

\mathcal{H}_2 -QUASI-OPTIMAL MODEL ORDER REDUCTION FOR QUADRATIC-BILINEAR CONTROL SYSTEMS*

PETER BENNER[†], PAWAN GOYAL[‡], AND SERKAN GUGERCIN[§]

Abstract. We investigate the optimal model reduction problem for large-scale quadratic-bilinear (QB) control systems. Our contributions are threefold. First, we discuss the variational analysis and the Volterra series formulation for QB systems. We then define the \mathcal{H}_2 -norm for a QB system based on the *kernels* of the underlying Volterra series and also propose a truncated \mathcal{H}_2 -norm. Next, we derive first-order necessary conditions for an optimal approximation, where optimality is measured in term of the truncated \mathcal{H}_2 -norm of the error system. We then propose an iterative model reduction algorithm, which upon convergence yields a reduced-order system that *approximately* satisfies the newly derived optimality conditions. We also discuss an efficient computation of the reduced Hessian, using the special Kronecker structure of the Hessian of the system. We illustrate the efficiency of the proposed method by means of several numerical examples resulting from semi-discretized nonlinear partial differential equations and show its competitiveness with the existing model reduction schemes for QB systems such as moment-matching methods and balanced truncation.

Key words. Model order reduction, quadratic-bilinear control systems, tensor calculus, \mathcal{H}_2 -norm, Sylvester equations, nonlinear partial differential equations.

AMS subject classifications. 15A69, 34C20, 41A05, 49M05, 93A15, 93C10, 93C15.

1. Introduction. Numerical simulation is a fundamental tool in the analysis of dynamical systems, and is required repeatedly in the computational control design, optimization and prediction thereof. Dynamical systems are generally governed by partial differential equations (PDEs), or ordinary differential equations (ODEs), or a combination of both. A high-fidelity approximation of the underlying physical phenomena requires a finely discretized mesh over the interesting spatial domain, leading to complex dynamical systems with a high-dimensional state-space. The simulations of such large-scale systems, however, impose a huge computational burden. This inspires *model order reduction* (MOR), which aims at constructing simple and reliable surrogate models such that their input-output behaviors approximate that of the original large-scale system accurately. These surrogate models can then be used in engineering studies, which make numerical simulations faster and efficient.

In recent decades, numerous theoretical and computational aspects for MOR of linear systems have been developed; e.g., see [1, 6, 14, 41]. These methods have been successfully applied in various fields, e.g., optimal control and PDE constrained optimization and uncertainty quantification; for example, see [16, 28]. In recent years, however, MOR of nonlinear systems has gained significant interest with a goal of extending the input-independent, optimal MOR techniques from linear systems to nonlinear ones. For example, MOR techniques for linear systems such as balanced truncation [1, 35], or the iterative rational Krylov method (IRKA) [23], have been extended to a special class of nonlinear systems, so-called *bilinear systems*, in which nonlinearity arises from the product of the state and input [7, 10, 20, 43]. However,

*Submitted to the editors May 17, 2017.

[†]Max Planck Institute for Dynamics of Complex Technical Systems, Sandtorstraße 1, 39106 Magdeburg, Germany; (benner@mpi-magdeburg.mpg.de).

[‡]Corresponding author. Max Planck Institute for Dynamics of Complex Technical Systems, Sandtorstraße 1, 39106 Magdeburg, Germany; (goyalp@mpi-magdeburg.mpg.de).

[§]Department of Mathematics, Virginia Tech, Blacksburg, VA 24061-0123, United States; (gugercin@math.vt.edu).

there are numerous open and difficult problems for other important classes of nonlinear systems. In this article, we address another vital class of nonlinear systems, called *quadratic-bilinear* (QB) systems. These are of the form:

$$(1.1) \quad \Sigma : \begin{cases} \dot{x}(t) = Ax(t) + H(x(t) \otimes x(t)) + \sum_{k=1}^m N_k x(t) u_k(t) + Bu(t), \\ y(t) = Cx(t), \quad x(0) = 0, \end{cases}$$

where $x(t) \in \mathbb{R}^n$, $u(t) \in \mathbb{R}^m$ and $y(t) \in \mathbb{R}^p$ are the states, inputs, and outputs of the systems, respectively, u_k is the k th component of u , and n is the state dimension which is generally very large. Furthermore, $A, N_k \in \mathbb{R}^{n \times n}$ for $k \in \{1, \dots, m\}$, $H \in \mathbb{R}^{n \times n^2}$, $B \in \mathbb{R}^{n \times m}$, and $C \in \mathbb{R}^{p \times n}$.

There is a variety of applications where the system inherently contains a quadratic nonlinearity, which can be modeled in the QB form (1.1), e.g., Burgers' equation. Moreover, a large class of smooth nonlinear systems, involving combinations of elementary functions like exponential, trigonometric, and polynomial functions, etc. can be equivalently rewritten as QB systems (1.1) as shown in [9, 22]. This is achieved by introducing some new appropriate state variables to simplify the nonlinearities present in the underlying control system and deriving differential equations, corresponding to the newly introduced variables, or by using appropriate algebraic constraints. When algebraic constraints are introduced in terms of the state and the newly defined variables, the system contains algebraic equations along with differential equations. Such systems are called *differential-algebraic equations* (DAEs) or *descriptor systems* [33]. MOR procedures for DAEs become inevitably more complicated, even in the linear and bilinear settings; e.g., see [11, 24]. In this article, we restrict ourselves to QB ODE systems and leave MOR for QB descriptor systems as a future research topic.

For a given QB system (1.1) Σ of dimension n , our aim is to construct a reduced-order system

$$(1.2) \quad \hat{\Sigma} : \begin{cases} \hat{\dot{x}}(t) = \hat{A}\hat{x}(t) + \hat{H}(\hat{x}(t) \otimes \hat{x}(t)) + \sum_{k=1}^m \hat{N}_k \hat{x}(t) u_k(t) + \hat{B}u(t), \\ \hat{y}(t) = \hat{C}\hat{x}(t), \quad \hat{x}(0) = 0, \end{cases}$$

where $\hat{A}, \hat{N}_k \in \mathbb{R}^{r \times r}$ for $k \in \{1, \dots, m\}$, $\hat{H} \in \mathbb{R}^{r \times r^2}$, $\hat{B} \in \mathbb{R}^{r \times m}$, and $\hat{C} \in \mathbb{R}^{p \times r}$ with $r \ll n$ such that the outputs of the system (1.1) and (1.2), y and \hat{y} , are very well approximated in a proper norm for all admissible inputs $u_k \in L^2[0, \infty[$.

Similar to the linear and bilinear cases, we construct the reduced-order system (1.2) via projection. Towards this goal, we construct two model reduction basis matrices $V, W \in \mathbb{R}^{n \times r}$ such that $W^T V$ is invertible. Then, the reduced matrices in (1.2) are given by

$$\begin{aligned} \hat{A} &= (W^T V)^{-1} W^T A V, & \hat{N}_k &= (W^T V)^{-1} W^T N_k V, \quad \text{for } k \in \{1, \dots, m\}, \\ \hat{H} &= (W^T V)^{-1} W^T H (V \otimes V), & \hat{B} &= (W^T V)^{-1} W^T B, \quad \text{and } \hat{C} = C V. \end{aligned}$$

It can be easily seen that the quality of the reduced-order system depends on the choice of the reduction subspaces $\text{Im } V$ and $\text{Im } W$. There exist various MOR approaches in the literature to determine these subspaces. One of the earlier and popular methods for nonlinear systems is trajectory-based MOR such as proper orthogonal decomposition (POD); e.g., see [2, 18, 27, 32]. This relies on the Galerkin projection

$\mathcal{P} = \mathcal{V}\mathcal{V}^T$, where \mathcal{V} is determined based on the dominate modes of the system dynamics. For efficient computations of nonlinear parts in the system, some advanced methodologies can be employed such as the empirical interpolation method (EIM), the best points interpolation method, the discrete empirical interpolation method (DEIM), see, e.g. [2, 4, 21, 18, 36]. Another widely used method for nonlinear system is the trajectory piecewise linear (TPWL) method, e.g., see [39]. For this method, the nonlinear system is replaced by a weighted sum of linear systems; these linear systems can then be reduced by using well-known methods for linear systems such as balanced truncation, or interpolation methods, e.g., see [1]. However, the above mentioned methods require some snapshots or solution trajectories of the original systems for particular inputs. This indicates that the resulting reduced-order system depends on the choice of inputs, which may make the reduced-order system inadequate in many applications such as control and optimization, where the variation of the input is inherent to the problem.

MOR methods, based on interpolation or moment-matching, have been extended from linear systems to QB systems, with the aim of capturing the input-output behavior of the underlying system independent of a training input. One-sided interpolatory projection for QB systems is studied in, e.g., [3, 22, 37, 38], and has been recently extended to a two-sided interpolatory projection in [8, 9]. These methods result in reduced systems that do not rely on the training data for a control input; see also the survey [6] for some related approaches. Thus, the determined reduced systems can be used in input-varying applications. In the aforementioned references, the authors have shown how to construct an interpolating reduced-order system for a given set of interpolation points. But it is still an open problem how to choose these interpolation points optimally with respect to an appropriate norm. Furthermore, the two-sided interpolatory projection method [9] is only applicable to single-input single-output systems, which is very restrictive, and additionally, the stability of the resulting reduced-order systems also remains another major issue.

Very recently, balanced truncation has been extended from linear/bilinear systems to QB systems [12]. This method first determines the states which are hard to control and observe and constructs the reduced model by truncating those states. Importantly, it results in locally Lyapunov stable reduced-order systems, and an appropriate order of the reduced system can be determined based on the singular values of the Gramians of the system. But unlike in the linear case the resulting reduced systems do not guarantee other desirable properties such as an a priori error bound. Moreover, in order to apply balanced truncation to QB systems, we require the solutions of four conventional Lyapunov equations, which could be computationally cumbersome in large-scale settings; though there have been many advancements in recent times related to computing the low-rank solutions of Lyapunov equations [15, 42].

In this paper, we study the \mathcal{H}_2 -optimal approximation problem for QB systems. Precisely, we show how to choose the model reduction bases in a two-sided projection interpolation framework for QB systems so that the reduced model is a local minimizer in an appropriate norm. Our main contributions are threefold. In Section 2, we derive various expressions and formulas related to Kronecker products, which are later heavily utilized in deriving optimality conditions. In Section 3, we first define the \mathcal{H}_2 -norm of the QB system (1.1) based on the *kernels* of its Volterra series (input/output mapping), and then derive an expression for a truncated \mathcal{H}_2 -norm for QB systems as well. Subsequently, based on a truncated \mathcal{H}_2 -norm of the error system, we derive first-order necessary conditions for optimal model reduction of QB systems. We then propose an iterative algorithm to construct reduced models that *approxi-*

mately satisfy the newly derived optimality conditions. Furthermore, we discuss an efficient alternative way to compute reduced Hessians as compared to the one proposed in [9]. In Section 4, we illustrate the efficiency of the proposed method for various semi-discretized nonlinear PDEs and compare it with existing methods such as balanced truncation [12] as well as the one-sided and two-sided interpolatory projection methods for QB systems [9, 22]. We conclude the paper with a short summary and potential future directions in Section 5.

Notation: Throughout the paper, we make use of the following notation:

- I_q denotes the identity matrix of size $q \times q$, and its p th column is denoted by e_p^q .
- $\text{vec}(\cdot)$ denotes vectorization of a matrix, and \mathcal{I}_m denotes $\text{vec}(I_m)$.
- $\text{tr}(\cdot)$ refers to the trace of a matrix.
- Using MATLAB[®] notation, we denote the j th column of the matrix A by $A(:, j)$.
- $\mathbf{0}$ is a zero matrix of appropriate size.
- We denote the full-order system (1.1) and reduced-order system (1.2) by Σ and $\widehat{\Sigma}$, respectively.

2. Tensor Matricizations and their Properties. We first review some basic concepts from tensor algebra. First, we note the following important properties of the $\text{vec}(\cdot)$ operator:

$$(2.1a) \quad \text{tr}(X^T Y) = \text{vec}(X)^T \text{vec}(Y) = \text{vec}(Y)^T \text{vec}(X), \text{ and}$$

$$(2.1b) \quad \text{vec}(XYZ) = (Z^T \otimes X) \text{vec}(Y).$$

Next, we review the concepts of matricization of a tensor. Since the Hessian H of the QB system in (1.1) is closely related to a 3rd order tensor, we focus on 3rd order tensors $\mathcal{X}^{n \times n \times n}$. However, most of the concepts can be extended to a general k th order tensor. Similar to how rows and columns are defined for a matrix, one can define a fiber of \mathcal{X} by fixing all indices except for one, e.g., $\mathcal{X}(:, i, j)$, $\mathcal{X}(i, :, j)$ and $\mathcal{X}(i, j, :)$. Mathematical operations involving tensors are easier to perform with its corresponding matrices. Therefore, there exists a very well-known process of unfolding a tensor into a matrix, called *matricization* of a tensor. For a 3-dimensional tensor, there are 3 different ways to unfold the tensor, depending on the mode- μ fibers that are used for the unfolding. If the tensor is unfolded using mode- μ fibers, it is called the mode- μ matricization of \mathcal{X} . We refer to [29, 31] for more details on these basic concepts of tensor theory.

In the following example, we illustrate how a 3rd order tensor $\mathcal{X} \in \mathbb{R}^{n \times n \times n}$ can be unfolded into different matrices.

EXAMPLE 2.1. Consider a 3rd order tensor $\mathcal{X}^{n \times n \times n}$ whose frontal slices are given by matrices $X_i \in \mathbb{R}^{n \times n}$ as shown in Figure 2.1.

Then, its mode- μ matricizations, $\mu \in \{1, 2, 3\}$, are given by

$$\begin{aligned} \mathcal{X}^{(1)} &= [X_1, X_2, \dots, X_n], \quad \mathcal{X}^{(2)} = [X_1^T, X_2^T, \dots, X_n^T], \text{ and} \\ \mathcal{X}^{(3)} &= [\text{vec}(X_1), \text{vec}(X_2), \dots, \text{vec}(X_n)]^T. \end{aligned}$$

Similar to the matrix-matrix product, one can also perform a tensor-matrix or tensor-tensor multiplication. Of particular interest, we here only consider tensor-matrix multiplications, which can be performed by means of matricizations; see, e.g., [31]. For a given tensor $\mathcal{X} \in \mathbb{R}^{n \times n \times n}$ and a matrix $\mathcal{A} \in \mathbb{R}^{n_1 \times n}$, the μ -mode matrix product is denoted by $\mathcal{X} \times_{\mu} \mathcal{A} =: \mathcal{Y}$, i.e., $\mathcal{Y} \in \mathbb{R}^{n_1 \times n \times n}$ for $\mu = 1$. In the case

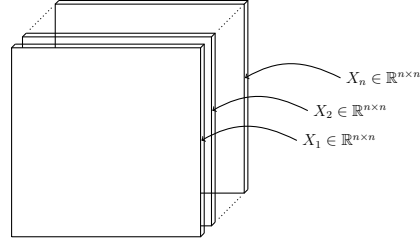


FIGURE 2.1. Representation of a tensor using the frontal slices [31].

of the μ -mode matrix multiplication, the mode- μ fiber is multiplied with the matrix \mathcal{A} , which can be written as

$$\mathcal{Y} = \mathcal{X} \times_{\mu} \mathcal{A} \Leftrightarrow \mathcal{Y}^{(\mu)} = \mathcal{A} \mathcal{X}^{(\mu)}.$$

Furthermore, if a tensor is given as

$$(2.2) \quad \mathcal{Z} = \mathcal{X} \times_1 \mathcal{A} \times_2 \mathcal{B} \times_3 \mathcal{C},$$

where $\mathcal{A} \in \mathbb{R}^{n_1 \times n}$, $\mathcal{B} \in \mathbb{R}^{n_2 \times n}$ and $\mathcal{C} \in \mathbb{R}^{n_3 \times n}$, then the mode- μ matricizations of \mathcal{Z} satisfy:

$$(2.3) \quad \mathcal{Z}^{(1)} = \mathcal{A} \mathcal{X}^{(1)} (\mathcal{C} \otimes \mathcal{B})^T, \quad \mathcal{Z}^{(2)} = \mathcal{B} \mathcal{X}^{(2)} (\mathcal{C} \otimes \mathcal{A})^T, \quad \mathcal{Z}^{(3)} = \mathcal{C} \mathcal{X}^{(3)} (\mathcal{B} \otimes \mathcal{A})^T.$$

Using these properties of the tensor products, we now introduce our first result on tensor matricizations.

LEMMA 2.2. Consider tensors $\mathcal{X}, \mathcal{Y} \in \mathbb{R}^{n \times n \times n}$ and let $\mathcal{X}^{(i)}$ and $\mathcal{Y}^{(i)}$ denote, respectively, their mode- i matricizations. Then,

$$\text{tr} \left(\mathcal{X}^{(1)} (\mathcal{Y}^{(1)})^T \right) = \text{tr} \left(\mathcal{X}^{(2)} (\mathcal{Y}^{(2)})^T \right) = \text{tr} \left(\mathcal{X}^{(3)} (\mathcal{Y}^{(3)})^T \right).$$

Proof. We begin by denoting the i th frontal slice of \mathcal{X} and \mathcal{Y} by X_i and Y_i , respectively; see Figure 2.1. Thus,

$$\begin{aligned} \text{tr} \left(\mathcal{X}^{(1)} (\mathcal{Y}^{(1)})^T \right) &= \text{tr} \left([X_1, X_2, \dots, X_n] [Y_1, Y_2, \dots, Y_n]^T \right) \\ &= \sum_{i=1}^n \text{tr} (X_i Y_i^T) = \sum_{i=1}^n \text{tr} (X_i^T Y_i) \\ &= \text{tr} \left([X_1^T, X_2^T, \dots, X_n^T] [Y_1^T, Y_2^T, \dots, Y_n^T]^T \right) = \text{tr} \left(\mathcal{X}^{(2)} (\mathcal{Y}^{(2)})^T \right). \end{aligned}$$

Furthermore, since $\text{tr} (X^T Y) = \text{vec} (X)^T \text{vec} (Y)$, this allows us to write

$$\text{tr} \left(\mathcal{X}^{(1)} (\mathcal{Y}^{(1)})^T \right) = \sum_{i=1}^n \text{tr} (X_i^T Y_i) = \sum_{i=1}^n \text{vec} (X_i)^T \text{vec} (Y_i).$$

Since the i th rows of $\mathcal{X}^{(3)}$ and $\mathcal{Y}^{(3)}$ are given by $\text{vec} (X_i)^T$ and $\text{vec} (Y_i)^T$, respectively, this means $\sum_{i=1}^n \text{vec} (X_i)^T \text{vec} (Y_i) = \text{tr} \left(\mathcal{X}^{(3)} (\mathcal{Y}^{(3)})^T \right)$. This concludes the proof. \square

Recall that the Hessian H in the QB system (1.1) is of size $n \times n^2$; thus, it can be interpreted as an unfolding of a tensor $\mathcal{H}^{n \times n \times n}$. Without loss of generality, we assume the Hessian H to be the mode-1 matricization of \mathcal{H} , i.e., $H = \mathcal{H}^{(1)}$. Also, we assume \mathcal{H} to be symmetric. This means that for given vectors u and v ,

$$(2.4) \quad H(u \otimes v) = \mathcal{H}^{(1)}(u \otimes v) = \mathcal{H}^{(1)}(v \otimes u) = H(v \otimes u).$$

This provides us additional information that the other two matricization modes of \mathcal{H} are the same, i.e.,

$$(2.5) \quad \mathcal{H}^{(2)} = \mathcal{H}^{(3)}.$$

In general, it is not necessary that the Hessian H (mode-1 matricization of the tensor \mathcal{H}), obtained from the discretization of the governing PDEs satisfies (2.4). However, as shown in [9], the Hessian H can be modified in such a way that the modified Hessian \tilde{H} satisfies (2.4) without any change in the dynamics of the system; thus, for the rest of the paper, without loss of generality, we assume that the tensor \mathcal{H} is symmetric.

The additional property that the Hessian is symmetric will allow us to derive some new relationships between matricizations and matrices, that will prove to be crucial ingredients in simplifying the expressions arising in the derivation of optimality conditions in Section 3.

