v}}_j - \mathbf{B}\mathbf{b}_j \quad \text{and} \quad \tilde{\boldsymbol{\xi}}_j = (\sigma_j\mathbf{E} - \mathbf{A})^T\tilde{\mathbf{w}}_j - \mathbf{C}^T\mathbf{c}_j.$$

Let the residual matrices  $\mathbf{R}_b$  and  $\mathbf{R}_c$ , and the rank  $2r$  matrix  $\mathbf{F}_{2r}$  be as defined in (44) and (45), respectively. Then, the inexact interpolatory reduced-order model

$$\tilde{\mathcal{H}}_r(s) = \tilde{\mathbf{C}}_r(s\tilde{\mathbf{E}}_r - \tilde{\mathbf{A}}_r)^{-1}\tilde{\mathbf{B}}_r \quad (58)$$

is an exact Hermite bitangential interpolant for the perturbed full-order model

$$\tilde{\mathcal{H}}(s) = \mathbf{C}(s\mathbf{E} - (\mathbf{A} + \mathbf{F}_{2r}))^{-1}\mathbf{B}. \quad (59)$$

Moreover, the reduced-order quantities satisfy

$$\begin{aligned} \tilde{\mathbf{E}}_r &= -\mathbb{L}[\tilde{\mathcal{H}}, \mathcal{S}], & \tilde{\mathbf{B}}_r &= \begin{bmatrix} \mathbf{c}_1^T \tilde{\mathcal{H}}(\sigma_1) \\ \vdots \\ \mathbf{c}_r^T \tilde{\mathcal{H}}(\sigma_r) \end{bmatrix}, & \tilde{\mathbf{C}}_r &= [\tilde{\mathcal{H}}(\sigma_1)\mathbf{b}_1, \dots, \tilde{\mathcal{H}}(\sigma_r)\mathbf{b}_r]. \end{aligned} \quad (60)$$

where  $\mathbb{L}[\tilde{\mathcal{H}}, \mathcal{S}]$  and  $\mathbb{M}[\tilde{\mathcal{H}}, \mathcal{S}]$  are the Loewner matrices associated with the dynamical systems  $\tilde{\mathcal{H}}(s)$  and  $s\tilde{\mathcal{H}}(s)$  respectively, and the interpolation data  $\mathcal{S}$  as defined in (55) and (56).

Corollary 6.1 reveals that the inexact reduced-order model quantities have exactly the same structure as their exact counterparts. The interpolation data  $\mathcal{S}$  is the same in both cases; the only difference is that  $\mathcal{H}(s)$  is replaced by  $\tilde{\mathcal{H}}(s)$  in the construction that yields the Loewner-matrix structure. The preservation of this structure is independent of the accuracy to which the linear systems are solved. In the case where  $\mathbf{E} = \mathbf{I}$ , the structure of the exact and inexact reduced-models becomes even simpler:

**Corollary 6.2.** *Assume the hypotheses of Theorem 6.1 with  $\mathbf{E} = \mathbf{I}$ . Then the exact interpolant  $\mathcal{H}_r(s) = \mathbf{C}_r(s\mathbf{I}_r - \mathbf{A}_r)^{-1}\mathbf{B}_r$  satisfies*

$$\mathbf{A}_r = \Sigma - \mathbf{Q}\mathbf{B}, \quad \mathbf{B}_r = \mathbf{Q}, \quad \text{and} \quad \mathbf{C}_r = [\mathbf{H}(\sigma_1)\mathbf{b}_1, \dots, \mathbf{H}(\sigma_r)\mathbf{b}_r] \quad (61)$$

where

$$\mathbf{Q} = (\mathbf{W}_r^T \mathbf{V}_r)^{-1} \mathbf{W}_r^T \mathbf{B}, \quad \Sigma = \text{diag}(\sigma_1, \dots, \sigma_r) \quad \text{and} \quad \mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_r]. \quad (62)$$