LEMMA 2.3. *Let $\mathcal{H} \in \mathbb{R}^{n \times n \times n}$ be a 3rd order tensor, satisfying (2.4) and (2.5), and consider matrices $\mathcal{A}, \mathcal{B}, \mathcal{C} \in \mathbb{R}^{n \times n}$. Then,*

$$(2.6) \quad \mathcal{H}^{(1)}(\mathcal{B} \otimes \mathcal{C}) \left(\mathcal{H}^{(1)} \right)^T = \mathcal{H}^{(1)}(\mathcal{C} \otimes \mathcal{B}) \left(\mathcal{H}^{(1)} \right)^T,$$

and

$$\begin{aligned} (\text{vec } \mathcal{B})^T \text{vec} \left(\mathcal{H}^{(2)}(\mathcal{C} \otimes \mathcal{A})(\mathcal{H}^{(2)})^T \right) &= (\text{vec } \mathcal{C})^T \text{vec} \left(\mathcal{H}^{(2)}(\mathcal{B} \otimes \mathcal{A})(\mathcal{H}^{(2)})^T \right) \\ &= (\text{vec } \mathcal{A})^T \text{vec} \left(\mathcal{H}^{(1)}(\mathcal{C} \otimes \mathcal{B})(\mathcal{H}^{(1)})^T \right). \end{aligned}$$

Proof. We begin by proving the relation in (2.6). The order in the Kronecker product can be changed via pre- and post-multiplication by appropriate permutation matrices; see [26, Sec. 3]. Thus,

$$\mathcal{B} \otimes \mathcal{C} = S(\mathcal{C} \otimes \mathcal{B})S^T,$$

where S is a permutation matrix, given as $S = \sum_{i=1}^n ((e_i^n)^T \otimes I_n \otimes e_i^n)$. We can then write

$$(2.7) \quad \mathcal{H}^{(1)}(\mathcal{B} \otimes \mathcal{C}) \left(\mathcal{H}^{(1)} \right)^T = \mathcal{H}^{(1)}S(\mathcal{C} \otimes \mathcal{B}) \left(\mathcal{H}^{(1)}S \right)^T.$$

We now manipulate the term $\mathcal{H}^{(1)}S$:

$$(2.8) \quad \mathcal{H}^{(1)}S = \sum_{i=1}^n \mathcal{H}^{(1)}((e_i^n)^T \otimes I_n \otimes e_i^n).$$

Furthermore, we can write I_n in the Kronecker form as

$$(2.9) \quad I_n = \sum_{j=1}^n (e_j^n)^T \otimes e_j^n,$$

and since for a vector $f \in \mathbb{R}^q$, $f^T \otimes f = ff^T$, we can write (2.9) in another form as

$$(2.10) \quad I_n = \sum_{j=1}^n e_j^n (e_j^n)^T.$$

Substituting these relations in (2.8) leads to

$$(2.11) \quad \begin{aligned} \mathcal{H}^{(1)}S &= \sum_{i=1}^n \sum_{j=1}^n \mathcal{H}^{(1)}((e_i^n)^T \otimes (e_j^n)^T \otimes e_j^n \otimes e_i^n) \\ &= \sum_{i=1}^n \sum_{j=1}^n \mathcal{H}^{(1)}(e_j^n \otimes e_i^n) ((e_i^n)^T \otimes (e_j^n)^T) \quad (\because \text{for } f \in \mathbb{R}^q, f^T \otimes f = ff^T) \\ &= \sum_{i=1}^n \sum_{j=1}^n \mathcal{H}^{(1)}(e_i^n \otimes e_j^n)((e_i^n)^T \otimes (e_j^n)^T). \quad (\because \text{the relation (2.4)}) \end{aligned}$$

Next, we use a tensor-multiplication property in the above equation, that is,

$$(2.12) \quad (\mathcal{P} \otimes \mathcal{Q})(\mathcal{R} \otimes \mathcal{S}) = (\mathcal{P}\mathcal{R} \otimes \mathcal{Q}\mathcal{S}),$$

where $\mathcal{P}, \mathcal{Q}, \mathcal{R}$ and \mathcal{S} are of compatible dimensions. Using the Kronecker product property (2.12) in (2.11), we obtain

$$\begin{aligned} \mathcal{H}^{(1)}S &= \mathcal{H}^{(1)} \left(\sum_{i=1}^n e_i^n (e_i^n)^T \otimes \sum_{j=1}^n e_j^n (e_j^n)^T \right) \\ &= \mathcal{H}^{(1)}(I_n \otimes I_n) = \mathcal{H}^{(1)}. \end{aligned} \quad (\text{from (2.10)})$$

Substituting the above relation in (2.7) proves (2.6). For the second part, we utilize the trace property (2.1a) to obtain

$$(\text{vec}(\mathcal{B}))^T \text{vec} \left(\mathcal{H}^{(2)}(\mathcal{C} \otimes \mathcal{A})(\mathcal{H}^{(2)})^T \right) = \text{tr} \left(\underbrace{\mathcal{B}^T \mathcal{H}^{(2)}(\mathcal{C} \otimes \mathcal{A})(\mathcal{H}^{(2)})^T}_{\mathcal{L}^{(2)}} \right),$$

where $\mathcal{L}^{(2)} \in \mathbb{R}^{n \times n^2}$ can be considered as a mode-2 matricization of a tensor $\mathcal{L}^{n \times n \times n}$. Using Lemma 2.2 and the relations (2.3), we obtain

$$\begin{aligned} \text{tr} \left(\mathcal{L}^{(2)} \left(\mathcal{H}^{(2)} \right)^T \right) &= \text{tr} \left(\mathcal{L}^{(3)} \left(\mathcal{H}^{(3)} \right)^T \right) = \text{tr} \left(\mathcal{C}^T \mathcal{H}^{(3)}(\mathcal{B} \otimes \mathcal{A}) \left(\mathcal{H}^{(3)} \right)^T \right) \\ &= \text{tr} \left(\mathcal{C}^T \mathcal{H}^{(2)}(\mathcal{B} \otimes \mathcal{A}) \left(\mathcal{H}^{(2)} \right)^T \right) \quad (\text{using (2.5)}) \\ &= (\text{vec}(\mathcal{C}))^T \text{vec} \left(\mathcal{H}^{(2)}(\mathcal{B} \otimes \mathcal{A}) \left(\mathcal{H}^{(2)} \right)^T \right). \end{aligned}$$

Furthermore, we also have

$$\begin{aligned} \text{tr} \left(\mathcal{L}^{(2)} \left(\mathcal{H}^{(2)} \right)^T \right) &= \text{tr} \left(\mathcal{L}^{(1)} \left(\mathcal{H}^{(1)} \right)^T \right) = \text{tr} \left(\mathcal{A}^T \mathcal{H}^{(1)}(\mathcal{C} \otimes \mathcal{B}) \left(\mathcal{H}^{(1)} \right)^T \right) \\ &= (\text{vec}(\mathcal{A}))^T \text{vec} \left(\mathcal{H}^{(1)}(\mathcal{C} \otimes \mathcal{B}) \left(\mathcal{H}^{(1)} \right)^T \right), \end{aligned}$$

which completes the proof. \square

Next, we prove the connection between a permutation matrix and the Kronecker product.

LEMMA 2.4. *Consider two matrices $X, Y \in \mathbb{R}^{n \times m}$. Define the permutation matrix $T_{(n,m)} \in \{0, 1\}^{n^2 m^2 \times n^2 m^2}$ as*

$$(2.13) \quad T_{(n,m)} = I_m \otimes [I_m \otimes e_1^n, \dots, I_m \otimes e_n^n] \otimes I_n.$$

Then,

$$\text{vec}(X \otimes Y) = T_{(n,m)} (\text{vec}(X) \otimes \text{vec}(Y)).$$

Proof. Let us denote the i th columns of X and Y by x_i and y_i , respectively. We can then write

$$(2.14) \quad \text{vec}(X \otimes Y) = \begin{bmatrix} \text{vec}(x_1 \otimes Y) \\ \vdots \\ \text{vec}(x_m \otimes Y) \end{bmatrix}.$$

Now we concentrate on the i th block row of $\text{vec}(X \otimes Y)$, which, using (2.1b) and (2.12), can be written as

$$(2.15) \quad \begin{aligned} \text{vec}(x_i \otimes Y) &= \text{vec}((x_i \otimes I_n)Y) = (I_m \otimes x_i \otimes I_n) \text{vec}(Y) \\ &= \left(I_m \otimes [x_i^{(1)} e_1^n + \dots + x_i^{(n)} e_n^n] \otimes I_n \right) \text{vec}(Y), \end{aligned}$$

where $x_i^{(j)}$ is the (j, i) th entry of the matrix X . An alternative way to write (2.15) is

$$\begin{aligned} \text{vec}(x_i \otimes Y) &= [I_m \otimes e_1^n \otimes I_n, \dots, I_m \otimes e_n^n \otimes I_n] (x_i \otimes I_{nm}) \text{vec}(Y) \\ &= ([I_m \otimes e_1^n, \dots, I_m \otimes e_n^n] \otimes I_n) (x_i \otimes \text{vec}(Y)). \end{aligned}$$

This yields

$$\begin{aligned} \text{vec}(X \otimes Y) &= \begin{bmatrix} ([I_m \otimes e_1^n, \dots, I_m \otimes e_n^n] \otimes I_n) (x_1 \otimes \text{vec}(Y)) \\ \vdots \\ ([I_m \otimes e_1^n, \dots, I_m \otimes e_n^n] \otimes I_n) (x_m \otimes \text{vec}(Y)) \end{bmatrix} \\ &= (I_m \otimes [I_m \otimes e_1^n, \dots, I_m \otimes e_n^n] \otimes I_n) \begin{bmatrix} x_1 \otimes \text{vec}(Y) \\ \vdots \\ x_m \otimes \text{vec}(Y) \end{bmatrix} \\ &= (I_m \otimes [I_m \otimes e_1^n, \dots, I_m \otimes e_n^n] \otimes I_n) \left(\begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \otimes \text{vec}(Y) \right) \\ &= (I_m \otimes [I_m \otimes e_1^n, \dots, I_m \otimes e_n^n] \otimes I_n) (\text{vec}(X) \otimes \text{vec}(Y)), \end{aligned}$$

which proves the assertion. \square

Lemma 2.4 will be utilized in simplifying the error expressions in the next section.

3. \mathcal{H}_2 -Norm for QB Systems and Optimality Conditions. In this section, we first define the \mathcal{H}_2 -norm for the QB systems (1.1) and its truncated version. Then, based on a truncated \mathcal{H}_2 measure, we derive first-order necessary conditions for optimal model reduction. These optimality conditions will naturally lead to a numerical algorithm to construct quasi-optimal reduced models for QB systems those are independent of any training data. The proposed model reduction framework extends the optimal \mathcal{H}_2 methodology from linear [23] and bilinear systems [7, 20] to QB nonlinear systems.

3.1. \mathcal{H}_2 -norm of QB systems. In order to define the \mathcal{H}_2 -norm for QB systems and its truncated version, we first require the input/output representation for QB systems. In other words, we aim at obtaining the solution of QB systems with the help of *Volterra* series as done for bilinear systems, e.g. as in [40, Section 3.1]. For this, one can utilize the *variational analysis* [40, Section 3.4]. As a first step, we consider the external force $u(t)$ in (1.1) that is multiplied by a scalar factor α . And since the QB system falls under the class of linear-analytic systems, we can write the solution $x(t)$ of (1.1) as

$$x(t) = \sum_{s=1}^{\infty} \alpha^s x_s(t),$$

where $x_s(t) \in \mathbb{R}^n$. Substituting the expression of $x(t)$ in terms of $x_s(t)$ and replacing input $u(t)$ by $\alpha u(t)$ in (1.1), we get

$$(3.1) \quad \left(\sum_{s=1}^{\infty} \alpha^s \dot{x}_s(t) \right) = A \left(\sum_{s=1}^{\infty} \alpha^s x_s(t) \right) + H \left(\sum_{s=1}^{\infty} \alpha^s x_s(t) \right) \otimes \left(\sum_{s=1}^{\infty} \alpha^s x_s(t) \right) + \sum_{k=1}^m \alpha N_k \sum_{s=1}^{\infty} \alpha^s x_s(t) u_k(t) + \alpha B u(t).$$

By comparing the coefficients of $\alpha^i, i \in \{1, 2, \dots\}$, leads to an infinite cascade of coupled linear systems as follows:

$$(3.2) \quad \begin{aligned} \dot{x}_1(t) &= Ax_1(t) + Bu(t), \\ \dot{x}_2(t) &= Ax_2(t) + Hx_1(t) \otimes x_1(t) + \sum_{k=1}^m N_k x_1 u_k(t), \\ \dot{x}_s(t) &= Ax_s(t) + \sum_{\substack{i,j \geq 1 \\ i+j=s}} Hx_i(t) \otimes x_j(t) + \sum_{k=1}^m N_k x_{s-1}(t) u_k(t), \quad s \geq 3. \end{aligned}$$

Then, let $\alpha = 1$ so that $x(t) = \sum_{s=1}^{\infty} x_s(t)$, where $x_s(t)$ solves a coupled linear differential equation (3.2). The equation for $x_1(t)$ corresponds to a linear system, thus allowing us to write the expression for $x_1(t)$ as a convolution:

$$(3.3a) \quad x_1(t) = \int_0^t e^{At_1} B u(t - t_1) dt_1.$$

Using the expression for $x_1(t)$, we can obtain an explicit expression for $x_2(t)$:

$$\begin{aligned} x_2(t) &= \int_0^t \int_0^{t-t_3} \int_0^{t-t_3} e^{At_3} H e^{At_2} B \otimes e^{At_1} B (u(t-t_2-t_3)) \otimes u(t-t_1-t_3) dt_1 dt_2 dt_3 \\ &\quad + \sum_{k=1}^m \int_0^t \int_0^{t-t_2} e^{At_2} N_k e^{At_1} B u(t-t_1-t_2) u_k(t-t_2) dt_1 dt_2. \end{aligned}$$

Similarly, one can write down explicit expressions for $x_s(t)$, $s \geq 3$ as well, but the notation and expression become tedious, and we skip them for brevity. Then, we can write the output $y(t)$ of the QB system as $y(t) = \sum_{s=1}^{\infty} C x_s(t)$, leading to the input/output mapping as:

$$\begin{aligned} (3.5) \quad y(t) &= \int_0^t C e^{At_1} B u(t-t_1) dt_1 + \\ &\quad \int_0^t \int_0^{t-t_3} \int_0^{t-t_3} e^{At_3} H (e^{At_2} B \otimes e^{At_1} B) (u(t-t_2-t_3)) \otimes u(t-t_1-t_3) dt_1 dt_2 dt_3 + \\ &\quad \int_0^t \int_0^{t-t_2} e^{At_2} [N_1, \dots, N_m] (I_m \otimes e^{At_1} B) (u(t-t_2) \otimes u(t-t_1-t_2)) dt_1 dt_2 + \dots \end{aligned}$$

Examining the structure of (3.5) reveals that the kernels $f_i(t_1, \dots, t_i)$ of (3.5) is given by the recurrence formula

$$(3.6) \quad f_i(t_1, \dots, t_i) = C g_i(t_1, \dots, t_i),$$

where

$$\begin{aligned} (3.7) \quad g_1(t_1) &= e^{At_1} B, \\ g_2(t_1, t_2) &= e^{At_2} [N_1, \dots, N_m] (I_m \otimes e^{At_1} B), \\ g_i(t_1, \dots, t_i) &= e^{At_i} [H [g_1(t_1) \otimes g_{i-2}(t_2, \dots, t_{i-1}), \dots, g_{i-2}(t_1, \dots, t_{i-2}) \otimes g_1(t_{i-1})], \\ &\quad [N_1, \dots, N_m] (I_m \otimes g_{i-1})], \quad i \geq 3. \end{aligned}$$

As shown in [43], the \mathcal{H}_2 -norm of a bilinear system can be defined in terms of an infinite series of kernels, corresponding to the input/output mapping. Similarly, we, in the following, define the \mathcal{H}_2 -norm of a QB system based on these kernels.

DEFINITION 3.1. *Consider the QB system (1.1) whose kernels of Volterra series are defined in (3.6). Then, the \mathcal{H}_2 -norm of the QB system is given by*

$$(3.8) \quad \|\Sigma\|_{\mathcal{H}_2} := \sqrt{\text{tr} \left(\sum_{i=1}^{\infty} \int_0^{\infty} \dots \int_0^{\infty} f_i(t_1, \dots, t_i) f_i^T(t_1, \dots, t_i) dt_1 \dots dt_i \right)}.$$

This definition of the \mathcal{H}_2 -norm is not suitable for computation. Fortunately, we can find an alternative way, to compute the norm in a numerically efficient way using

matrix equations. We know from the case of linear and bilinear systems that the \mathcal{H}_2 -norms of these systems can be computed in terms of the certain system Gramians. This is also the case for QB systems. The algebraic Gramians for QB systems is recently studied in [12]. So, in the following, we extend such relations between the \mathcal{H}_2 -norm, see Definition 3.1 and the systems Gramians to QB systems.

LEMMA 3.2. *Consider a QB system with a stable matrix A , and let P and Q , respectively, be the controllability and observability Gramians of the system, which are the unique positive semidefinite solutions of the following quadratic-type Lyapunov equations:*

$$(3.9) \quad AP + PA^T + H(P \otimes P)H^T + \sum_{k=1}^m N_k P N_k^T + BB^T = 0, \quad \text{and}$$

$$(3.10) \quad A^T Q + QA + \mathcal{H}^{(2)}(P \otimes Q) \left(\mathcal{H}^{(2)} \right)^T + \sum_{k=1}^m N_k^T Q N_k + C^T C = 0.$$

Assuming the \mathcal{H}_2 -norm of a QB system exists, i.e., the infinite series in (3.8) converges, then the \mathcal{H}_2 -norm of a QB system can be computed as

$$(3.11) \quad \|\Sigma\|_{\mathcal{H}_2} := \sqrt{\text{tr}(CPC^T)} = \sqrt{\text{tr}(B^TQB)}.$$

Proof. We begin with the definition of the \mathcal{H}_2 -norm of a QB system, that is,

$$\begin{aligned} \|\Sigma_{QB}\|_{\mathcal{H}_2} &= \sqrt{\text{tr} \left(\sum_{i=1}^{\infty} \int_0^{\infty} \cdots \int_0^{\infty} f_i(t_1, \dots, t_i) f_i^T(t_1, \dots, t_i) dt_1 \dots dt_i \right)} \\ &= \sqrt{\text{tr} \left(C \left(\sum_{i=1}^{\infty} \int_0^{\infty} \cdots \int_0^{\infty} g_i(t_1, \dots, t_i) g_i^T(t_1, \dots, t_i) dt_1 \dots dt_i \right) C^T \right)}, \end{aligned}$$

where $f_i(t_1, \dots, t_i)$ and $g_i(t_1, \dots, t_i)$ are defined in (3.6) and (3.7), respectively. It is shown in [12] that

$$(3.12) \quad \left(\sum_{i=1}^{\infty} \int_0^{\infty} \cdots \int_0^{\infty} g_i(t_1, \dots, t_i) g_i^T(t_1, \dots, t_i) dt_1 \dots dt_i \right) = P,$$

where P solves (3.9) if the sequence with infinite terms converges. Thus,

$$\|\Sigma\|_{\mathcal{H}_2} = \sqrt{\text{tr}(CPC^T)}.$$

Next, we prove that $\text{tr}(CPC^T) = \text{tr}(B^TQB)$, where Q solves (3.10). Making use of Kronecker properties (2.1), we can write $\text{tr}(CPC^T)$ as

$$\text{tr}(CPC^T) = \mathcal{I}_p^T(C \otimes C) \text{vec}(P).$$

Taking $\text{vec}(\cdot)$ both sides of (3.9), we obtain

$$(3.13) \quad \left(A \otimes I_n + I_n \otimes A + \sum_{k=1}^m N_k \otimes N_k \right) \text{vec}(P) + (H \otimes H) \text{vec}(P \otimes P) + (B \otimes B) \mathcal{I}_m = 0.$$

Using [Lemma 2.4](#) in the above equation and performing some simple manipulations yield an expression for $\text{vec}(P)$ as

$$\text{vec}(P) = \mathcal{G}^{-1}(B \otimes B)\mathcal{I}_m =: P_v,$$

where

$$\mathcal{G} = - \left(A \otimes I_n + I_n \otimes A + \sum_{k=1}^m N_k \otimes N_k + (H \otimes H)T_{(n,n)}(I_{n^2} \otimes \text{vec}(P)) \right).$$

Thus,

$$(3.14) \quad \text{tr}(CPC^T) = \mathcal{I}_p^T(C \otimes C)\mathcal{G}^{-1}(B \otimes B)\mathcal{I}_m = \mathcal{I}_m^T(B^T \otimes B^T)\mathcal{G}^{-T}(C^T \otimes C^T)\mathcal{I}_p.$$

Now, let $Q_v = \mathcal{G}^{-T}(C^T \otimes C^T)\mathcal{I}_p^T$. As a result, we obtain

$$\begin{aligned} (C^T \otimes C^T)\mathcal{I}_p^T &= \text{vec}(C^T C) = \mathcal{G}^T Q_v \\ &= - \left(A^T \otimes I_n + I_n \otimes A^T + \sum_{k=1}^m N_k^T \otimes N_k^T \right) Q_v \\ &\quad + \left((H \otimes H)T_{(n,n)}(I_{n^2} \otimes P_v) \right)^T Q_v. \end{aligned}$$

Next, we consider a matrix \tilde{Q} such that $\text{vec}(\tilde{Q}) = Q_v$, which further simplifies the above equation as

$$(3.15) \quad \text{vec}(C^T C) = - \text{vec} \left(A^T \tilde{Q} + \tilde{Q} A + \sum_{k=1}^m N_k^T \tilde{Q} N_k \right) - \left((H \otimes H)T_{(n,n)}(I_{n^2} \otimes P_v) \right)^T Q_v.$$

Now, we focus on the transpose of the second part of [\(3.15\)](#), that is,

$$\begin{aligned} &Q_v^T (H \otimes H)T_{(n,n)}(I_{n^2} \otimes \text{vec}(P)) \\ &= Q_v^T (H \otimes H)T_{(n,n)} \left[e_1^{n^2} \otimes \text{vec}(P), \dots, e_{n^2}^{n^2} \otimes \text{vec}(P) \right] \\ &= Q_v^T (H \otimes H) \left[\text{vec}(\Psi_1 \otimes P), \dots, \text{vec}(\Psi_{n^2} \otimes P) \right] =: \Xi, \quad (\text{using } \text{Lemma 2.4}) \end{aligned}$$

where $\Psi_i \in \mathbb{R}^{n \times n}$ is such that $e_i^{n^2} = \text{vec}(\Psi_i)$. Using [\(2.1\)](#) and [Lemma 2.3](#), we further analyze the above equation:

$$\begin{aligned} \Xi &= \text{vec}(\tilde{Q})^T \left[\text{vec}(H(\Psi_1 \otimes P)H^T), \dots, \text{vec}(H(\Psi_{n^2} \otimes P)H^T) \right] \\ &= \text{vec}(\tilde{Q})^T \left[\text{vec}(H(P \otimes \Psi_1)H^T), \dots, \text{vec}(H(P \otimes \Psi_{n^2})H^T) \right] \\ &= \left[\text{vec}(\Psi_1)^T \text{vec} \left(\mathcal{H}^{(2)}(P \otimes \tilde{Q})(\mathcal{H}^{(2)})^T \right), \dots, \right. \\ &\quad \left. \text{vec}(\Psi_{n^2})^T \text{vec} \left(\mathcal{H}^{(2)}(P \otimes \tilde{Q})(\mathcal{H}^{(2)})^T \right) \right] \\ &= \left[\left(e_1^{n^2} \right)^T \text{vec} \left(\mathcal{H}^{(2)}(P \otimes \tilde{Q})(\mathcal{H}^{(2)})^T \right), \dots, \left(e_{n^2}^{n^2} \right)^T \text{vec} \left(\mathcal{H}^{(2)}(P \otimes \tilde{Q})(\mathcal{H}^{(2)})^T \right) \right] \\ &= \left(\text{vec} \left(\mathcal{H}^{(2)}(P \otimes \tilde{Q})(\mathcal{H}^{(2)})^T \right) \right)^T. \end{aligned}$$

Substituting this relation in (3.15) yields

$$\text{vec}(C^T C) = -\text{vec}\left(A^T \tilde{Q} + \tilde{Q} A + \sum_{k=1}^m N_k^T \tilde{Q} N_k + \mathcal{H}^{(2)}(P \otimes \tilde{Q})(\mathcal{H}^{(2)})^T\right),$$

which shows that \tilde{Q} solves (3.10) as well. Since it is assumed that Eq. (3.10) has a unique solution, thus $\tilde{Q} = Q$. Replacing $\mathcal{G}^{-T}(C^T \otimes C^T)\mathcal{I}_p^T$ by $\text{vec}(Q)$ in (3.14) and using (2.1) result in

$$\text{tr}(CPC^T) = \mathcal{I}_m^T(B^T \otimes B^T)\text{vec}(Q) = \text{tr}(B^T QB).$$

This concludes the proof. \square

It can be seen that if H is zero, the expression (3.11) boils down to the \mathcal{H}_2 -norm of bilinear systems, and if N_k are also set to zero then it provides us the \mathcal{H}_2 -norm of stable linear systems as one would expect.