Assume the hypotheses of Corollary 6.1 with  $\mathbf{E} = \mathbf{I}$ . Then, the inexact interpolant  $\tilde{\mathcal{H}}_r(s) = \tilde{\mathbf{C}}_r(s\mathbf{I}_r - \tilde{\mathbf{A}}_r)^{-1}\tilde{\mathbf{B}}_r$  satisfies

$$\tilde{\mathbf{A}}_r = \Sigma - \tilde{\mathbf{Q}}\mathbf{B}, \quad \tilde{\mathbf{B}}_r = \tilde{\mathbf{Q}}, \quad \text{and} \quad \tilde{\mathbf{C}}_r = [\tilde{\mathcal{H}}(\sigma_1)\mathbf{b}_1, \dots, \tilde{\mathcal{H}}(\sigma_r)\mathbf{b}_r] \quad (63)$$

where

$$\tilde{\mathbf{Q}} = (\tilde{\mathbf{W}}_r^T \tilde{\mathbf{V}}_r)^{-1} \tilde{\mathbf{W}}_r^T \mathbf{B}, \quad (64)$$

$\tilde{\mathcal{H}}(s)$  is the perturbed full-order model as in (59) with  $\mathbf{E} = \mathbf{I}$ , and  $\Sigma$  and  $\mathbf{B}$  are as defined in (62).

Corollary 6.2 illustrates that in the case of  $\mathbf{E} = \mathbf{I}$ , both of the reduced system matrices,  $\mathbf{A}_r$  and  $\tilde{\mathbf{A}}_r$ , are perturbations of rank  $\min(r, m, p)$  to the diagonal matrix of interpolation points,  $\Sigma$ .

## References

- [1] Kapil Ahuja. Recycling Bi-Lanczos algorithms: BiCG, CGS, BiCGSTAB. Master's thesis, Virginia Tech, Blacksburg, Virginia, August 2009.
- [2] A.C. Antoulas. *Approximation of Large-Scale Dynamical Systems (Advances in Design and Control)*. Society for Industrial and Applied Mathematics Philadelphia, PA, USA, 2005.
- [3] A.C. Antoulas, C.A. Beattie, and S. Gugercin. Interpolatory model reduction of large-scale dynamical systems. In J. Mohammadpour and K. Grigoriadis, editors, *Efficient Modeling and Control of Large-Scale Systems*. Springer-Verlag, 2010.
- [4] R. Barrett, M. Berry, T.F. Chan, J. Demmel, J.M. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. Van der Vorst. *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*. Society for Industrial Mathematics, 1994.
- [5] C.A. Beattie and S. Gugercin. Krylov-based minimization for optimal  $\mathcal{H}_2$  model reduction. *46th IEEE Conference on Decision and Control*, pages 4385–4390, Dec. 2007.
- [6] C.A. Beattie and S. Gugercin. Interpolatory projection methods for structure-preserving model reduction. *Systems & Control Letters*, 58(3):225–232, 2009.
- [7] C.A. Beattie and S. Gugercin. A trust region method for optimal  $\mathcal{H}_2$  model reduction. *48th IEEE Conference on Decision and Control*, Dec. 2009.
- [8] P. Benner. Solving large-scale control problems. *Control Systems Magazine, IEEE*, 24(1):44–59, 2004.
- [9] P. Benner and J. Saak. Efficient numerical solution of the LQR-problem for the heat equation. *Proc. Appl. Math. Mech*, 4(1):648–649, 2004.
- [10] A. Bunse-Gerstner, D. Kubalinska, G. Vossen, and D. Wilczek.  $\mathcal{H}_2$ -optimal model reduction for large scale discrete dynamical MIMO systems. *Journal of Computational and Applied Mathematics*, 2009. doi:10.1016/j.cam.2008.12.029.