REMARK 3.3. *In Lemma 3.2, we have assumed that the solutions of (3.9) and (3.10) exist, and are unique and positive semidefinite. Equivalently, the series appearing in the definition of the \mathcal{H}_2 -norm is finite (see Definition 3.1); hence, the \mathcal{H}_2 -norm exists. Naturally, the stability of the matrix A is necessary for the existence of Gramians, and a detailed study of such solutions of (3.9) and (3.10) has been carried out in [12]. However, as for bilinear systems, these Gramians may not have these desired properties when $\|N_k\|$ and $\|H\|$ are large.*

Nonetheless, a solution of these problems can be obtained via rescaling of the system as done in the bilinear case [19]. For this, we need to rescale the input variable $u(t)$ as well as the state vector $x(t)$. More precisely, we can replace $x(t)$ and $u(t)$ by $\gamma x(t) =: \tilde{x}(t)$ and $\gamma u(t) =: \tilde{u}(t)$ in (1.1). This leads to

$$(3.16) \quad \begin{aligned} \gamma \dot{\tilde{x}}(t) &= \gamma A \tilde{x}(t) + \gamma^2 H(\tilde{x}(t) \otimes \tilde{x}(t)) + \gamma^2 \sum_{k=1}^m N_k \tilde{x}(t) \tilde{u}_k(t) + \gamma B \tilde{u}(t), \\ y(t) &= \gamma C \tilde{x}(t), \quad \tilde{x}(0) = 0. \end{aligned}$$

For $\gamma \neq 0$, we get a scaled system as follows:

$$(3.17) \quad \begin{aligned} \dot{\tilde{x}}(t) &= A \tilde{x}(t) + (\gamma H)(\tilde{x}(t) \otimes \tilde{x}(t)) + \sum_{k=1}^m (\gamma N_k) \tilde{x}(t) \tilde{u}_k(t) + B \tilde{u}(t), \\ \tilde{y}(t) &= C \tilde{x}(t), \quad \tilde{x}(0) = 0, \end{aligned}$$

where $\tilde{y}(t) = y(t)/\gamma$. Now, if you compare the systems (1.1) and (3.17), the input/output mapping is the same, and the outputs of the systems differ by the factor γ . Since interpolation-like model reduction methods aim at capturing the input/output mapping, we can also use (3.17) for constructing the subspaces V and W .

Our primary aim is to determine a reduced-order system which minimizes the \mathcal{H}_2 -norm of the error system. From the derived \mathcal{H}_2 -norm expression for the QB system, it is clear that the true \mathcal{H}_2 -norm has a quite complicated structure as defined in (3.8) and does not lend itself well to deriving necessary conditions for optimality. Therefore, to simplify the problem, we focus only on the first three leading terms of the series (3.5). The main reason for considering the first three terms is that it is the minimum number of terms containing contributions from all the system matrices

(A, H, N_k, B, C) ; in other words, linear, bilinear and quadratic terms are already contained in these first three kernels. Our approach is also inspired by [20] where a truncated \mathcal{H}_2 norm is defined for bilinear systems and used to construct high-fidelity reduced-order models minimizing corresponding error measures. Therefore, based on these three leading kernels, we define a truncated \mathcal{H}_2 -norm for QB systems, denoted by $\|\Sigma\|_{\mathcal{H}_2^{(\mathcal{T})}}$. Precisely, in terms of kernels, a truncated norm can be defined as follows:

$$(3.18) \quad \|\Sigma\|_{\mathcal{H}_2^{(\mathcal{T})}} := \sqrt{\text{tr} \left(\sum_{i=1}^3 \int_0^\infty \cdots \int_0^\infty \tilde{f}_i(t_1, \dots, t_i) \left(\tilde{f}_i(t_1, \dots, t_i) \right)^T dt_1 \cdots dt_i \right)},$$

where

$$(3.19) \quad \tilde{f}_i(t_1, \dots, t_i) = C \tilde{g}_i(t_1, \dots, t_i), \quad i \in \{1, 2, 3\},$$

and

$$\begin{aligned} \tilde{g}_1(t_1) &= e^{At_1} B, & \tilde{g}_2(t_1, t_2) &= e^{At_2} [N_1, \dots, N_m] (I_m \otimes e^{At_1} B), \\ \tilde{g}_3(t_1, t_2, t_3) &= e^{At_3} H (e^{At_2} B \otimes e^{At_1} B). \end{aligned}$$

Analogous to the \mathcal{H}_2 -norm of the QB system, a truncated \mathcal{H}_2 -norm of QB systems can be determined by truncated controllability and observability Gramians associated with the QB system, denoted by $P_{\mathcal{T}}$ and $Q_{\mathcal{T}}$, respectively [12]. These truncated Gramians are the unique solutions of the following Lyapunov equations, provided A is stable:

$$(3.20a) \quad AP_{\mathcal{T}} + P_{\mathcal{T}}A^T + \sum_{k=1}^m N_k P_l N_k^T + H(P_l \otimes P_l)H^T + BB^T = 0,$$

$$(3.20b) \quad A^T Q_{\mathcal{T}} + Q_{\mathcal{T}}A + \sum_{k=1}^m N_k^T Q_l N_k + \mathcal{H}^{(2)}(P_l \otimes Q_l) \left(\mathcal{H}^{(2)} \right)^T + C^T C = 0,$$

where $\mathcal{H}^{(2)}$ is the mode-2 matricization of the QB Hessian, and P_l and Q_l are the unique solutions of the following Lyapunov equations:

$$(3.21a) \quad AP_l + P_l A^T + BB^T = 0,$$

$$(3.21b) \quad A^T Q_l + Q_l A + C^T C = 0.$$

Note that the stability of the matrix A is a necessary and sufficient condition for the existence of truncated Gramians (3.20). In what follows, we show the connection between the truncated \mathcal{H}_2 -norm and the defined truncated Gramians for QB systems.

LEMMA 3.4. *Let Σ be the QB system (1.1) with a stable A matrix. Then, the truncated \mathcal{H}_2 -norm based on the first three terms of the Volterra series is given by*

$$\|\Sigma\|_{\mathcal{H}_2^{(\mathcal{T})}} = \sqrt{\text{tr}(CP_{\mathcal{T}}C^T)} = \sqrt{\text{tr}(B^T Q_{\mathcal{T}} B)},$$

where $P_{\mathcal{T}}$ and $Q_{\mathcal{T}}$ are truncated controllability and observability Gramians of the system, satisfying (3.20).

Proof. First, we note that (3.20) and (3.21) are standard Lyapunov equations. As A is assumed to be stable, these equations have unique solutions [5]. Next, let \mathcal{R}_i be

$$\mathcal{R}_i = \int_0^\infty \cdots \int_0^\infty \tilde{f}_i(t_1, \dots, t_i) \left(\tilde{f}_i(t_1, \dots, t_i) \right)^T dt_1 \cdots dt_i,$$

where $\tilde{f}_i(t_1, \dots, t_i)$ are defined in (3.19). Thus, $\|\Sigma\|_{\mathcal{H}_2(\mathcal{T})}^2 = \text{tr} \left(C \left(\sum_{i=1}^3 \mathcal{R}_i \right) C^T \right)$. It is shown in [12] that $\sum_{i=1}^3 \mathcal{R}_i = P_{\mathcal{T}}$ solves the Lyapunov equation (3.20a). Hence,

$$\|\Sigma\|_{\mathcal{H}_2(\mathcal{T})}^2 = \text{tr} (C P_{\mathcal{T}} C^T).$$

Next, we show that $\text{tr} (C P_{\mathcal{T}} C^T) = \text{tr} (B^T Q_{\mathcal{T}} B)$. For this, we use the trace property (2.1b) to obtain:

$$\text{tr} (C P_{\mathcal{T}} C^T) = (\mathcal{I}_p)^T (C \otimes C) \text{vec} (P_{\mathcal{T}}) \quad \text{and} \quad \text{tr} (B^T Q_{\mathcal{T}} B) = (\text{vec} (Q_{\mathcal{T}}))^T (B \otimes B) \mathcal{I}_m.$$

Applying $\text{vec}(\cdot)$ to both sides of (3.20) results in

$$\begin{aligned} \text{vec} (P_{\mathcal{T}}) &= \mathcal{L}^{-1} \left((B \otimes B) \mathcal{I}_m + \sum_{k=1}^m (N_k \otimes N_k) \mathcal{L}^{-1} (B \otimes B) \mathcal{I}_m \right. \\ &\quad \left. + \text{vec} (H (P_l \otimes P_l) H^T) \right), \quad \text{and} \\ \text{vec} (Q_{\mathcal{T}}) &= \mathcal{L}^{-T} \left((C \otimes C)^T \mathcal{I}_p + \sum_{k=1}^m (N_k \otimes N_k)^T \mathcal{L}^{-T} (C \otimes C)^T \mathcal{I}_p \right. \\ &\quad \left. + \text{vec} \left(\mathcal{H}^{(2)} (P_l \otimes Q_l) \left(\mathcal{H}^{(2)} \right)^T \right) \right), \end{aligned}$$

where $\mathcal{L} = -(A \otimes I_n + I_n \otimes A)$, and P_l and Q_l solve (3.21). Thus,

$$\begin{aligned} (3.22) \quad \text{tr} (B^T Q_{\mathcal{T}} B) &= \left((\mathcal{I}_p)^T (C \otimes C) + (\mathcal{I}_p)^T (C \otimes C) \mathcal{L}^{-1} \sum_{k=1}^m (N_k \otimes N_k) \right. \\ &\quad \left. + \left(\text{vec} \left(\mathcal{H}^{(2)} (P_l \otimes Q_l) \left(\mathcal{H}^{(2)} \right)^T \right) \right)^T \right) \mathcal{L}^{-1} (B \otimes B) \mathcal{I}_m. \end{aligned}$$

Since P_l and Q_l are the unique solutions of (3.21a) and (3.21b), this gives $\text{vec} (P_l) = \mathcal{L}^{-1} (B \otimes B) \mathcal{I}_m$ and $\text{vec} (Q_l) = \mathcal{L}^{-T} (C \otimes C)^T \mathcal{I}_p$. This implies that

$$\begin{aligned} &\left(\text{vec} \left(\mathcal{H}^{(2)} (P_l \otimes Q_l) \left(\mathcal{H}^{(2)} \right)^T \right) \right)^T \text{vec} (P_l) \\ &= \text{vec} (P_l)^T \text{vec} \left(\mathcal{H}^{(2)} (P_l \otimes Q_l) \left(\mathcal{H}^{(2)} \right)^T \right) \\ &= \text{vec} (Q_l)^T \text{vec} (H (P_l \otimes P_l) H^T) \quad \text{(using Lemma 2.2)} \\ &= (\mathcal{I}_p)^T (C \otimes C) \mathcal{L}^{-1} \text{vec} (H (P_l \otimes P_l) H^T). \end{aligned}$$

Substituting the above relation in (3.22) yields

$$\begin{aligned} \text{tr}(B^T Q_{\mathcal{T}} B) &= (\mathcal{I}_p)^T (C \otimes C) \mathcal{L}^{-1} \left((B \otimes B) \mathcal{I}_m + \sum_{k=1}^m (N_k \otimes N_k) \mathcal{L}^{-1} (B \otimes B) \mathcal{I}_m \right. \\ &\quad \left. + \text{vec} \left(H(P_l \otimes P_l) H^T \right) \right) \\ &= (\mathcal{I}_p)^T (C \otimes C) \text{vec}(P_{\mathcal{T}}) = \text{tr}(C P_{\mathcal{T}} C^T). \end{aligned}$$

This concludes the proof. \square

REMARK 3.5. *One can consider the first \mathcal{M} number of terms of the corresponding Volterra series and, based on these \mathcal{M} kernels, another truncated \mathcal{H}_2 -norm can be defined. However, this significantly increases the complexity of the problem. In this paper, we stick to the truncated \mathcal{H}_2 -norm for the QB system that depends on the first three terms of the input/output mapping. And we intend to construct reduced-order systems (1.2) such that this truncated \mathcal{H}_2 -norm of the error system is minimized. Another motivation for the derived truncated \mathcal{H}_2 -norm for QB systems is that for bilinear systems, the authors in [20] showed that the \mathcal{H}_2 -optimal model reduction based on a truncated \mathcal{H}_2 -norm (with only two terms of the Volterra series of a bilinear system) also mimics the accuracy of the true \mathcal{H}_2 -optimal approximation very closely.*

3.2. Optimality conditions based on a truncated \mathcal{H}_2 -norm. We now derive necessary conditions for optimal model reduction based on the truncated \mathcal{H}_2 -norm of the error system. First, we define the QB error system. For the full QB model Σ in (1.1) and the reduced QB model $\hat{\Sigma}$ in (1.2), we can write the error system as

$$\begin{aligned} (3.23) \quad \begin{bmatrix} \dot{x}(t) \\ \dot{\hat{x}}(t) \end{bmatrix} &= \underbrace{\begin{bmatrix} A & \mathbf{0} \\ \mathbf{0} & \hat{A} \end{bmatrix}}_{A^e} \underbrace{\begin{bmatrix} x(t) \\ \hat{x}(t) \end{bmatrix}}_{x^e(t)} + \underbrace{\begin{bmatrix} H(x(t) \otimes x(t)) \\ \hat{H}(\hat{x}(t) \otimes \hat{x}(t)) \end{bmatrix}}_{N_k^e} + \sum_{k=1}^m \underbrace{\begin{bmatrix} N_k & \mathbf{0} \\ \mathbf{0} & \hat{N}_k \end{bmatrix}}_{N_k^e} \begin{bmatrix} x(t) \\ \hat{x}(t) \end{bmatrix} u_k(t) + \underbrace{\begin{bmatrix} B \\ \hat{B} \end{bmatrix}}_{B^e} u(t), \\ y^e(t) = y(t) - \hat{y}(t) &= \underbrace{\begin{bmatrix} C & -\hat{C} \end{bmatrix}}_{C^e} \begin{bmatrix} x^T(t) & \hat{x}^T(t) \end{bmatrix}^T, \quad x^e(0) = 0. \end{aligned}$$

It can be seen that the error system (3.23) is not in the conventional QB form due to the absence of the quadratic term $x^e(t) \otimes x^e(t)$. However, we can rewrite the system (3.23) into a regular QB form by using an appropriate Hessian of the error system (3.23) as follows:

$$(3.24) \quad \Sigma^e := \begin{cases} \dot{x}^e(t) = A^e x^e(t) + H^e(x^e(t) \otimes x^e(t)) + \sum_{k=1}^m N_k^e x^e(t) u_k(t) + B^e u(t), \\ y^e(t) = C^e x^e(t), \quad x^e(0) = 0, \end{cases}$$

where $H^e = \begin{bmatrix} H\mathcal{F} \\ \hat{H}\hat{\mathcal{F}} \end{bmatrix}$ with $\mathcal{F} = [I_n \ \mathbf{0}] \otimes [I_n \ \mathbf{0}]$ and $\hat{\mathcal{F}} = [\mathbf{0} \ I_r] \otimes [\mathbf{0} \ I_r]$. Next, we write the truncated \mathcal{H}_2 -norm, as defined in Lemma 3.4, for the error system (3.24). For the existence of this norm for the system (3.24), it is necessary to assume that the matrix A^e is stable, i.e., the matrices A and \hat{A} are stable. Further, we assume that the matrix \hat{A} is diagonalizable. Then, by performing basic algebraic manipulations and making use of Lemma 2.4, we obtain the expression for the error functional \mathcal{E} based on the truncated \mathcal{H}_2 -norm of the error system (3.24) as shown next.

COROLLARY 3.6. *Let Σ be the original system, having a stable matrix A , and let $\widehat{\Sigma}$ be the reduced-order original system, having a stable and diagonalizable matrix \widehat{A} . Then,*

$$(3.25) \quad \mathcal{E}^2 := \|\Sigma^e\|_{\mathcal{H}_2^{(\tau)}}^2 = (\mathcal{I}_p)^T (C^e \otimes C^e) (-A^e \otimes I_{n+r} - I_{n+r} \otimes A^e)^{-1} \left((B^e \otimes B^e) \mathcal{I}_m \right. \\ \left. + \sum_{k=1}^m (N_k^e \otimes N_k^e) \text{vec}(P_l^e) + \text{vec} \left(H^e (P_l^e \otimes P_l^e) (H^e)^T \right) \right),$$

where P_l^e solves

$$A^e P_l^e + P_l^e (A^e)^T + B^e (B^e)^T = 0.$$

Furthermore, let $\widehat{A} = R\Lambda R^{-1}$ be the spectral decomposition of \widehat{A} , and define $\widetilde{B} = R^{-1}\widehat{B}$, $\widetilde{C} = \widehat{C}R$, $\widetilde{N}_k = R^{-1}\widehat{N}_k R$ and $\widetilde{H} = R^{-1}\widehat{H}(R \otimes R)$. Then, the error can be rewritten as

$$(3.26) \quad \mathcal{E}^2 = (\mathcal{I}_p)^T \left(\widetilde{C}^e \otimes \widetilde{C}^e \right) \left(-\widetilde{A}^e \otimes I_{n+r} - I_{n+r} \otimes \widetilde{A}^e \right)^{-1} \left((\widetilde{B}^e \otimes \widetilde{B}^e) \mathcal{I}_m \right. \\ \left. + \sum_{k=1}^m \left(\widetilde{N}_k^e \otimes \widetilde{N}_k^e \right) \mathcal{P}_l + \left(\widetilde{H}^e \otimes \widetilde{H}^e \right) T_{(n+r, n+r)}(\mathcal{P}_l \otimes \mathcal{P}_l) \right),$$

where

$$(3.27) \quad \widetilde{A}^e = \begin{bmatrix} A & \mathbf{0} \\ \mathbf{0} & \Lambda \end{bmatrix}, \widetilde{N}_k^e = \begin{bmatrix} N_k & \mathbf{0} \\ \mathbf{0} & \widetilde{N}_k \end{bmatrix}, \widetilde{H}^e = \begin{bmatrix} HF \\ \widetilde{H}\widetilde{F} \end{bmatrix}, \widetilde{B}^e = \begin{bmatrix} B \\ \widetilde{B} \end{bmatrix}, \widetilde{C}^e = \begin{bmatrix} C^T \\ -\widetilde{C}^T \end{bmatrix}^T, \\ \mathcal{P}_l = \begin{bmatrix} \mathcal{P}_l^{(1)} \\ \mathcal{P}_l^{(2)} \end{bmatrix} = \begin{bmatrix} \left(-A \otimes I_{n+r} - I_n \otimes \widetilde{A}^e \right)^{-1} \left(B \otimes \widetilde{B}^e \right) \mathcal{I}_m \\ \left(-\Lambda \otimes I_{n+r} - I_r \otimes \widetilde{A}^e \right)^{-1} \left(\widetilde{B} \otimes \widetilde{B}^e \right) \mathcal{I}_m \end{bmatrix}, \text{ and}$$

$$T_{(n+r, n+r)} = I_{n+r} \otimes [I_{n+r} \otimes e_1^{n+r}, \dots, I_{n+r} \otimes e_{n+r}^{n+r}] \otimes I_{n+r}.$$

The above spectral decomposition for \widehat{A} is computationally useful in simplifying the expressions as we will see later. It reduces the optimization variables by $r(r-1)$ since Λ becomes a diagonal matrix without changing the value of the cost function (this is a state-space transformation on the reduced model, which does not change the input-output map). Even though it limits the reduced-order systems to those only having diagonalizable \widehat{A} , as observed in the linear [23] and bilinear cases [7, 20]. It is extremely rare in practice that the optimal \mathcal{H}_2 models will have a non-diagonalizable \widehat{A} ; therefore, this diagonalizability assumption does not incur any restriction from a practical perspective.

Our aim is to choose the optimization variables Λ , \widetilde{B} , \widetilde{C} , \widetilde{N}_k and \widetilde{H} such that the $\|\Sigma - \widehat{\Sigma}\|_{\mathcal{H}_2^{(\tau)}}$, i.e., equivalently the error expression (3.26), is minimized. Before we proceed further, we introduce a particular permutation matrix

$$(3.28) \quad M_{pqr} = \begin{bmatrix} I_p \otimes \begin{bmatrix} I_q \\ \mathbf{0} \end{bmatrix} & I_p \otimes \begin{bmatrix} \mathbf{0} \\ I_r \end{bmatrix} \end{bmatrix},$$

which will prove helpful in simplifying the expressions related to the Kronecker product of block matrices. For example, consider matrices $\mathcal{A} \in \mathbb{R}^{p \times p}$, $\mathcal{B} \in \mathbb{R}^{q \times q}$ and

$\mathcal{C} \in \mathbb{R}^{r \times r}$. Then, the following relation holds:

$$M_{pqr}^T \left(\mathcal{A} \otimes \begin{bmatrix} \mathcal{B} & \mathbf{0} \\ \mathbf{0} & \mathcal{C} \end{bmatrix} \right) M_{pqr} = \begin{bmatrix} \mathcal{A} \otimes \mathcal{B} & \mathbf{0} \\ \mathbf{0} & \mathcal{A} \otimes \mathcal{C} \end{bmatrix}.$$

Similar block structures can be found in the error expression \mathcal{E} in [Corollary 3.6](#), which can be simplified analogously. Moreover, due to the presence of many Kronecker products, it would be convenient to derive necessary conditions for optimality in the Kronecker formulation itself. Furthermore, these conditions can be easily translated into a theoretically equivalent framework of Sylvester equations, which are more concise, are more easily interpretable, and, more importantly, automatically lead to an effective numerical algorithm for model reduction. To this end, let $V_i \in \mathbb{R}^{n \times r}$ and $W_i \in \mathbb{R}^{n \times r}$, $i \in \{1, 2\}$, be the solutions of the following standard Sylvester equations:

$$(3.29a) \quad V_1(-\Lambda) - AV_1 = B\tilde{B}^T,$$

$$(3.29b) \quad W_1(-\Lambda) - A^T W_1 = C^T \tilde{C},$$

$$(3.29c) \quad V_2(-\Lambda) - AV_2 = \sum_{k=1}^m N_k V_1 \tilde{N}_k^T + H(V_1 \otimes V_1) \tilde{H}^T,$$

$$(3.29d) \quad W_2(-\Lambda) - A^T W_2 = \sum_{k=1}^m N_k^T W_1 \tilde{N}_k + 2 \cdot \mathcal{H}^{(2)}(V_1 \otimes W_1)(\tilde{\mathcal{H}}^{(2)})^T,$$

where $\Lambda, \tilde{N}_k, \tilde{B}$ and \tilde{C} are as defined in [Corollary 3.6](#). Furthermore, we define trial and test bases $V \in \mathbb{R}^{n \times r}$ and $W \in \mathbb{R}^{n \times r}$ as

$$(3.30) \quad V = V_1 + V_2 \quad \text{and} \quad W = W_1 + W_2.$$

We also define $\hat{V} \in \mathbb{R}^{r \times r}$ and $\hat{W} \in \mathbb{R}^{r \times r}$ (that will appear in optimality conditions as we see later) as follows:

$$(3.31) \quad \hat{V} = \hat{V}_1 + \hat{V}_2 \quad \text{and} \quad \hat{W} = \hat{W}_1 + \hat{W}_2,$$

where $\hat{V}_i \in \mathbb{R}^{r \times r}$, $\hat{W}_i \in \mathbb{R}^{r \times r}$, $i \in \{1, 2\}$ are the solutions of the set of equations in [\(3.29\)](#) but with the original system's state-space matrices being replaced with the reduced-order system ones, for example, A with \hat{A} and B with \hat{B} , etc. Next, we present first order necessary conditions for optimality, which aim at minimizing the error expression [\(3.26\)](#). The following theorem extends the truncated \mathcal{H}_2 optimal conditions from the bilinear case to the much more general quadratic-bilinear nonlinearities.

THEOREM 3.7. *Let Σ and $\hat{\Sigma}$ be the original and reduced-order systems as defined in [\(1.1\)](#) and [\(1.2\)](#), respectively. Let $\Lambda = R^{-1} \hat{A} R$ be the spectral decomposition of \hat{A} , and define $\tilde{H} = R^{-1} \hat{H}(R \otimes R)$, $\tilde{N}_k = R^{-1} \hat{N}_k R$, $\tilde{C} = \hat{C} R$, $\tilde{B} = R^{-1} \hat{B}$. If $\hat{\Sigma}$ is a reduced-order system that minimize the truncated \mathcal{H}_2 -norm of the error system [\(3.24\)](#) subject to \hat{A} being diagonalizable, then $\hat{\Sigma}$ satisfies the following conditions:*

$$(3.32a)$$

$$\text{tr} \left(C V e_i^r (e_j^p)^T \right) = \text{tr} \left(\hat{C} \hat{V} e_i^r (e_j^p)^T \right), \quad i \in \{1, \dots, r\}, \quad j \in \{1, \dots, p\},$$

$$(3.32b)$$

$$\text{tr} \left(B^T W e_i^r (e_j^m)^T \right) = \text{tr} \left(\hat{B}^T \hat{W} e_i^r (e_j^m)^T \right), \quad i \in \{1, \dots, r\}, \quad j \in \{1, \dots, m\},$$

(3.32c)

$$(W_1(:, i))^T N_k V_1(:, j) = (\widehat{W}_1(:, i))^T \widehat{N}_k \widehat{V}_1(:, j), \quad i, j \in \{1, \dots, r\}, \quad k \in \{1, \dots, m\},$$

(3.32d)

$$(W_1(:, i))^T H(V_1(:, j) \otimes V_1(:, l)) = (\widehat{W}_1(:, i))^T \widehat{H}(\widehat{V}_1(:, j) \otimes \widehat{V}_1(:, l)), \\ i, j, l \in \{1, \dots, r\},$$

(3.32e)

$$(W_1(:, i))^T V(:, i) + (W_2(:, i))^T V_1(:, i) = (\widehat{W}_1(:, i))^T \widehat{V}(:, i) + (\widehat{W}_2(:, i))^T \widehat{V}_1(:, i), \\ i \in \{1, \dots, r\}.$$

Proof. The proof is given in [Appendix B](#). \square

3.3. Truncated quadratic-bilinear iterative rational Krylov algorithm.

The remaining challenge is now to develop a numerically efficient model reduction algorithm to construct a reduced QB system satisfying first-order optimality conditions in [Theorem 3.7](#). However, as in the linear [\[23\]](#) and bilinear [\[7, 20\]](#) cases, since the optimality conditions involve the matrices $V, W, \widehat{V}, \widehat{W}$, which depend on the reduced-order system matrices we are trying to construct, it is not a straightforward task to determine a reduced-order system directly that satisfies all the necessary conditions for optimality, i.e., [\(3.32a\)–\(3.32e\)](#). We propose [Algorithm 3.1](#), which upon convergence leads to reduced-order systems that *approximately* satisfy first-order necessary conditions for optimality of [Theorem 3.7](#). Throughout the paper, we denote the algorithm by truncated QB-IRKA, or TQB-IRKA.

Algorithm 3.1 TQB-IRKA for quadratic-bilinear systems.

Input: The system matrices: $A, H, N_1, \dots, N_m, B, C$.

Output: The reduced matrices: $\widehat{A}, \widehat{H}, \widehat{N}_1, \dots, \widehat{N}_m, \widehat{B}, \widehat{C}$.

- 1: Symmetrize the Hessian H and determine its mode-2 matricization $\mathcal{H}^{(2)}$.
 - 2: Make an initial (random) guess of the reduced matrices $\widehat{A}, \widehat{H}, \widehat{N}_1, \dots, \widehat{N}_m, \widehat{B}, \widehat{C}$.
 - 3: **while** not converged **do**
 - 4: Perform the spectral decomposition of \widehat{A} and define:

$$\Lambda = R^{-1} \widehat{A} R, \quad \widetilde{N}_k = R^{-1} \widehat{N}_k R, \quad \widetilde{H} = R^{-1} \widehat{H} (R \otimes R), \quad \widetilde{B} = R^{-1} \widehat{B}, \quad \widetilde{C} = \widehat{C} R.$$
 - 5: Compute mode-2 matricization $\widetilde{\mathcal{H}}^{(2)}$.
 - 6: Solve for V_1 and V_2 :

$$-V_1 \Lambda - A V_1 = B \widetilde{B}^T, \\ -V_2 \Lambda - A V_2 = H(V_1 \otimes V_1) \widetilde{H}^T + \sum_{k=1}^m N_k V_1 \widetilde{N}_k^T.$$
 - 7: Solve for W_1 and W_2 :

$$-W_1 \Lambda - A^T W_1 = C^T \widetilde{C}, \\ -W_2 \Lambda - A^T W_2 = 2 \cdot \mathcal{H}^{(2)}(V_1 \otimes W_1) (\widetilde{\mathcal{H}}^{(2)})^T + \sum_{k=1}^m N_k^T W_1 \widetilde{N}_k.$$
 - 8: Compute V and W :

$$V := V_1 + V_2, \quad W := W_1 + W_2.$$
 - 9: $V = \text{orth}(V), \quad W = \text{orth}(W)$.
 - 10: Determine the reduced matrices:

$$\widehat{A} = (W^T V)^{-1} W^T A V, \quad \widehat{H} = (W^T V)^{-1} W^T H (V \otimes V), \\ \widehat{N}_k = (W^T V)^{-1} W^T N_k V, \quad \widehat{B} = (W^T V)^{-1} W^T B, \quad \widehat{C} = C V.$$
 - 11: **end while**
-

REMARK 3.8. *Ideally, the word upon convergence means that the reduced-order models quantities \widehat{A} , \widehat{H} , \widehat{N}_k , \widehat{B} , \widehat{C} in Algorithm 3.1 stagnate. In a numerical implementation, one can check the stagnation based on the change of eigenvalues of the reduced matrix \widehat{A} and terminate the algorithm once the relative change in the eigenvalues of \widehat{A} is less than machine precision. However, in our all numerical experiments, we run TQB-IRKA until the relative change in the eigenvalues of \widehat{A} is less than 10^{-5} . We observe that the quality of reduced-order systems does not change significantly, thereafter; as in the cases of IRKA, B-IRKA, and TB-IRKA.*

Our next goal is to show how the reduced-order system resulting from TQB-IRKA upon convergence relates to first-order optimality conditions (3.32). As a first step, we provide explicit expressions showing how farther away the resulting reduced-order system is from satisfying the optimality conditions. Later, based on these expressions, we discuss how the reduced-order systems, obtained from TQB-IRKA for weakly nonlinear QB systems, satisfy the optimality condition with small perturbations. We also illustrate using our numerical examples that the reduced-order system satisfies optimality conditions quite accurately.

THEOREM 3.9. *Let Σ be a QB system (1.1), and let $\widehat{\Sigma}$ be the reduced-order QB system (1.2), computed by TQB-IRKA upon convergence. Let $V_{\{1,2\}}, W_{\{1,2\}}, V$ and W be the projection matrices that solve (3.29) using the converged reduced-order system. Similarly, let $\widehat{V}_{\{1,2\}}, \widehat{W}_{\{1,2\}}, \widehat{V}$ and \widehat{W} that solve (3.31). Also, assume that $\sigma(\widehat{A}) \cap \sigma(-\Pi A) = \emptyset$ and $\sigma(\widehat{A}) \cap \sigma(-\Pi^T A^T) = \emptyset$, where $\Pi = V(W^T V)^{-1} W^T$ and $\sigma(\cdot)$ denotes the eigenvalue spectrum of a matrix. Furthermore, assume $\Pi_v = V_1(W^T V_1)^{-1} W^T$ and $\Pi_w = W_1(V^T W_1)^{-1} V^T$ exist. Then, the reduced-order system $\widehat{\Sigma}$ satisfies the following relations:*

(3.33a)

$$\operatorname{tr}\left(CV e_i^r (e_j^p)^T\right) = \operatorname{tr}\left(\widehat{C}\widehat{V} e_i^r (e_j^p)^T\right) + \epsilon_C^{(i,j)}, \quad i \in \{1, \dots, r\}, \quad j \in \{1, \dots, p\},$$

(3.33b)

$$\operatorname{tr}\left(B^T W e_i^r (e_j^m)^T\right) = \operatorname{tr}\left(\widehat{B}^T \widehat{W} e_i^r (e_j^m)^T\right) + \epsilon_B^{(i,j)}, \quad i \in \{1, \dots, r\}, \quad j \in \{1, \dots, m\},$$

(3.33c)

$$(W_1(:, i))^T N_k V_1(:, j) = (\widehat{W}_1(:, i))^T \widehat{N}_k \widehat{V}_1(:, j) + \epsilon_N^{(i,j,k)},$$

$$i, j \in \{1, \dots, r\}, \quad k \in \{1, \dots, m\},$$

(3.33d)

$$(W_1(:, i))^T H(V_1(:, j) \otimes V_1(:, l)) = (\widehat{W}_1(:, i))^T \widehat{H}(\widehat{V}_1(:, j) \otimes \widehat{V}_1(:, l)) + \epsilon_H^{(i,j,l)},$$

$$i, j, l \in \{1, \dots, r\},$$

(3.33e)

$$(W_1(:, i))^T V(:, i) + (W_2(:, i))^T V_1(:, i) = (\widehat{W}_1(:, i))^T \widehat{V}(:, i) + (\widehat{W}_2(:, i))^T \widehat{V}_1(:, i) + \epsilon_\lambda^{(i)},$$

$$i \in \{1, \dots, r\}.$$

where

$$\epsilon_C^{(i,j)} = -\operatorname{tr}\left(CV \Gamma_v e_i^r (e_j^p)^T\right),$$

$$\epsilon_B^{(i,j)} = -\operatorname{tr}\left(B^T W (W^T V)^{-T} \Gamma_w e_i^r (e_j^m)^T\right),$$

$$\epsilon_N^{(i,j,k)} = (\epsilon_w(:, i))^T N_k (V_1(:, j) - \epsilon_v(:, j)) + (W_1(:, i))^T N_k (\epsilon_v(:, j)),$$

$$\begin{aligned} \epsilon_H^{(i,j,l)} &= (W_1(:,i) - \epsilon_w(:,i))^T H(\epsilon_v(:,j) \otimes (V_1(:,l) - \epsilon_v(:,l)) + V_1(:,j) \otimes \epsilon_v(:,l)) \\ &\quad + (\epsilon_w(:,i))^T H((V(:,j) - \epsilon_v(:,j)) \otimes (V_1(:,l) - \epsilon_v(:,l))), \text{ and} \\ \epsilon_\lambda^{(i)} &= -\left(\widehat{W}(:,i)\right)^T \Gamma_v(:,i) - (\Gamma_w(:,i))^T \left(\widehat{V}(:,i) - \Gamma_v(:,i)\right) \\ &\quad - (W_2(:,i))^T V_2(:,i) + (\widehat{W}_2(:,i))^T \widehat{V}_2(:,i), \end{aligned}$$

in which ϵ_v , ϵ_w , Γ_v and Γ_w , respectively, solve

(3.34a)

$$\epsilon_v \Lambda + \Pi A \epsilon_w = (\Pi - \Pi_v)(A V_1 + B \widetilde{B}^T),$$

(3.34b)

$$\epsilon_w \Lambda + (A \Pi)^T \epsilon_w = (\Pi^T - \Pi_w)(A^T W_1 + C^T \widetilde{C}),$$

(3.34c)

$$\Gamma_v \Lambda + \widehat{A} \Gamma_v = -(W^T V)^{-1} W^T \left(\sum_{k=1}^m N_k \epsilon_v \widetilde{N}_k^T + H(\epsilon_v \otimes (V_1 + \epsilon_v) + V_1 \otimes \epsilon_v) \widetilde{H}^T \right)$$

(3.34d)

$$\Gamma_w \Lambda + \widehat{A}^T \Gamma_w = V^T \left(\sum_{k=1}^m N_k^T \epsilon_w \widetilde{N}_k + \mathcal{H}^{(2)}(\epsilon_v \otimes (W_1 + \epsilon_w) + V_1 \otimes \epsilon_w) \left(\mathcal{H}^{(2)} \right)^T \right).$$

Proof. The proof is given in [Appendix C](#). \square

REMARK 3.10. In [Theorem 3.9](#), we have presented measures, e.g., the distance between $\text{tr}(C V e_i^r (e_j^p)^T)$ and $\text{tr}(\widehat{C} \widehat{V} e_i^r (e_j^p)^T)$, denoted $\epsilon_C^{(i,j)}$ with which the reduced-order system, via TQB-IRKA, satisfies the optimality conditions [\(3.32\)](#). Even though [Theorem 3.9](#) does not provide a guarantee for smallness of these distances, we indeed observe in all the numerical examples in [Section 4](#) that these deviations from the exact optimality conditions are rather small. Here, we provide an intuition for these numerical results. Recall that V_1 and V_2 solve the Sylvester equations [\(3.29a\)](#) and [\(3.29c\)](#), respectively, and the right-hand side for V_2 is quadratic in H and N_k . So, for a weakly nonlinear QB system, meaning $\|H\|$ and $\|N_k\|$ are small w.r.t. $\|B\|$, the matrix V_2 will be relatively small compared to the matrix V_1 . Hence, the subspace V is expected to be close to V_1 . Thus, one could anticipate that the projectors $\Pi = V(W^T V)W^T$ and $\Pi_v = V_1(W^T V_1)W^T$ will be close to each other. As a result, the right-hand side of the Sylvester equation [\(3.34a\)](#) will be small, and hence thus so is ϵ_v . In a similar way, one can argue that ϵ_w in [\(3.34b\)](#) will be small. Therefore, it can be shown that in the case of weakly nonlinear QB systems [\(1.1\)](#), all ϵ 's in [\(3.33\)](#) such as $\epsilon_C^{(i,j)}$ will be very small.

Indeed, the situation in practice proves much better. We observe in our numerical results (see [Section 4](#)) that even for strongly nonlinear QB systems, i.e., $\|H\|$ and $\|N_k\|$ are comparable or even much larger than $\|B\|$, [Algorithm 3.1](#) yields reduced-order systems which satisfy the optimality conditions [\(3.32\)](#) almost exactly with negligible perturbations.

REMARK 3.11. [Algorithm 3.1](#) can be seen as an extension of the truncated B-IRKA with the truncation index 2 [[20](#), Algo. 2] from bilinear systems to QB systems. In [[20](#)], the truncation index \mathcal{N} , which denotes the number of terms in the underlying Volterra series for bilinear systems, is free, and as $\mathcal{N} \rightarrow \infty$, all the perturbations go

to zero. However, it is shown in [20] that in most cases, a small N , for example, 2 or 3 is enough to satisfy all optimality conditions very closely. In our case, a similar convergence will occur if we let the number of terms in the underlying Volterra series of the QB system grow; however, this is not numerically feasible since the subsystems in the QB case become rather complicated after the first three terms. Indeed, because of this, [22] and [9] have considered the interpolation of multivariate transfer functions corresponding to only the first two subsystems. Moreover, even in the case of balanced truncation for QB systems [12], it is shown by means of numerical examples that the truncated Gramians for QB systems based on the first three terms of the underlying Volterra series produce quantitatively very accurate reduced-order systems. Our numerical examples show that this is the case here as well.

REMARK 3.12. So far in all of our discussions, we have assumed that the reduced matrix \hat{A} is diagonalizable. This is a reasonable assumption since non-diagonalizable matrices lie in a set of the Lebesgue measure zero. The probability of entering this set by any numerical algorithm including TQB-IRKA is zero with respect to the Lebesgue measure. Thus, TQB-IRKA can be considered safe in this regard.