- [11] K. Gallivan, A. Vandendorpe, and P. Van Dooren. Model reduction via truncation: an interpolation point of view. *Linear Algebra and Its Applications*, 375:115–134, 2003.
- [12] K. Glover. All optimal Hankel-norm approximations of linear multi-variable systems and their  $L^\infty$ -error bounds. *International Journal of Control*, 39(6):1115–1193, 1984.
- [13] E. Grimme. *Krylov Projection Methods for Model Reduction*. PhD thesis, Coordinated-Science Laboratory, University of Illinois at Urbana-Champaign, 1997.
- [14] S. Gugercin. An iterative rational Krylov algorithm (IRKA) for optimal  $\mathcal{H}_2$  model reduction. In *Householder Symposium XVI*, Seven Springs Mountain Resort, PA, USA, May 2005.
- [15] S. Gugercin, A.C. Antoulas, and C.A. Beattie. A rational Krylov iteration for optimal  $\mathcal{H}_2$  model reduction. In *Proceedings of MTNS*, volume 2006, 2006.
- [16] S. Gugercin, A.C. Antoulas, and C.A. Beattie.  $\mathcal{H}_2$  model reduction for large-scale linear dynamical systems. *SIAM Journal on Matrix Analysis and Applications*, 30(2):609–638, 2008.
- [17] Y. Halevi. Frequency weighted model reduction via optimal projection. *Automatic Control, IEEE Transactions on*, 37(10):1537–1542, 1992.
- [18] D. Hyland and D. Bernstein. The optimal projection equations for model reduction and the relationships among the methods of Wilson, Skelton, and Moore. *Automatic Control, IEEE Transactions on*, 30(12):1201–1211, 1985.
- [19] J.G. Korvink and E.B. Rudnyi. Oberwolfach benchmark collection. In *Dimension reduction of large-scale systems: proceedings of a workshop held in Oberwolfach, Germany, October 19-25, 2003*, page 311. Springer Verlag, 2005.
- [20] D. Kubalinska, A. Bunse-Gerstner, G. Vossen, and D. Wilczek.  $\mathcal{H}_2$ -optimal interpolation based model reduction for large-scale systems. In *Proceedings of the 16<sup>th</sup> International Conference on System Science*, Poland, 2007.

- [21] Y. Liu and B.D.O. Anderson. Singular perturbation approximation of balanced systems. *International Journal of Control*, 50(4):1379–1405, 1989.
- [22] A.J. Mayo and A.C. Antoulas. A framework for the solution of the generalized realization problem. *Linear Algebra and Its Applications*, 425(2-3):634–662, 2007.
- [23] L. Meier III and D. Luenberger. Approximation of linear constant systems. *Automatic Control, IEEE Transactions on*, 12(5):585–588, 1967.
- [24] B. Moore. Principal component analysis in linear systems: Controllability, observability, and model reduction. *Automatic Control, IEEE Transactions on*, 26(1):17–32, 1981.
- [25] C. Mullis and R. Roberts. Synthesis of minimum roundoff noise fixed point digital filters. *Circuits and Systems, IEEE Transactions on*, 23(9):551–562, 1976.
- [26] J.T. Spanos, M.H. Milman, and D.L. Mingori. A new algorithm for  $L_2$  optimal model reduction. *Automatica (Journal of IFAC)*, 28(5):897–909, 1992.
- [27] D.B. Szyld. The many proofs of an identity on the norm of oblique projections. *Numerical Algorithms*, 42(3):309–323, 2006.
- [28] A. van der Sluis. Condition numbers and equilibration of matrices. *Numerische Mathematik*, 14(1):14–23, 1969.
- [29] P. van Dooren, K.A. Gallivan, and P.A. Absil.  $\mathcal{H}_2$ -optimal model reduction of MIMO systems. *Applied Mathematics Letters*, 2008.
- [30] DA Wilson. Optimum solution of model-reduction problem. *Proc. IEE*, 117(6):1161–1165, 1970.
- [31] W.Y. Yan and J. Lam. An approximate approach to  $\mathcal{H}_2$  optimal model reduction. *Automatic Control, IEEE Transactions on*, 44(7):1341–1358, 1999.
- [32] D. Zigic, LT Watson, and C. Beattie. Contragredient transformations applied to the optimal projection equations. *Linear Algebra and Its Applications*, 188:665–676, 1993.