Furthermore, throughout the analysis, it has been assumed that the reduced matrix \hat{A} is Hurwitz. However, in case \hat{A} is not Hurwitz, then the truncated \mathcal{H}_2 -norm of the error system will be unbounded; thus the reduced-order systems indeed cannot be (locally) optimal. Nonetheless, a mathematical study to ensure the stability from \mathcal{H}_2 iterative schemes are still under investigation even for linear systems. However, a simple fix to this problem is to reflect the unstable eigenvalues of \hat{A} in every step back to the left-half plane. See also [30] for a more involved approach to stabilize a reduced order system.

REMARK 3.13. Thus far, we have used $E = I$ for QB systems; however, in the case of $E \neq I$, but nonetheless a nonsingular matrix, we can still employ Algorithm 3.1. One obvious way is to invert E , but this is inadmissible in the large-scale setting. Moreover, the resulting matrices may be dense, making the algorithm computationally expensive. Nevertheless, Algorithm 3.1 can be employed without inverting E . For this, we need to modify steps 6 and 7 in Algorithm 3.1 as follows:

$$\begin{aligned} -EV_1\Lambda - AV_1 &= B\tilde{B}^T, \\ -EV_2\Lambda - AV_2 &= H(V_1 \otimes V_1)\tilde{H}^T + \sum_{k=1}^m N_k V_1 \tilde{N}_k^T, \\ -E^T W_1 \Lambda - A^T W_1 &= C^T \tilde{C}, \\ -E^T W_2 \Lambda - A^T W_2 &= 2 \cdot \mathcal{H}^{(2)}(V_1 \otimes W_1)(\tilde{\mathcal{H}}^{(2)})^T + \sum_{k=1}^m N_k^T W_1 \tilde{N}_k, \end{aligned}$$

and replace $(W^T V)^{-1}$ with $(W^T E V)^{-1}$, assuming $W^T E V$ is invertible while determining the reduced system matrices in step 9 of Algorithm 3.1. Then, the modified iterative algorithm with the matrix E also provides a reduced-order system, approximately satisfying optimality conditions. We skip the rigorous proof for the $E \neq I$ case, but it can be proven along the lines of $E = I$. Indeed, one does not even need to invert $W^T E V$ by letting the reduced QB have a reduced E term as $W^T E V$, and the spectral decomposition of \hat{A} is replaced by a generalized eigenvalue decomposition of \hat{E} and \hat{A} . However, to keep the notation of Algorithm 3.1 simple, we omit these details.

3.4. Computational issues. The main bottleneck in applying TQB-IRKA is the computations of the reduced matrices, especially the computational cost related to $\hat{H} := W^T H(V \otimes V)$ that needs to be evaluated at each iteration. Regarding

Algorithm 3.2 Computation of the Hessian of the reduced QB system [9].

- 1: Determine $\mathcal{Y} \in \mathbb{R}^{r \times n \times n}$, such that $\mathcal{Y}^{(1)} = W^T H$.
 - 2: Determine $\mathcal{Z} \in \mathbb{R}^{r \times r \times n}$, such that $\mathcal{Z}^{(2)} = V^T \mathcal{Y}^{(2)}$.
 - 3: Determine $\mathcal{X} \in \mathbb{R}^{r \times r \times r}$, such that $\mathcal{X}^{(3)} = V^T \mathcal{Z}^{(3)}$.
 - 4: Then, the reduced Hessian is $\hat{H} = \mathcal{X}^{(1)}$.
-

this, there is an efficient method, proposed in [9], utilizing the properties of tensor matricizations, which we summarize in [Algorithm 3.2](#).

[Algorithm 3.2](#) avoids the highly undesirable explicit formulation of $V \otimes V$ for large-scale systems to compute the reduced Hessian, and the algorithm does not rely on any particular structure of the Hessian. However, a QB system obtained using semi-discretization of a PDEs usually leads to a Hessian which has a particular structure related to that particular PDE and the choice of the discretization method.

Therefore, we propose another efficient way to compute \hat{H} that utilizes a particular structure of the Hessian, arising from the governing PDEs or ODEs. Generally, the term $H(x \otimes x)$ in the QB system (1.1) can be written as $H(x \otimes x) = \sum_{j=1}^p (\mathcal{A}^{(j)} x) \circ (\mathcal{B}^{(j)} x)$, where \circ denotes the Hadamard product, and $\mathcal{A}^{(j)}$ and $\mathcal{B}^{(j)}$ are sparse matrices, depending on the nonlinear operators in PDEs and discretization scheme, and p is generally a very small integer; for instance, it is equal to 1 in case of Burgers' equations. Furthermore, using the i th rows of $\mathcal{A}^{(j)}$ and $\mathcal{B}^{(j)}$, we can construct the i th row of the Hessian:

$$H(i, :) = \sum_{j=1}^p \mathcal{A}^{(j)}(i, :) \otimes \mathcal{B}^{(j)}(i, :),$$

where $H(i, :)$, $\mathcal{A}^{(j)}(i, :)$ and $\mathcal{B}^{(j)}(i, :)$ represent i th row of the matrix H , $\mathcal{A}^{(j)}$ and $\mathcal{B}^{(j)}$, respectively. This clearly shows that there is a particular Kronecker structure of the Hessian H , which can be used in order to determine \hat{H} . Using the Chafee-Infante equation as an example, we illustrate how the structure of the Hessian (Kronecker structure) can be exploited to determine \hat{H} efficiently.

EXAMPLE 3.14. *Here, we consider the Chafee-Infante equation, which is discretized over the spatial domain via a finite difference scheme. The MOR problem for this example is considered in [Subsection 4.1](#), where one can also find the governing equations and boundary conditions. For this particular example, the Hessian (after having rewritten the system into the QB form) is given by*

$$\begin{aligned} H(i, :) &= -\frac{1}{2}e_i^n \otimes e_{k+i}^n - \frac{1}{2}e_{k+i}^n \otimes e_i^n, \quad i \in \{1, \dots, k\}, \\ H(i, :) &= -2(e_i^n \otimes e_i^n) + e_{i-k}^n \otimes [X(i-k, :)\ \mathbf{0}] + [X(i-k, :)\ \mathbf{0}] \otimes e_{i-k}^n, \\ &\quad i \in \{k+1, \dots, n\}, \end{aligned}$$

where k is the number of grid points, $n = 2k$, and $H(i, :)$ is the i th row vector of the matrix H . $X(i, :)$ also denotes the i th row vector of the matrix $X \in \mathbb{R}^{k \times k}$

$$(3.35) \quad X = \begin{bmatrix} 0 & 1 & & \\ 1 & 0 & \ddots & \\ & \ddots & \ddots & 1 \\ & & & 1 & 0 \end{bmatrix}.$$

Algorithm 3.3 Computation of the reduced Hessian for Chafee-Infante example.

- 1: **Input:** $V, W \in \mathbb{R}^{2k \times r}$, $X \in \mathbb{R}^{k \times k}$ (as defined in (3.35))
 - 2: Compute $V_x := XV(1:k, :)$, where $V(1:k, :)$ denotes the first k row vectors of V .
 - 3: **for** $i = 1 : k$ **do**
 - 4: $H_v(i, :) = -\frac{1}{2}V(i, :) \otimes V(k+i, :) - \frac{1}{2}V(i, :) \otimes V(k+i, :)$,
 $H_v(k+i, :) = -2(V(i, :) \otimes V(i, :)) + V(i, :) \otimes V_x(i, :) + V_x(i, :) \otimes V(i, :)$,
 where $H_v(q, :)$ is the q th row vector of H_v and the same holds for other matrices.
 - 5: **end for**
 - 6: Then, the reduced Hessian is $\hat{H} = W^T H_v$.
-

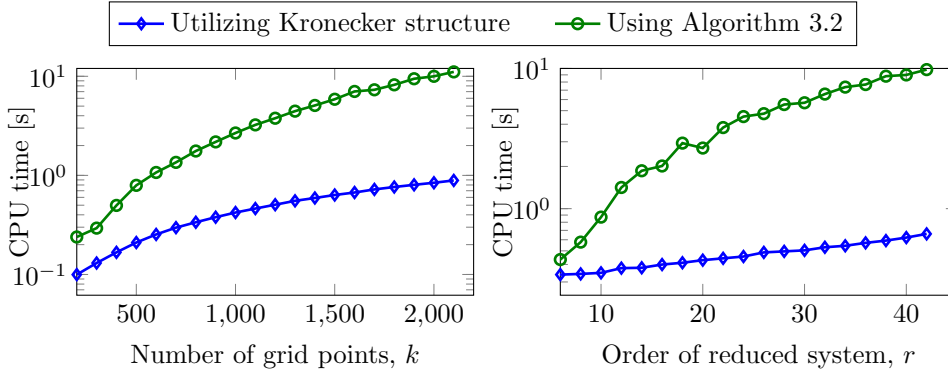


FIGURE 3.1. The left figure shows the computational time for $\hat{H} := W^T H(V \otimes V)$ by varying the number of grid points in the spatial domain by fixing the order of the reduced system to $r = 20$. In the right figure, we show the computational time for different orders of the reduced system using a fix number of grid points, $k = 1000$.

The Kronecker representation of each row of the matrix H allows us to compute each row of $H_v := H(V \otimes V)$ by selecting only the required rows of V . This way, we can determine H_v efficiently for large-scale settings, and then multiply with W^T to obtain the desired reduced Hessian. We note down the step in Algorithm 3.3 that shows how one can determine the reduced Hessian for the Chafee-Infante example.

In order to show the effectiveness of the proposed methodology that uses the special Kronecker structure of the Hessian H , we compute $\hat{H} = W^T H(V \otimes V)$ for different orders of original and reduced-order systems and show the required CPU-time to compute it in Figure 3.1.

These figures show that if one utilizes the present Kronecker structure of the Hessian H , then the computational cost does not increase as rapidly as one would use Algorithm 3.2 to compute the reduced Hessian. So, our conclusion here is that it is worth to exploit the Kronecker structure of the Hessian of the system for an efficient computation of \hat{H} in the large-scale settings.

4. Numerical Results. In this section, we illustrate the efficiency of the proposed model reduction method TQB-IRKA for QB systems by means of several semi-discretized nonlinear PDEs, and compare it with the existing MOR techniques, such as one-sided and two-sided subsystem-based interpolatory projection methods [9, 22, 38], balanced truncation (BT) for QB systems [12] and POD, e.g., see [27, 32]. We iterate Algorithm 3.1 until the relative change in the eigenvalues of \hat{A} becomes smaller

than a tolerance level. We set this tolerance to $1e-5$. Moreover, we determine interpolation points for the one-sided and two-sided interpolatory projection methods using IRKA [23] on the corresponding linear part, which appear to be a good set of interpolation points as shown in [9]. All the simulations were done on a board with 4 Intel[®] Xeon[®] E7-8837 CPUs with a 2.67-GHz clock speed using MATLAB 8.0.0.783 (R2012b). Some more details related to the numerical examples are as follows:

1. All original and reduced-order systems are integrated by the routine `ode15s` in MATLAB with relative error tolerance of 10^{-8} and absolute error tolerance of 10^{-10} .
2. We measure the output at 500 equidistant points within the time interval $[0, T]$, where T is defined in each numerical example.
3. In order to employ BT, we need to solve four standard Lyapunov equations. For this, we use `mess_lyap.m` from M.E.S.S.-1.0.1 which is based on one of the latest ADI methods proposed in [13].
4. We initialize TQB-IRKA (Algorithm 3.1) by choosing an arbitrary reduced system while ensuring \hat{A} is Hurwitz and diagonalizable. One can also first compute the projection matrices V and W by employing IRKA on the corresponding linear system and, using it, we can compute the initial reduced-order system. But we observe that random reduced systems give almost the same convergence rate or even better sometimes. Therefore, we stick to a random initial guess selection.
5. Since POD can be applied to a general nonlinear system, we apply POD to the original nonlinear system, without transforming it into a QB system as we observe that this way, POD yields better reduced systems.
6. One of the aims of the numerical examples is to determine the residuals in Theorem 3.9. For this, we first define $\Phi_C^e \in \mathbb{R}^{r \times p}$, $\Phi_B^e \in \mathbb{R}^{r \times m}$, $\Phi_N^e \in \mathbb{R}^{r \times r \times m}$, $\Phi_H^e \in \mathbb{R}^{r \times r \times r}$ and $\Phi_\Lambda^e \in \mathbb{R}^r$ such that $\epsilon_C^{(i,j)}$ is the (i, j) th entry of Φ_C^e , $\epsilon_B^{(i,j)}$ is the (i, j) th entry of Φ_B^e , $\epsilon_N^{(i,j,k)}$ is the (i, j, k) th entry of Φ_N^e , $\epsilon_H^{(i,j,k)}$ is the (i, j, k) th entry of Φ_H^e and $\epsilon_\Lambda^{(i)}$ is i th entry of Φ_Λ^e . Furthermore, we define Φ_C , Φ_B , Φ_N , Φ_H and Φ_Λ to be the terms on the left hand side of equations (3.33a) – (3.33e) in Theorem 3.9, e.g., the (i, j) th entry of Φ_C is $\text{tr}(C V e_i^r (e_j^p)^T)$. As a result, we define a relative perturbation measures as follows:

$$(4.1) \quad \mathcal{E}_C = \frac{\|\Phi_C^e\|_2}{\|\Phi_C\|_2}, \quad \mathcal{E}_B = \frac{\|\Phi_B^e\|_2}{\|\Phi_B\|_2}, \quad \mathcal{E}_N = \frac{\|\Phi_N^{e1}\|_2}{\|\Phi_N^{(1)}\|_2}, \quad \mathcal{E}_H = \frac{\|\Phi_H^{e1}\|_2}{\|\Phi_H^{(1)}\|_2}, \quad \mathcal{E}_\Lambda = \frac{\|\Phi_\Lambda^e\|_2}{\|\Phi_\Lambda\|_2},$$

where $\Phi_{\{N,H\}}^{(1)}$ and $\Phi_{\{N,H\}}^{e1}$, respectively are mode-1 matricizations of the tensors $\Phi_{\{N,H\}}$ and $\Phi_{\{N,H\}}^e$.

7. We also address a numerical issue which one might face while employing Algorithm 3.1. In step 8 of Algorithm 3.1, we need to take a sum of two matrices V_1 and V_2 , but if H and N_k are too large, then the norm of V_2 can be much larger than that of V_1 . Thus, a direct sum might reduce the effect of V_1 . As a remedy we propose to use a scaling factor γ for H and N_k , thus resulting in the matrices V_1 and V_2 which have almost the same order of norm. We have already noted in Remark 3.3 that this scaling does not change the input-output behavior.

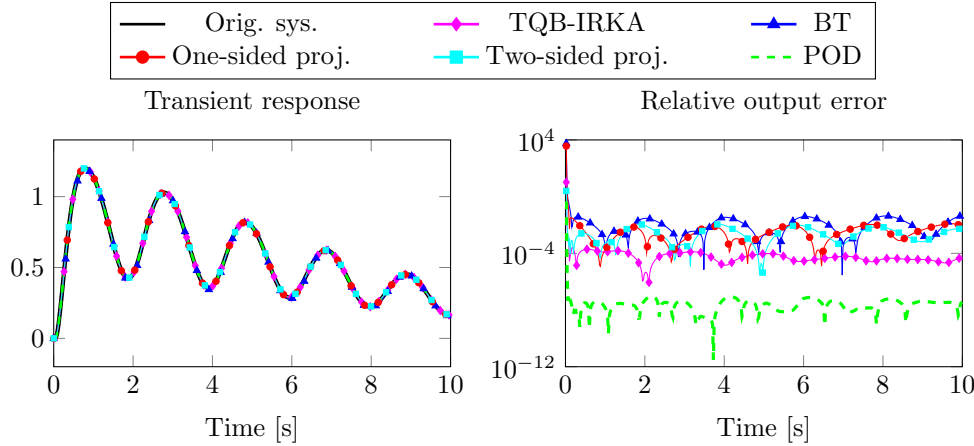


FIGURE 4.1. *Chafee-Infante*: comparison of responses for the boundary control input $u_1(t) = (1 + \sin(\pi t)) \exp(-t/5)$.

4.1. One dimensional Chafee-Infante equation. Here, we consider the one-dimensional Chafee-Infante (Allen-Cahn) equation whose governing equation, initial condition and boundary controls are given by

$$(4.2) \quad \begin{aligned} \dot{v} + v^3 &= v_{xx} + v, & (0, L) \times (0, T), & & v(0, \cdot) &= u(t), & (0, T), \\ v_x(L, \cdot) &= 0, & (0, T), & & v(x, 0) &= 0, & (0, L). \end{aligned}$$

MOR for this system has been considered in various articles; see, e.g., [9, 12]. The governing equation (4.2) contains a cubic nonlinearity, which can then be rewritten into QB form as shown in [9]. For more details on the system, we refer to [17, 25]. Next, we utilize a finite difference scheme by using k equidistant points over the length, resulting in a semi-discretized QB system of order $2k$. The output of our interest is the response at the right boundary, i.e., $v(L, t)$, and we set the number of grid points to $k = 500$, leading to an order $n = 1000$ QB system.

We construct reduced-order systems of order $r = 10$ using TQB-IRKA, BT, one-sided and two-sided interpolatory projection methods, and POD. Having initialized TQB-IRKA randomly, it takes 9 iterations to converge, and for this example, we choose the scaling factor $\gamma = 0.01$. We compute the reduced Hessian as shown in Algorithm 3.3. For the POD based approximation, we collect 500 snapshots of the true solution for the training input $u_1(t) = (1 + \sin(\pi t)) \exp(-t/5)$ and compute the projection by taking the 10 dominant basis vectors.

In order to compare the quality of these reduced-order systems with respect to the original system, we first simulate them using the same training input used to construct the POD basis, i.e., $u_1(t) = (1 + \sin(\pi t)) \exp(-t/5)$. We plot the transient responses and relative output errors for the input in Figure 4.1. As expected, since we are comparing the reduced models for the same forcing term used for POD, Figure 4.1 shows that the POD approximation outperforms the other methods for the input u_1 . However, the interpolatory methods also provide adequate reduced-order systems for $u_1(t)$ even though the reduction is performed without any knowledge of $u_1(t)$.

To test the robustness of the reduced systems, we compare the time-domain simulations of the reduced systems with the original one in Figure 4.2 for a slightly

Input	TQB-IRKA	BT	One-sided	Two-sided	POD
$u_1(t)$	$6.54 \cdot 10^{-5}$	$1.40 \cdot 10^{-2}$	$4.30 \cdot 10^{-3}$	$3.51 \cdot 10^{-3}$	$2.87 \cdot 10^{-8}$
$u_2(t)$	$1.63 \cdot 10^{-3}$	$1.43 \cdot 10^{-2}$	$4.59 \cdot 10^{-1}$	$6.65 \cdot 10^{-3}$	$6.70 \cdot 10^{-2}$

TABLE 4.1

Chafee-Infante: the mean relative errors of the output.

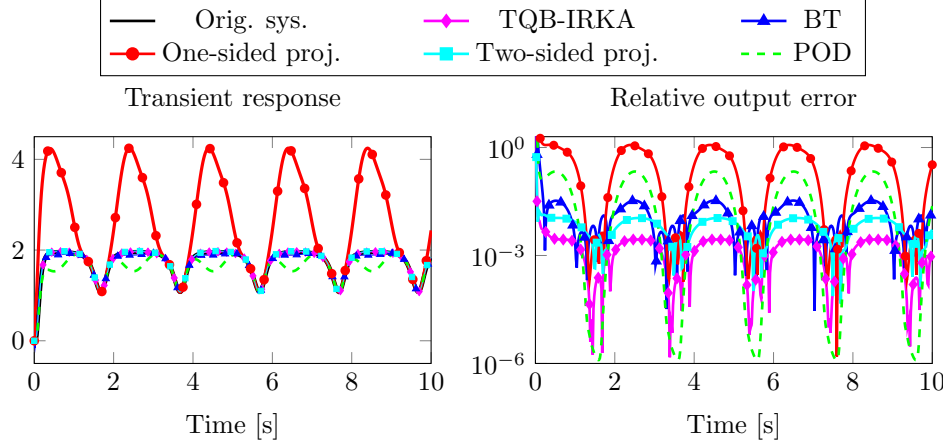


FIGURE 4.2. *Chafee-Infante: comparison of responses for the boundary control input $u_2(t) = 25(1 + \sin(\pi t))$.*

\mathcal{E}_C	\mathcal{E}_B	\mathcal{E}_N	\mathcal{E}_H	\mathcal{E}_λ
1.35×10^{-8}	8.85×10^{-12}	8.84×10^{-16}	1.77×10^{-13}	1.44×10^{-11}

TABLE 4.2

Chafee-Infante: perturbations to the optimality conditions.

different input, namely $u_2(t) = 25(1 + \sin(\pi t))$. First, observe that the POD approximation fails to reproduce the system’s dynamics for the input u_2 accurately as POD is mainly an input-dependent algorithm. Moreover, the one-sided interpolatory projection method also performs worse for the input u_2 . On the other hand, TQB-IRKA, BT, and the two-sided interpolatory projection method, all yield very accurate reduced-order systems of comparable qualities; TQB-IRKA produces marginally better reduced systems. Once again it is important to emphasize that neither $u_1(t)$ nor $u_2(t)$ have entered the model reduction procedure in TQB-IRKA. To give a quantitative comparison of the reduced systems for both inputs, u_1 and u_2 , we report the mean relative errors in Table 4.1 as well, which also provides us a similar information.

In Theorem 3.9, we have presented the quantities, denoted by ϵ_X , where $X = \{C, B, N, H, \lambda\}$ which measure how far the reduced-order system upon convergence of TQB-IRKA is from satisfying the optimality conditions (3.32). These quantities can be computed as shown in (4.1), and are listed in Table 4.2, which shows a very small magnitude of these perturbations. In Remark 3.10, we have argued that for a weakly nonlinear QB system, we expect these quantities to be small. However, even for this example with strong nonlinearity, i.e., $\|H\|$ and $\|N_k\|$ are not small at all, the reduced-order system computed by TQB-IRKA satisfies the optimality conditions (3.32) very accurately. This result also strongly supports the discussion of Remark 3.11 that a

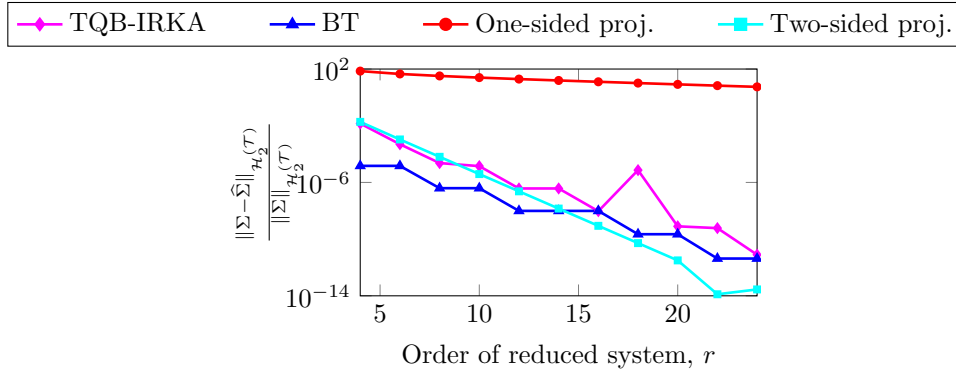


FIGURE 4.3. *Chafee-Infante*: comparison of the truncated \mathcal{H}_2 -norm of the error system, having obtained reduced systems of different orders via different methods.

small truncation index is expected to be enough in many cases.

Furthermore, since TQB-IRKA approximately minimizes the truncated \mathcal{H}_2 -norm of the error system, i.e. $\|\Sigma - \hat{\Sigma}\|_{\mathcal{H}_2^{(\mathcal{T})}}$, we also compare the truncated \mathcal{H}_2 -norm of the error system in Figure 4.3, where the $\hat{\Sigma}$ are constructed by various methods of different orders. As mentioned before, the reduced-order systems obtained via POD preserve the structure of the original nonlinearities; therefore, the truncated \mathcal{H}_2 -norm definition, given in Lemma 3.4, does not apply.

Figure 4.3 indicates that the reduced-order systems obtained via one-sided interpolatory projection perform worst in the truncated \mathcal{H}_2 -norm, whereas the quality of the reduced-order systems obtained by TQB-IRKA, BT, and the two-sided interpolatory projection method are comparable with respect to the truncated \mathcal{H}_2 -norm. This also potentially explains why the reduced systems obtained via TQB-IRKA, BT and two-sided interpolatory projection are of the same quality for both control inputs u_1 and u_2 , and the one-sided interpolatory projection method provides the worst reduced systems, see Figure 4.2. However, it is important to emphasize that unlike for linear dynamical systems, the \mathcal{H}_2 -norm and the L^∞ -norm of the output for nonlinear systems, including QB systems, are not directly connected. This can be seen in Figure 4.3 as for reduced systems of order $r = 10$, even though BT yields the smallest truncated \mathcal{H}_2 -norm of the error system, but in time-domain simulations for inputs u_1 and u_2 , it is not the best in terms of the L^∞ -norm of the output error. Nevertheless, the truncated \mathcal{H}_2 -norm of the error system is still a robust indicator for the quality of the reduced system.

4.2. A nonlinear RC ladder. We consider a nonlinear RC ladder, which consists of capacitors and nonlinear I-V diodes. The characteristics of the I-V diode are governed by exponential nonlinearities, which can also be rewritten in QB form. For a detailed description of the dynamics of this electronic circuit, we refer to [3, 22, 34, 37, 38]. We set the number of capacitors in the ladder to $k = 500$, resulting in a QB system of order $n = 1000$. Note that the matrix A of the obtained QB system has eigenvalues at zero; therefore, the truncated \mathcal{H}_2 -norm does not exist. Moreover, BT also cannot be employed as we need to solve Lyapunov equations that require a stable A matrix. Thus, we shift the matrix A to $A_s := A - 0.01I_n$ to determine the projection matrices for TQB-IRKA and BT but we project the original system matrices.

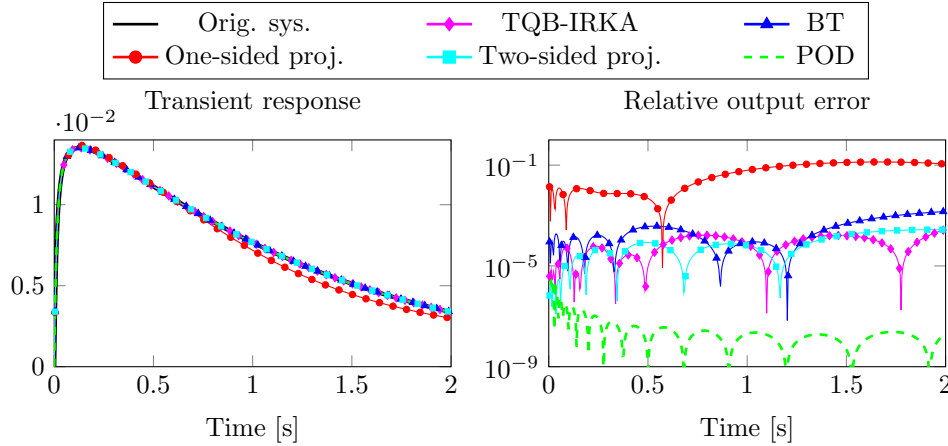


FIGURE 4.4. An RC circuit: comparison of responses for the input $u_1(t) = e^{-t}$.

Input	TQB-IRKA	BT	One-sided	Two-sided	POD
$u_1(t)$	$8.82 \cdot 10^{-5}$	$3.67 \cdot 10^{-4}$	$6.50 \cdot 10^{-2}$	$1.01 \cdot 10^{-4}$	$7.24 \cdot 10^{-8}$
$u_2(t)$	$1.12 \cdot 10^{-3}$	$2.15 \cdot 10^{-3}$	$2.32 \cdot 10^{-1}$	$7.80 \cdot 10^{-4}$	$7.8 \cdot 10^{-3}$

TABLE 4.3

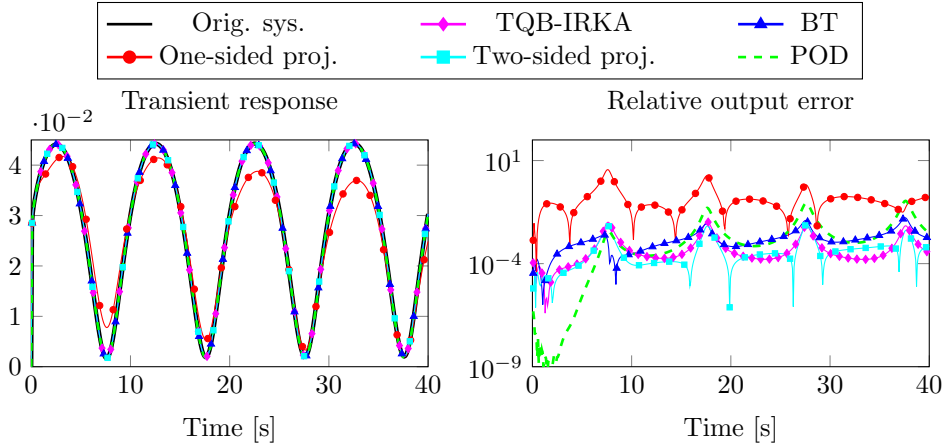
An RC circuit: the mean absolute errors of the output.

We construct reduced-order systems of order $r = 10$ using all five different methods. In this example as well, we initialize TQB-IRKA randomly and it converges after 27 iterations. We choose the scaling factor $\gamma = 0.01$. For this example, we determine the reduced Hessian by exploiting the particular structure of the Hessian. In order to compute reduced-order systems via POD, we first obtain 500 snapshots of the true solution for the training input $u_1 = e^{-t}$ and then use the 10 dominant modes to determine the projection.

We first compare the accuracy of these reduced systems for the same training input $u_1(t) = e^{-t}$ that is also used to compute the POD basis. Figure 4.4 shows the transient responses and relative errors of the output for the input u_1 . As one would expect, POD outperforms all other methods since the control input u_1 is the same as the training input for POD. Nonetheless, TQB-IRKA, BT, and two-sided projection also yield very good reduced-order systems, considering they are obtained without any prior knowledge of the input.

We also test the reduced-order systems for a different input than the training input, precisely, $u_2 = 2.5(\sin(\pi t/5) + 1)$. Figure 4.5 shows the transient responses and relative errors of the output for the input u_2 . We observe that POD does not perform as good as the other methods, such as TQB-IRKA, BT and two-sided projection methods, and the one-sided projection method completely fails to capture the system dynamics for the input u_2 . This can also be observed from Table 4.3, where the mean relative errors of the outputs are reported.

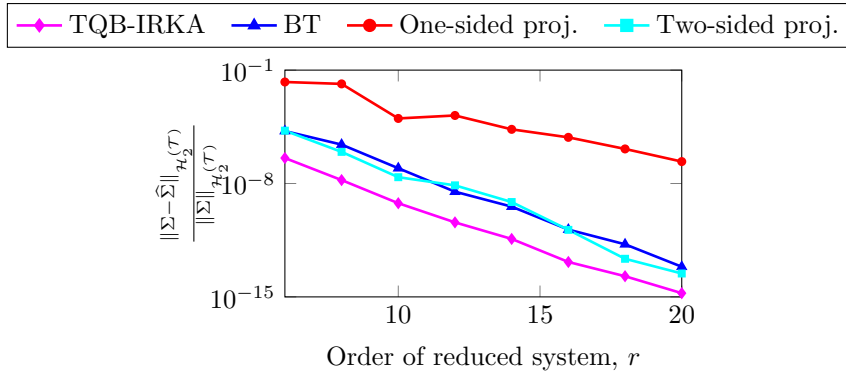
Further, we compute the quantities as defined in (4.1) using the reduced system of order $r = 10$ obtained upon convergence of TQB-IRKA and list them in Table 4.4. This also indicates that the obtained reduced-order system using TQB-IRKA satisfies all the optimality conditions (3.32) very accurately even though the nonlinear part of the system plays a significant role in the system dynamics.

FIGURE 4.5. An RC circuit: comparison of responses for the input $u_2 = 2.5(\sin(\pi t/5) + 1)$.

\mathcal{E}_C	\mathcal{E}_B	\mathcal{E}_N	\mathcal{E}_H	\mathcal{E}_λ
3.99×10^{-10}	4.68×10^{-8}	3.91×10^{-7}	3.37×10^{-8}	3.91×10^{-8}

TABLE 4.4

An RC circuit: perturbations to the optimality conditions.

FIGURE 4.6. An RC circuit: comparison of the truncated \mathcal{H}_2 -norm of the error system obtained via different methods of various orders.

Next, we also compare the truncated \mathcal{H}_2 -norm of the error system, i.e., $\|\Sigma - \widehat{\Sigma}\|_{\mathcal{H}_2(\mathcal{T})}$, in Figure 4.6, where the $\widehat{\Sigma}$ are constructed by various methods of different orders. The figure explains that TQB-IRKA yields better reduced-order systems with respect to the truncated \mathcal{H}_2 -norm.

At last, we mention here again that the reduced system obtained via POD retains the original exponential nonlinearities; therefore, we cannot use the same definition of the truncated \mathcal{H}_2 -norm as in Lemma 3.4 for such nonlinear systems. Hence, POD is omitted in Figure 4.6.

4.3. The FitzHugh-Nagumo (F-N) system. This example considers the F-N system, describing activation and deactivation dynamics of spiking neurons. This model is a simplification of the Hodgkin-Huxley neuron model. The dynamics of the

system is governed by the following nonlinear coupled PDEs:

$$\begin{aligned} \epsilon v_t(x, t) &= \epsilon^2 v_{xx}(x, t) + f(v(x, t)) - w(x, t) + q, \\ w_t(x, t) &= hv(x, t) - \gamma w(x, t) + q \end{aligned}$$

with a nonlinear function $f(v(x, t)) = v(v - 0.1)(1 - v)$ and initial and boundary conditions as follows:

$$\begin{aligned} v(x, 0) &= 0, & w(x, 0) &= 0, & x &\in (0, L), \\ v_x(0, t) &= i_0(t), & v_x(1, t) &= 0, & t &\geq 0, \end{aligned}$$

where $\epsilon = 0.015$, $h = 0.5$, $\gamma = 2$, $q = 0.05$, $L = 0.3$, and $i_0(t)$ is an actuator, acting as control input. The voltage and recovery voltage are denoted by v and w , respectively. MOR for this model has been considered in [8, 12, 18]. Furthermore, we also consider the same output as considered in [8, 12], which is the limit-cycle at the left boundary, i.e., $x = 0$. The system can be considered as having two inputs, namely q and $i_0(t)$; it has also two outputs, which are $v(0, t)$ and $w(0, t)$. This means that the system is a multi-input multi-output system as opposed to the two previous examples. We discretize the governing equations using a finite difference scheme. This leads to an ODE system, having cubic nonlinearity, which can then be transformed into the QB form. We consider $k = 300$ grid points, resulting in a QB system of order $3k = 900$.

We next determine reduced systems of order $r = 35$ using TQB-IRKA, BT, and POD. We choose the scaling factor $\gamma = 1$ in TQB-IRKA and it requires 26 iterations to converge. For this example, we also utilize the Kronecker structure of the Hessian to perform an efficient computation of the reduced Hessian. In order to apply POD, we first collect 500 snapshots of the original system for the time interval $(0, 10s]$ using $i_0(t) = 50(\sin(2\pi t) - 1)$ and then determine the projection based on the 35 dominant modes. The one-sided and two-sided subsystem-based interpolatory projection methods have major disadvantages in the MIMO QB case. The one-sided interpolatory projection approach of [22] can be applied to MIMO QB systems, however the dimension of the subspace V , and thus the dimension of the reduced model, increases quadratically due to the $V \otimes V$ term. As we mentioned in Section 1, two-sided interpolatory projection is only applicable to SISO QB systems. When the number of inputs and outputs are the same, which is the case in this example, one can still employ [9, Algo. 1] to construct a reduced system. This is exactly what we did here. However, it is important to note that even though the method can be applied numerically, it no longer ensures the theoretical subsystem interpolation property. Despite these drawbacks, for completeness of the comparison, we still construct reduced models using both one-sided and two-sided subsystem-based interpolatory projections.

Since the FHN system has two inputs and two outputs, each interpolation point yields 6 vectors in projection matrices V and W . Thus, in order to apply the two-sided projection, we use 6 linear \mathcal{H}_2 -optimal points and determine the reduced system of order 35 by taking the 35 dominant vectors. We do the same for the one-sided interpolatory projection method to compute the reduced-order system.

Next, we compare the quality of the reduced-order systems and plot the transient responses and the absolute errors of the outputs in Figure 4.7 for the training input $i_0(t) = 50(\sin(2\pi t) - 1)$.

As anticipated, POD provides a very good reduced-order system since the POD basis is constructed by using the same trajectory. Note that despite not reporting

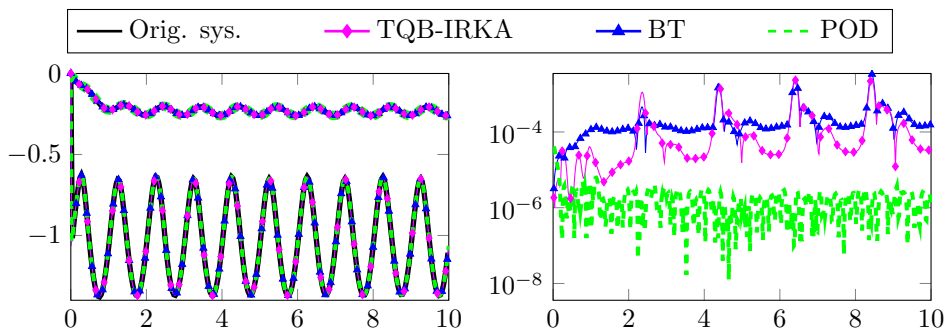


FIGURE 4.7. The FitzHugh-Nagumo system: comparison of the limit-cycle at the left boundary, $x = 0$ for $i_0(t) = 50(\sin(2\pi t) - 1)$.

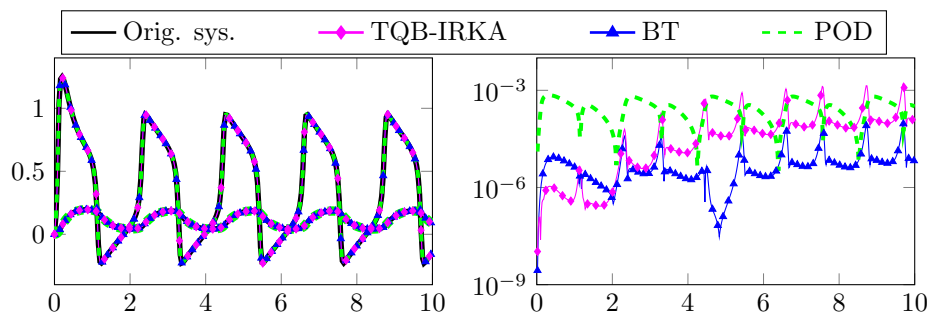


FIGURE 4.8. The FitzHugh-Nagumo system: comparison of the limit-cycle at the left boundary, $x = 0$ for $i_0(t) = 5 \cdot 10^4 t^3 \exp(-15t)$.

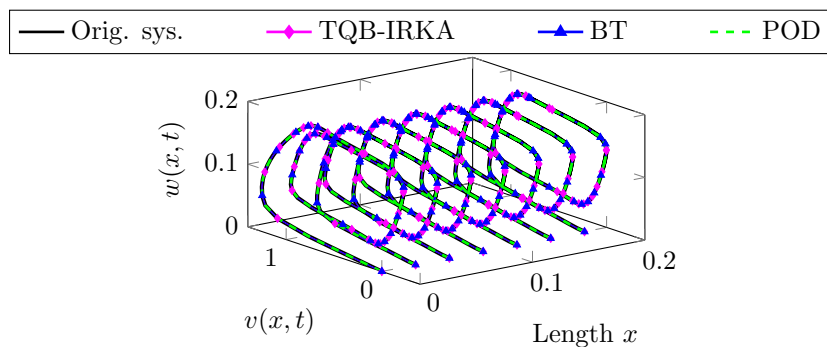


FIGURE 4.9. The FitzHugh-Nagumo system: limit-cycle behavior of the original and reduced-order systems in the spatial domain.

CPU times for the offline phases in this paper, due to the very different levels of the implementations used for the various methods, we would like to mention that in this example the construction of the POD basis with the fairly sophisticated MATLAB integrator `ode15s` takes roughly 1.5 more CPU time than constructing the TQB-IRKA reduced-order model with our vanilla implementation.

Between TQB-IRKA and BT, TQB-IRKA gives a marginally better reduced-order system as compared to BT for $i_0(t) = 50(\sin(2\pi t) - 1)$, but still both are

\mathcal{E}_C	\mathcal{E}_B	\mathcal{E}_N	\mathcal{E}_H	\mathcal{E}_λ
8.76×10^{-8}	7.35×10^{-9}	1.78×10^{-11}	4.27×10^{-9}	9.14×10^{-10}

TABLE 4.5

The FitzHugh-Nagumo system: perturbations to the optimality conditions.

very competitive. In contrast, the one-sided and two-sided interpolatory projection methods produce unstable reduced-order systems and are therefore omitted from the figures.

To test the robustness of the obtained reduced-order systems, we choose a different $i_0(t) = 5 \cdot 10^4 t^3 \exp(-15t)$ than the training one, and compare the transient responses in Figure 4.8. In the figure, we observe that BT performs the best among all methods for $i_0(t) = 5 \cdot 10^4 t^3 \exp(-15t)$, and POD and TQB-IRKA produce reduced-order systems of almost the same quality. One-sided and two-sided projection result in unstable reduced-order systems for $i_0(t) = 5 \cdot 10^4 t^3 \exp(-15t)$ as well. Furthermore, we also show the limit-cycles on the full space obtained from the original and reduced-order systems in Figure 4.9 for $i_0(t) = 5 \cdot 10^4 t^3 \exp(-15t)$, and observe that the reduced-order systems obtained using POD, TQB-IRKA, and BT, enable us to reproduce the limit-cycles, which is a typical neuronal dynamics as shown in Figures 4.7 and 4.9

As shown in [8], for particular interpolation points and higher-order moments, it might be possible to construct reduced-order systems via one-sided and two-sided interpolatory projection methods, which can reconstruct the limit-cycles. But as discussed in [8], stability of the reduced-order systems is highly sensitive to these specific choices and even slight modifications may lead to unstable systems. For the \mathcal{H}_2 linear optimal interpolation points selection we made here, the one-sided and two-sided approaches were not able to reproduce the limit-cycles; thus motivating the usage of TQB-IRKA and BT once again, especially for the MIMO case.

Moreover, we report how far the reduced system of order $r = 35$ due to TQB-IRKA is from satisfying the optimality conditions (3.32). For this, we compute the perturbations (4.1) and list them in Table 4.5. This clearly indicates that the reduced-order system almost satisfies all optimality conditions.

Lastly, we measure the truncated \mathcal{H}_2 -norm of the error systems, using the reduced-order systems obtained via different methods of various orders. We plot the relative truncated \mathcal{H}_2 -norm of the error systems in Figure 4.10. We observe that TQB-IRKA produces better reduced-order systems with respect to a truncated \mathcal{H}_2 -norm as compared to BT and one-sided projection. Furthermore, since we require stability of the matrix \hat{A} in the reduced QB system (1.2) to be able to compute a truncated \mathcal{H}_2 -norm of the error systems, we could not achieve this in the case of two-sided projection. For POD, we preserve the cubic nonlinearity in the reduced-order system; hence, the truncated \mathcal{H}_2 -norm definition in Lemma 3.4 does not apply for such systems. Thus, we cannot compute a truncated \mathcal{H}_2 -norm of the error system in the cases of the two-sided projection and POD, thereby these methods are not included in Figure 4.10.

5. Conclusions. In this paper, we have investigated the optimal model reduction problem for quadratic-bilinear control systems. We first have defined the \mathcal{H}_2 -norm for quadratic-bilinear systems based on the kernels of the underlying Volterra series and also proposed a truncated \mathcal{H}_2 -norm for the latter system class. We have then derived first-order necessary conditions, minimizing a truncated \mathcal{H}_2 -norm of the error system. These

necessary conditions, minimizing a truncated \mathcal{H}_2 -norm of the error system. These

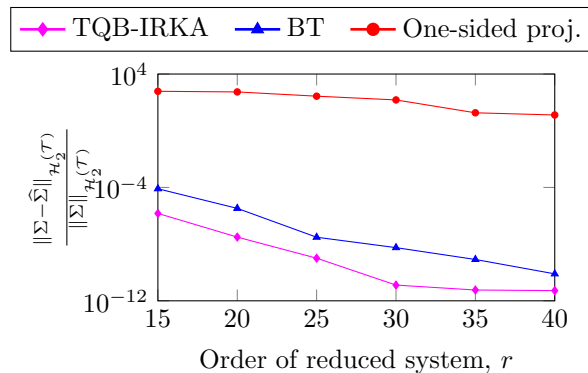


FIGURE 4.10. *The FitzHugh-Nagumo system: comparison of the truncated \mathcal{H}_2 -norm of the error system, having obtained reduced systems of different orders using various methods.*

optimality conditions lead to the proposed model reduction algorithm (TQB-IRKA), which iteratively constructs reduced order models that *approximately* satisfy the optimality conditions. We have also discussed the efficient computation of the reduced Hessian, utilizing the Kronecker structure of the Hessian of the QB system. Via several numerical examples, we have shown that TQB-IRKA outperforms the one-sided (including POD) and two-sided projection methods in most cases, and is comparable to balanced truncation. Especially for MIMO QB systems, TQB-IRKA and BT are the preferred methods of choice since the current framework of two-sided subspace interpolatory projection method is only applicable to SISO systems and the extension of the one-sided interpolatory projection method to MIMO QB systems yields reduced models whose dimension increases quadratically with the number of inputs. Moreover, our numerical experiments reveal that the reduced systems via TQB-IRKA and BT are more robust as compared to the one-sided and two-sided interpolatory projection methods in terms of stability of the reduced models although we do not have any theoretical justification of this observation yet.

As a future research topic, it would be interesting to use sophisticated tools from tensor theory to perform the computations related to Kronecker products efficiently and faster, thus accelerating the iteration steps in TQB-IRKA. Secondly, even though a stable random initialization has performed well in all of our numerical examples, a more educated initial guess could further improve the convergence of TQB-IRKA.

Acknowledgments. The authors would like to thank Dr. Tobias Breiten for providing the numerical examples, and MATLAB implementations for one-sided and two-sided interpolatory projection methods. We would also like to thank Dr. Patrick Kürschner for providing us MATLAB codes for solving Lyapunov equations via ADI. This research is supported by a research grant of the “International Max Planck Research School (IMPRS) for Advanced Methods in Process and System Engineering (Magdeburg)”. The work of S. Gugercin was supported in part by NSF through Grant DMS-1522616 and by the Alexander von Humboldt Foundation.

REFERENCES

- [1] A. C. ANTOUNAS, *Approximation of Large-Scale Dynamical Systems*, SIAM Publications, Philadelphia, PA, 2005.
- [2] P. ASTRID, S. WEILAND, K. WILLCOX, AND T. BACKX, *Missing point estimation in models*

- described by proper orthogonal decomposition, IEEE Trans. Autom. Control, 53 (2008), pp. 2237–2251.
- [3] Z. BAI, *Krylov subspace techniques for reduced-order modeling of large-scale dynamical systems*, Appl. Numer. Math., 43 (2002), pp. 9–44.
 - [4] M. BARRAULT, Y. MADAY, N. C. NGUYEN, AND A. T. PATERA, *An ‘empirical interpolation’ method: application to efficient reduced-basis discretization of partial differential equations*, C. R. Math. Acad. Sci. Paris, 339 (2004), pp. 667–672.
 - [5] R. H. BARTELS AND G. W. STEWART, *Solution of the matrix equation $AX + XB = C$* , Comm. ACM, 15 (1972), pp. 820–826.
 - [6] U. BAUR, P. BENNER, AND L. FENG, *Model order reduction for linear and nonlinear systems: a system-theoretic perspective*, Arch. Comput. Methods Eng., 21 (2014), pp. 331–358.
 - [7] P. BENNER AND T. BREITEN, *Interpolation-based \mathcal{H}_2 -model reduction of bilinear control systems*, SIAM J. Matrix Anal. Appl., 33 (2012), pp. 859–885.
 - [8] P. BENNER AND T. BREITEN, *Two-sided moment matching methods for nonlinear model reduction*, Preprint MPIMD/12-12, MPI Magdeburg, 2012. Available from <http://www.mpi-magdeburg.mpg.de/preprints/>.
 - [9] P. BENNER AND T. BREITEN, *Two-sided projection method for nonlinear model reduction*, SIAM J. Sci. Comput., 37 (2015), pp. B239–B260, doi:10.1137/14097255X.
 - [10] P. BENNER AND T. DAMM, *Lyapunov equations, energy functionals, and model order reduction of bilinear and stochastic systems*, SIAM J. Cont. Optim., 49 (2011), pp. 686–711, doi:10.1137/09075041X.
 - [11] P. BENNER AND P. GOYAL, *Multipoint interpolation of Volterra series and \mathcal{H}_2 -model reduction for a family of bilinear descriptor systems*, Systems Control Lett., 97 (2016), pp. 1–11.
 - [12] P. BENNER AND P. GOYAL, *Balanced truncation model order reduction for quadratic-bilinear control systems*, arXiv e-prints 1705.00160, 2017, <https://arxiv.org/abs/1705.00160>. math.OC.
 - [13] P. BENNER, P. KÜRSCHNER, AND J. SAAK, *Self-generating and efficient shift parameters in ADI methods for large Lyapunov and Sylvester equations*, Electron. Trans. Numer. Anal., 43 (2014), pp. 142–162.
 - [14] P. BENNER, V. MEHRMANN, AND D. C. SORESENSEN, *Dimension Reduction of Large-Scale Systems*, vol. 45 of Lect. Notes Comput. Sci. Eng., Springer-Verlag, Berlin/Heidelberg, Germany, 2005.
 - [15] P. BENNER AND J. SAAK, *Numerical solution of large and sparse continuous time algebraic matrix Riccati and Lyapunov equations: a state of the art survey*, GAMM-Mitt., 36 (2013), pp. 32–52.
 - [16] P. BENNER, E. SACHS, AND S. VOLKWEIN, *Model order reduction for PDE constrained optimization*, in Trends in PDE Constrained Optimization, G. Leugering, P. Benner, S. Engell, A. Griewank, H. Harbrecht, M. Hinze, R. Rannacher, and S. Ulbrich, eds., vol. 165 of International Series of Numerical Mathematics, Springer International Publishing, 2014, pp. 303–326, doi:10.1007/978-3-319-05083-6_19.
 - [17] N. CHAFEE AND E. F. INFANTE, *A bifurcation problem for a nonlinear partial differential equation of parabolic type*, Appl. Anal., 4 (1974), pp. 17–37.
 - [18] S. CHATURANTABUT AND D. C. SORESENSEN, *Nonlinear model reduction via discrete empirical interpolation*, SIAM J. Sci. Comput., 32 (2010), pp. 2737–2764.
 - [19] M. CONDON AND R. IVANOV, *Nonlinear systems-algebraic gramians and model reduction*, COMPEL, 24 (2005), pp. 202–219.
 - [20] G. FLAGG AND S. GUGERCIN, *Multipoint Volterra series interpolation and \mathcal{H}_2 optimal model reduction of bilinear systems*, SIAM J. Matrix Anal. Appl., 36 (2015), pp. 549–579.
 - [21] M. A. GREPL, Y. MADAY, N. C. NGUYEN, AND A. T. PATERA, *Efficient reduced-basis treatment of nonaffine and nonlinear partial differential equations*, ESAIM: Math. Model. Numer. Anal., 41 (2007), pp. 575–605.
 - [22] C. GU, *QLMOR: A projection-based nonlinear model order reduction approach using quadratic-linear representation of nonlinear systems*, IEEE Trans. Comput.-Aided Design Integr. Circuits Syst., 30 (2011), pp. 1307–1320.
 - [23] S. GUGERCIN, A. C. ANTOULAS, AND C. A. BEATTIE, *\mathcal{H}_2 model reduction for large-scale dynamical systems*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 609–638.
 - [24] S. GUGERCIN, T. STYKEL, AND S. WYATT, *Model reduction of descriptor systems by interpolatory projection methods*, SIAM J. Sci. Comput., 35 (2013), pp. B1010–B1033, doi:10.1137/130906635.
 - [25] E. HANSEN, F. KRAMER, AND A. OSTERMANN, *A second-order positivity preserving scheme for semilinear parabolic problems*, Appl. Numer. Math., 62 (2012), pp. 1428–1435.
 - [26] H. V. HENDERSON AND S. R. SEARLE, *The vec-permutation matrix, the vec operator and Kro-*

- necker products: A review*, Linear and Multilinear Algebra, 9 (1981), pp. 271–288.
- [27] M. HINZE AND S. VOLKWEIN, *Proper orthogonal decomposition surrogate models for nonlinear dynamical systems: Error estimates and suboptimal control*, in [14], pp. 261–306.
- [28] M. HINZE AND S. VOLKWEIN, *Error estimates for abstract linear-quadratic optimal control problems using proper orthogonal decomposition*, Comput. Optim. Appl., 39 (2008), pp. 319–345.
- [29] R. A. HORN AND C. R. JOHNSON, *Topics in Matrix Analysis*, Cambridge University Press, Cambridge.
- [30] M. KÖHLER, *On the closest stable descriptor system in the respective spaces RH_2 and RH_∞* , Linear Algebra Appl., 443 (2014), pp. 34–49.
- [31] T. G. KOLDA AND B. W. BADER, *Tensor decompositions and applications*, SIAM Rev., 51 (2009), pp. 455–500.
- [32] K. KUNISCH AND S. VOLKWEIN, *Proper orthogonal decomposition for optimality systems*, ESAIM: Math. Model. Numer. Anal., 42 (2008), pp. 1–23.
- [33] P. KUNKEL AND V. MEHRMANN, *Differential-algebraic equations: analysis and numerical solution*, European Mathematical Society, 2006.
- [34] P. LI AND L. T. PILEGGI, *Compact reduced-order modeling of weakly nonlinear analog and RF circuits*, IEEE Trans. Comput.-Aided Design Integr. Circuits Syst., 24 (2005), pp. 184–203.
- [35] B. C. MOORE, *Principal component analysis in linear systems: controllability, observability, and model reduction*, IEEE Trans. Autom. Control, AC-26 (1981), pp. 17–32.
- [36] N. NGUYEN, A. T. PATERA, AND J. PERAIRE, *A best points interpolation method for efficient approximation of parametrized functions*, Internat. J. Numer. Methods Engrg., 73 (2008), pp. 521–543.
- [37] J. R. PHILLIPS, *Projection frameworks for model reduction of weakly nonlinear systems*, in Proc. of the 2000 Design Automation Conference, 2000, pp. 184–189.
- [38] J. R. PHILLIPS, *Projection-based approaches for model reduction of weakly nonlinear, time-varying systems*, IEEE Trans. Comput.-Aided Design Integr. Circuits Syst., 22 (2003), pp. 171–187.
- [39] M. J. REWIEŃSKI, *A Trajectory Piecewise-Linear Approach to Model Order Reduction of Nonlinear Dynamical Systems*, PhD thesis, Massachusetts Institute of Technology, 2003.
- [40] W. J. RUGH, *Nonlinear System Theory*, The Johns Hopkins University Press, Baltimore, MD, 1981.
- [41] W. H. A. SCHILDERS, H. A. VAN DER VORST, AND J. ROMMES, *Model Order Reduction: Theory, Research Aspects and Applications*, Springer-Verlag, Berlin, Heidelberg, 2008.
- [42] V. SIMONCINI, *Computational methods for linear matrix equations*, SIAM Rev., 58 (2016), pp. 377–441.
- [43] L. ZHANG AND J. LAM, *On H_2 model reduction of bilinear systems*, Automatica, 38 (2002), pp. 205–216.

Appendix A. Important relations of the Kronecker products. In this section, we provide some relations between Kronecker products, which will remarkably simplify the optimality conditions in [Appendix B](#).

LEMMA A.1. [7, Lemma A.1] Consider $f(x) \in \mathbb{R}^{s \times n}$, $A(y) \in \mathbb{R}^{n \times n}$, $G \in \mathbb{R}^{n \times q}$ with $x, y \in \mathbb{R}$, and let $\mathcal{L}(y)$ be defined as

$$\mathcal{L}(y) = -A(y) \otimes I_n - I_n \otimes A(y).$$

If the functions f and A are differentiable with respect to x and y , respectively, then

$$\begin{aligned} \frac{\partial}{\partial x} [(\mathcal{I}_s)^T (f(x) \otimes f(x)) \mathcal{L}^{-1}(y) (G \otimes G) \mathcal{I}_q] \\ = 2(\mathcal{I}_s)^T \left(\left(\frac{\partial}{\partial x} f(x) \right) \otimes f(x) \right) \mathcal{L}^{-1}(y) (G \otimes G) \mathcal{I}_q. \end{aligned}$$

Moreover, let $X, Y \in \mathbb{R}^{n \times n}$ be symmetric matrices. Then,

$$\frac{\partial}{\partial y} \left[\text{vec}(X)^T \mathcal{L}^{-1}(y) \text{vec}(Y) \right] = 2 \cdot \text{vec}(X)^T \mathcal{L}^{-1}(y) \left(\frac{\partial}{\partial y} A(y) \otimes I_n \right) \mathcal{L}^{-1}(y) \text{vec}(Y).$$

LEMMA A.2. Let \mathcal{F} , $\widehat{\mathcal{F}}$ be defined as follows:

$$\mathcal{F} = [I_n \ \mathbf{0}] \otimes [I_n \ \mathbf{0}] \quad \text{and} \quad \widehat{\mathcal{F}} = [\mathbf{0} \ I_r] \otimes [\mathbf{0} \ I_r],$$

and consider a permutation matrix

$$(A.1) \quad M = \begin{bmatrix} M_{nnr} & \mathbf{0} \\ \mathbf{0} & M_{rnr} \end{bmatrix},$$

where M_{pqr} is defined in (3.28). Moreover, let the two column vectors x and y be partitioned as

$$x = [x_1^T \ x_2^T \ x_3^T \ x_4^T]^T \quad \text{and} \quad y = [y_1^T \ y_2^T \ y_3^T \ y_4^T]^T,$$

where $x_1, y_1 \in \mathbb{R}^{n^2}$, $x_{\{2,3\}}, y_{\{2,3\}} \in \mathbb{R}^{nr}$, and $x_4, y_4 \in \mathbb{R}^{r^2}$. Then, the following relations hold:

$$(A.2) \quad (\widehat{\mathcal{F}} \otimes \mathcal{F})T_{(n+r, n+r)}(M \otimes M)(x \otimes y) = T_{(n,r)}(x_3 \otimes y_3), \quad \text{and}$$

$$(A.3) \quad (\widehat{\mathcal{F}} \otimes \widehat{\mathcal{F}})T_{(n+r, n+r)}(M \otimes M)(x \otimes y) = T_{(r,r)}(x_4 \otimes y_4),$$

where $T_{(n,m)}$ is also a permutation matrix given by

$$T_{(n,m)} = I_m \otimes [I_m \otimes e_1^n, \dots, I_m \otimes e_n^n] \otimes I_n.$$

Proof. Let us begin by considering the following equation:

$$(\widehat{\mathcal{F}} \otimes \mathcal{F})T_{(n+r, n+r)} = [\mathbf{0} \ I_r \otimes [\mathbf{0} \ I_r] \otimes \mathcal{F}] (I_{n+r} \otimes \mathcal{G}),$$

where $\mathcal{G} = [I_{n+r} \otimes e_1^{n+r}, \dots, I_{n+r} \otimes e_{n+r}^{n+r}] \otimes I_{n+r}$. Next, we split I_{n+r} as $I_{n+r} = \begin{bmatrix} I_n & \mathbf{0} \\ \mathbf{0} & I_r \end{bmatrix}$, leading to

$$(A.4) \quad \begin{aligned} (\widehat{\mathcal{F}} \otimes \mathcal{F})T_{(n+r, n+r)} &= [\mathbf{0} \ I_r \otimes [\mathbf{0} \ I_r] \otimes \mathcal{F}] \begin{bmatrix} I_n \otimes \mathcal{G} & \mathbf{0} \\ \mathbf{0} & I_r \otimes \mathcal{G} \end{bmatrix} \\ &= [\mathbf{0} \ (I_r \otimes [\mathbf{0} \ I_r] \otimes \mathcal{F}) (I_r \otimes \mathcal{G})] \\ &= [\mathbf{0} \ I_r \otimes (([\mathbf{0} \ I_r] \otimes \mathcal{F}) \mathcal{G})]. \end{aligned}$$

Now, we investigate the following equation (a component of the previous equation):

$$([\mathbf{0} \ I_r] \otimes \mathcal{F}) \mathcal{G}_i =: \mathcal{L}_i,$$

where \mathcal{G}_i is i th block column of the matrix \mathcal{G} given by $\mathcal{G}_i = I_{n+r} \otimes e_i^{n+r} \otimes I_{n+r}$. This yields

$$\begin{aligned} \mathcal{L}_i &= ([\mathbf{0} \ I_r] \otimes \mathcal{F}) (I_{n+r} \otimes e_i^{n+r} \otimes I_{n+r}) \\ &= ([\mathbf{0} \ I_r] I_{n+r}) \otimes (\mathcal{F}(e_i^{n+r} \otimes I_{n+r})) = [\mathbf{0} \ I_r] \otimes (\mathcal{F}(e_i^{n+r} \otimes I_{n+r})). \end{aligned}$$

Assuming that $1 \leq i \leq n$, we can write \mathcal{L}_i as

$$\begin{aligned} \mathcal{L}_i &= [\mathbf{0} \ I_r] \otimes \left(\mathcal{F} \left(\begin{bmatrix} e_i^n \\ \mathbf{0} \end{bmatrix} \otimes I_{n+r} \right) \right) = [\mathbf{0} \ I_r] \otimes \left([I_n \otimes [I_n \ \mathbf{0}] \ \mathbf{0}] \begin{bmatrix} e_i^n \otimes I_{n+r} \\ \mathbf{0} \end{bmatrix} \right) \\ &= [\mathbf{0} \ I_r] \otimes [e_i^n \otimes [I_n \ \mathbf{0}]] = [\mathbf{0} \ I_r \otimes (e_i^n \otimes [I_n \ \mathbf{0}])]. \end{aligned}$$

Subsequently, we assume $n + r \geq i > n$, which leads to

$$\begin{aligned}\mathcal{L}_i &= [\mathbf{0} \quad I_r] \otimes \left(\mathcal{F} \left(\begin{bmatrix} \mathbf{0} \\ e_{i-n}^r \end{bmatrix} \otimes I_{n+r} \right) \right) \\ &= [\mathbf{0} \quad I_r] \otimes \left([I_n \otimes [I_n \quad \mathbf{0}] \quad \mathbf{0}] \begin{bmatrix} e_{i-n}^r \otimes I_{n+r} \\ \mathbf{0} \end{bmatrix} \right) = \mathbf{0}.\end{aligned}$$

Thus,

$$([\mathbf{0} \quad I_r] \otimes \mathcal{F}) \mathcal{G} = [\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n, \mathbf{0}] =: \mathcal{L}.$$

Inserting the above expression in (A.4) yields

$$(\widehat{\mathcal{F}} \otimes \mathcal{F}) T_{(n+r, n+r)} = [\mathbf{0} \quad I_r \otimes \mathcal{L}].$$

Now, we are ready to investigate the following term:

$$\begin{aligned}(\widehat{\mathcal{F}} \otimes \mathcal{F}) T_{(n+r, n+r)}(M \otimes M) &= [\mathbf{0} \quad I_r \otimes \mathcal{L}] \begin{bmatrix} M_{nnr} \otimes M & \mathbf{0} \\ \mathbf{0} & M_{rnr} \otimes M \end{bmatrix} \\ &= [\mathbf{0} \quad I_r \otimes \mathcal{L}] \begin{bmatrix} M_{nnr} \otimes M & \mathbf{0} \\ \mathbf{0} & M_{rnr} \otimes M \end{bmatrix} \\ &= [\mathbf{0} \quad (I_r \otimes \mathcal{L})(M_{rnr} \otimes M)].\end{aligned}$$

Further, we consider the second block column of the above relation and substitute for M_{nnr} and M_{rnr} using (3.28) to get

$$\begin{aligned}(\text{A.5}) \quad (I_r \otimes \mathcal{L})(M_{rnr} \otimes M) &= (I_r \otimes \mathcal{L}) \left[I_r \otimes \begin{bmatrix} I_n \\ \mathbf{0} \end{bmatrix} \otimes M \quad I_r \otimes \begin{bmatrix} \mathbf{0} \\ I_r \end{bmatrix} \otimes M \right] \\ &= \left[(I_r \otimes \mathcal{L}) \left(I_r \otimes \begin{bmatrix} I_n \\ \mathbf{0} \end{bmatrix} \otimes M \right) \quad (I_r \otimes \mathcal{L}) \left(I_r \otimes \begin{bmatrix} \mathbf{0} \\ I_r \end{bmatrix} \otimes M \right) \right].\end{aligned}$$

Our following task is to examine each block column of (A.5). We begin with the first block; this is

$$\begin{aligned}(I_r \otimes \mathcal{L}) \left(I_r \otimes \begin{bmatrix} I_n \\ \mathbf{0} \end{bmatrix} \otimes M \right) &= I_r \otimes \left(\mathcal{L} \left(\begin{bmatrix} I_n \\ \mathbf{0} \end{bmatrix} \otimes M \right) \right) = I_r \otimes \left(\mathcal{L} \begin{bmatrix} I_n \otimes M \\ \mathbf{0} \end{bmatrix} \right) \\ &= I_r \otimes [\mathcal{L}_1 M, \dots, \mathcal{L}_n M].\end{aligned}$$

We next aim to simplify the term $\mathcal{L}_i M$, which appears in the previous equation:

$$\begin{aligned}(\text{A.6}) \quad \mathcal{L}_i M &= [\mathbf{0} \quad I_r \otimes [e_j^n \otimes [I_n \quad \mathbf{0}]]] \begin{bmatrix} M_{nnr} & \mathbf{0} \\ \mathbf{0} & M_{rnr} \end{bmatrix} \\ &= [\mathbf{0} \quad (I_r \otimes [e_j^n \otimes [I_n \quad \mathbf{0}]]) M_{rnr}] \\ &= \left[\mathbf{0} \quad (I_r \otimes e_j^n \otimes [I_n \quad \mathbf{0}]) \begin{bmatrix} I_r \otimes \begin{bmatrix} I_n \\ \mathbf{0} \end{bmatrix} & I_r \otimes \begin{bmatrix} \mathbf{0} \\ I_r \end{bmatrix} \end{bmatrix} \right] \\ &= [\mathbf{0} \quad (I_r \otimes e_j^n \otimes I_n) \quad \mathbf{0}] := \mathcal{X}_i.\end{aligned}$$

The second block column of (A.5) can be studied in a similar fashion, and it can be shown that

$$(I_r \otimes \mathcal{L}) \left(I_r \otimes \begin{bmatrix} \mathbf{0} \\ I_r \end{bmatrix} \otimes M \right) = \mathbf{0}.$$

Summing up all these expressions, we obtain

$$(\widehat{\mathcal{F}} \otimes \mathcal{F})T_{(n+r,n+r)}(M \otimes M) = \begin{bmatrix} \mathbf{0} & (I_r \otimes [\mathcal{X}_1, \dots, \mathcal{X}_n]) & \mathbf{0} \end{bmatrix},$$

where \mathcal{X}_i is defined in (A.6). This gives

$$(A.7) \quad (\widehat{\mathcal{F}} \otimes \mathcal{F})T_{(n+r,n+r)}(M \otimes M)(x \otimes y) = \begin{bmatrix} \mathbf{0} & I_r \otimes [\mathcal{X}_1, \dots, \mathcal{X}_n] & \mathbf{0} \end{bmatrix} (x \otimes y) \\ = (I_r \otimes [\mathcal{X}_1, \dots, \mathcal{X}_n]) (x_3 \otimes y).$$

Next, we define another permutation

$$\mathcal{Q} = \left[\underbrace{I_r \otimes I_n \otimes \begin{bmatrix} I_{n^2} \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}}_{\mathcal{Q}_1} \quad \underbrace{I_r \otimes I_n \otimes \begin{bmatrix} \mathbf{0} \\ I_{nr} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}}_{\mathcal{Q}_2} \quad \underbrace{I_r \otimes I_n \otimes \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ I_{nr} \\ 0 \end{bmatrix}}_{\mathcal{Q}_3} \quad \underbrace{I_r \otimes I_n \otimes \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \\ I_{r^2} \end{bmatrix}}_{\mathcal{Q}_4} \right],$$

which allows us to write

$$(x_3 \otimes y) = \mathcal{Q} \begin{bmatrix} x_3 \otimes y_1 \\ x_3 \otimes y_2 \\ x_3 \otimes y_3 \\ x_3 \otimes y_4 \end{bmatrix}.$$

Substituting this into (A.7) results in

$$(\widehat{\mathcal{F}} \otimes \mathcal{F})T_{(n+r,n+r)}(M \otimes M)(x \otimes y) \\ = (I_r \otimes [\mathcal{X}_1, \dots, \mathcal{X}_n]) \begin{bmatrix} \mathcal{Q}_1 & \mathcal{Q}_2 & \mathcal{Q}_3 & \mathcal{Q}_4 \end{bmatrix} \begin{bmatrix} x_3 \otimes y_1 \\ x_3 \otimes y_2 \\ x_3 \otimes y_3 \\ x_3 \otimes y_4 \end{bmatrix}.$$

Now, it can be easily verified that $(I_r \otimes [\mathcal{X}_1, \dots, \mathcal{X}_n]) \begin{bmatrix} \mathcal{Q}_1 & \mathcal{Q}_2 & \mathcal{Q}_4 \end{bmatrix} = 0$. Thus, we obtain

$$(\widehat{\mathcal{F}} \otimes \mathcal{F})T_{(n+r,n+r)}(M \otimes M)(x \otimes y) = (I_r \otimes [\mathcal{X}_1, \dots, \mathcal{X}_n]) \mathcal{Q}_3 (x_3 \otimes y_3) \\ = (I_r \otimes [\mathcal{X}_1, \dots, \mathcal{X}_n]) \left(I_r \otimes I_n \otimes \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ I_{nr} \\ \mathbf{0} \end{bmatrix} \right) (x_3 \otimes y_3) \\ = (I_r \otimes [I_r \otimes e_1^n \otimes I_n, \dots, I_r \otimes e_1^n \otimes I_n]) (x_3 \otimes y_3) = T_{(n,r)}(x_3 \otimes y_3).$$

One can prove the relation (A.2) in a similar manner. However, for the brevity of the paper, we omit it. This concludes the proof. \square

We will find similar expressions as (A.2) and (A.3) in Appendix B, where we then make use of Lemma A.2 to simplify them.

Appendix B. Proof of Theorem 3.7.

Optimality conditions with respect to \tilde{C} . We start with deriving the optimality conditions by taking the derivative of the error functional \mathcal{E} (3.26) with respect

to \tilde{C} . By using [Lemma A.1](#), we obtain

$$\begin{aligned} \frac{\partial \mathcal{E}^2}{\partial \tilde{C}_{ij}} &= 2(\mathcal{I}_p)^T \left([\mathbf{0} \quad -e_i^p(e_j^r)^T] \otimes \tilde{C}^e \right) \left(-\tilde{A}^e \otimes I_{n+r} - I_{n+r} \otimes \tilde{A}^e \right)^{-1} \\ &\quad \left(\left(\tilde{B}^e \otimes \tilde{B}^e \right) \mathcal{I}_m + \sum_{k=1}^m \left(\tilde{N}_k^e \otimes \tilde{N}_k^e \right) \mathcal{P}_l + \left(\tilde{H}^e \otimes \tilde{H}^e \right) T_{(n+r, n+r)}(\mathcal{P}_l \otimes \mathcal{P}_l) \right), \end{aligned}$$

where \mathcal{P}_l is defined in [\(3.27\)](#). On simplification, we get

$$\begin{aligned} \frac{\partial \mathcal{E}^2}{\partial \tilde{C}_{ij}} &= 2(\mathcal{I}_p)^T \left(-e_i^p(e_j^r)^T \otimes \tilde{C}^e \right) \left(-\Lambda \otimes I_{n+r} - I_r \otimes \tilde{A}^e \right)^{-1} \left(\left(\tilde{B} \otimes \tilde{B}^e \right) \mathcal{I}_m \right. \\ &\quad \left. + \sum_{k=1}^m \left(\tilde{N}_k \otimes \tilde{N}_k^e \right) \mathcal{P}_l^{(2)} + \left(\tilde{H} \hat{\mathcal{F}} \otimes \tilde{H}^e \right) T_{(n+r, n+r)}(\mathcal{P}_l \otimes \mathcal{P}_l) \right), \\ &= 2(\mathcal{I}_p)^T \left(-e_i^p(e_j^r)^T \otimes \tilde{C}^e \right) \left(M_{rnr} (-\mathcal{J}_\Lambda - \mathcal{J}_A) M_{rnr}^T \right)^{-1} \left(\left(\tilde{B} \otimes \tilde{B}^e \right) \mathcal{I}_m \right. \\ \text{(B.1)} \quad &\left. + \sum_{k=1}^m \left(\tilde{N}_k \otimes \tilde{N}_k^e \right) \mathcal{P}_l^{(2)} + \left(\tilde{H} \hat{\mathcal{F}} \otimes \tilde{H}^e \right) T_{(n+r, n+r)}(\mathcal{P}_l \otimes \mathcal{P}_l) \right), \end{aligned}$$

where

$$\mathcal{J}_\Lambda = \begin{bmatrix} \Lambda \otimes I_n & \mathbf{0} \\ \mathbf{0} & \Lambda \otimes I_r \end{bmatrix}, \quad \mathcal{J}_A = \begin{bmatrix} I_r \otimes A & \mathbf{0} \\ \mathbf{0} & I_r \otimes \Lambda \end{bmatrix}, \quad \text{and}$$

$\mathcal{P}_l^{(2)}$ is the lower block row of \mathcal{P}_l as shown in [\(3.27\)](#). Furthermore, since M_{rnr} is a permutation matrix, this implies $M_{rnr} M_{rnr}^T = I$. Using this relation in [\(B.1\)](#), we obtain

$$\begin{aligned} \frac{\partial \mathcal{E}^2}{\partial \tilde{C}_{ij}} &= 2(\mathcal{I}_p)^T \left(-e_i^p(e_j^r)^T \otimes [C \quad -\tilde{C}] \right) M_{rnr} (-\mathcal{J}_\Lambda - \mathcal{J}_A)^{-1} \left(M_{rnr}^T \left(\tilde{B} \otimes \tilde{B}^e \right) \mathcal{I}_m \right. \\ &\quad \left. + M_{rnr}^T \sum_{k=1}^m \left(\tilde{N}_k \otimes \tilde{N}_k^e \right) \mathcal{P}_1^{(2)} + M_{rnr}^T \left(\tilde{H} \hat{\mathcal{F}} \otimes \tilde{H}^e \right) T_{(n+r, n+r)}(\mathcal{P}_l \otimes \mathcal{P}_l) \right) \\ &= 2(\mathcal{I}_p)^T \left(\left[-e_i^p(e_j^r)^T \otimes C \quad e_i e_j^T \otimes \tilde{C} \right] \right) (-\mathcal{J}_\Lambda - \mathcal{J}_A)^{-1} \left(\left[\begin{array}{c} \tilde{B} \otimes B \\ \tilde{B} \otimes \tilde{B} \end{array} \right] \mathcal{I}_m \right. \\ &\quad \left. + \sum_{k=1}^m \left[\begin{array}{cc} \tilde{N}_k \otimes N_k & \mathbf{0} \\ \mathbf{0} & \tilde{N}_k \otimes \tilde{N}_k \end{array} \right] M_{rnr}^T \mathcal{P}_1^{(2)} \right. \\ \text{(B.2)} \quad &\left. + \left[\left(\tilde{H} \hat{\mathcal{F}} \otimes H \mathcal{F} \right) T_{(n+r, n+r)}(M \otimes M)(M^T \otimes M^T)(\mathcal{P}_l \otimes \mathcal{P}_l) \right] \right), \end{aligned}$$

where M is the permutation matrix defined in [\(A.1\)](#). The multiplication of M^T and \mathcal{P}_l yields

$$M^T \mathcal{P}_l = \begin{bmatrix} M_{nnr} \mathcal{P}_l^{(1)} \\ M_{rnr} \mathcal{P}_l^{(2)} \end{bmatrix} = [p_1^T \quad p_2^T \quad p_3^T \quad p_4^T]^T =: \tilde{\mathcal{P}}_l,$$

where

(B.3)

$$\begin{aligned} p_1 &= (-A \otimes I_n - I_n \otimes A)^{-1} (B \otimes B) \mathcal{I}_m, & p_2 &= (-A \otimes I_r - I_n \otimes \Lambda)^{-1} (B \otimes \tilde{B}) \mathcal{I}_m, \\ p_3 &= (-\Lambda \otimes I_n - I_r \otimes A)^{-1} (\tilde{B} \otimes B) \mathcal{I}_m, & p_4 &= (-\Lambda \otimes I_r - I_r \otimes \Lambda)^{-1} (\tilde{B} \otimes \tilde{B}) \mathcal{I}_m. \end{aligned}$$

Moreover, note that $p_3 = \text{vec}(V_1)$, where V_1 solves (3.29a). Applying the result of Lemma A.2 in (B.2) yields

$$\begin{aligned} \frac{\partial \mathcal{E}^2}{\partial \tilde{C}_{ij}} &= 2(\mathcal{I}_p)^T (e_i^p (e_j^r)^T \otimes C) (-\Lambda \otimes I_n - I_r \otimes A)^{-1} \left((\tilde{B} \otimes B) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes N_k) p_3 \right. \\ &\quad \left. + (\tilde{H} \otimes H) T_{(n,r)} p_3 \otimes p_3 \right) - 2(\mathcal{I}_p)^T \left(e_i^p (e_j^r)^T \otimes \tilde{C} \right) \left(-\Lambda \otimes I_n - I_r \otimes \tilde{A} \right)^{-1} \\ &\quad \times \left((\tilde{B} \otimes \tilde{B}) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes \tilde{N}_k) p_4 + (\tilde{H} \otimes \tilde{H}) T_{(r,r)} (p_4 \otimes p_4) \right) \\ &= 2(\mathcal{I}_p)^T (e_i^p (e_j^r)^T \otimes C) (-\Lambda \otimes I_n - I_r \otimes A)^{-1} \left((\tilde{B} \otimes B) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes N_k) p_3 \right. \\ &\quad \left. + (\tilde{H} \otimes H) T_{(n,r)} p_3 \otimes p_3 \right) - 2(\mathcal{I}_p)^T \left(e_i^p (e_j^r)^T \otimes \hat{C} \right) \left(-\Lambda \otimes I_n - I_r \otimes \hat{A} \right)^{-1} \\ \text{(B.4)} &\quad \times \left((\tilde{B} \otimes \hat{B}) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes \hat{N}_k) \hat{p}_4 + (\tilde{H} \otimes \hat{H}) T_{(r,r)} (\hat{p}_4 \otimes \hat{p}_4) \right), \end{aligned}$$

where $\hat{p}_4 = (-\Lambda \otimes I_r - I_r \otimes \hat{A})^{-1} (\tilde{B} \otimes \hat{B}) \mathcal{I}_m = \text{vec}(\hat{V}_1)$, where \hat{V}_1 is as defined in (3.31). Setting (B.4) equal to zero results in a necessary condition with respect to \tilde{C} as follows:

$$\begin{aligned} &(\mathcal{I}_p)^T (e_i^p (e_j^r)^T \otimes C) (-\Lambda \otimes I_n - I_r \otimes A)^{-1} \left((\tilde{B} \otimes B) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes N_k) p_3 \right. \\ &\quad \left. + (\tilde{H} \otimes H) T_{(n,r)} (p_3 \otimes p_3) \right) \\ \text{(B.5)} &= (\mathcal{I}_p)^T \left(e_i^p (e_j^r)^T \otimes \hat{C} \right) \left(-\Lambda \otimes I_n - I_r \otimes \hat{A} \right)^{-1} \\ &\quad \times \left((\tilde{B} \otimes \hat{B}) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes \hat{N}_k) \hat{p}_4 + (\tilde{H} \otimes \hat{H}) T_{(r,r)} (\hat{p}_4 \otimes \hat{p}_4) \right). \end{aligned}$$

Now, we first manipulate the left-side of the above equation (B.5). Using Lemma 2.4 and (2.1), we get

$$\begin{aligned} &(\mathcal{I}_p)^T (e_i^p (e_j^r)^T \otimes C) (-\Lambda \otimes I_n - I_r \otimes A)^{-1} \left((\tilde{B} \otimes B) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes N_k) p_3 \right. \\ &\quad \left. + (\tilde{H} \otimes H) T_{(n,r)} (p_3 \otimes p_3) \right) \\ &= (\mathcal{I}_p)^T (e_i^p (e_j^r)^T \otimes C) (-\Lambda \otimes I_n - I_r \otimes A)^{-1} \left(\text{vec} \left(B \tilde{B}^T \right) + \sum_{k=1}^m \text{vec} \left(N_k V_1 \tilde{N}_k^T \right) \right. \\ &\quad \left. + (\tilde{H} \otimes H) \text{vec} (V_1 \otimes V_1) \right) \end{aligned}$$

$$\begin{aligned}
&= (\mathcal{I}_p)^T (e_i^p (e_j^r)^T \otimes C) (-\Lambda \otimes I_n - I_r \otimes A)^{-1} \left(\text{vec} \left(B\tilde{B}^T + \sum_{k=1}^m N_k V_1 \tilde{N}_k^T \right) \right. \\
&\quad \left. + \text{vec} \left(H(V_1 \otimes V_1) \tilde{H}^T \right) \right) \\
&= (\mathcal{I}_p)^T (e_i^p (e_j^r)^T \otimes C) (\text{vec}(V_1) + \text{vec}(V_2)) = \text{tr} \left(C(V_1 + V_2) e_j^r (e_i^p)^T \right) \\
&= \text{tr} \left(C V e_j^r (e_i^p)^T \right),
\end{aligned}$$

where V_2 solves (3.29c) and $V = V_1 + V_2$. Using the similar steps, we can show that the right-side of (B.5) is equal to $\text{tr} \left(\hat{C} \hat{V} e_j^r (e_i^p)^T \right)$, where \hat{V} is defined in (3.31). Therefore, (B.5) is the same as (3.32a).

Necessary conditions with respect to Λ . By utilizing Lemma A.1, we aim at deriving the necessary condition with respect to the i th diagonal entry of Λ . We differentiate \mathcal{E} w.r.t. λ_i to obtain

$$\begin{aligned}
\frac{\partial \mathcal{E}^2}{\partial \lambda_i} &= 2(\mathcal{I}_p)^T \left(\tilde{C}^e \otimes \tilde{C}^e \right) \mathcal{L}_e^{-1} \mathbb{E} \mathcal{L}_e^{-1} \left(\left(\tilde{B}^e \otimes \tilde{B}^e \right) \mathcal{I}_m + \sum_{k=1}^m \left(\tilde{N}_k^e \otimes \tilde{N}_k^e \right) \mathcal{P}_l \right. \\
&\quad \left. + \left(\tilde{H}^e \otimes \tilde{H}^e \right) T_{(n+r, n+r)}(\mathcal{P}_l \otimes \mathcal{P}_l) \right) + (\mathcal{I}_p)^T \left(\tilde{C}^e \otimes \tilde{C}^e \right) \mathcal{L}_e^{-1} \\
&\quad \times \left(2 \sum_{k=1}^m \left(\tilde{N}_k^e \otimes \tilde{N}_k^e \right) \mathcal{L}_e^{-1} \mathbb{E} \mathcal{P}_l + 4 \left(\tilde{H}^e \otimes \tilde{H}^e \right) T_{(n+r, n+r)} \left((\mathcal{L}_e^{-1} \mathbb{E} \mathcal{P}_l) \otimes \mathcal{P}_l \right) \right),
\end{aligned}$$

where

$$\mathcal{L}_e = - \left(\tilde{A}^e \otimes I_{n+r} + I_{n+r} \otimes \tilde{A}^e \right) \quad \text{and} \quad \mathbb{E} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & e_i^r (e_i^r)^T \end{bmatrix} \otimes I_{n+r}.$$

Performing some algebraic calculations gives rise to the following expression:

$$\begin{aligned}
\frac{\partial \mathcal{E}^2}{\partial \lambda_i} &= 2(\mathcal{I}_p)^T \left(-\tilde{C} \otimes \tilde{C}^e \right) \mathcal{Z}_e^{-1} \Xi_{n+r} \mathcal{Z}_e^{-1} \left(\left(\tilde{B} \otimes \tilde{B}^e \right) \mathcal{I}_m + \sum_{k=1}^m \left(\tilde{N}_k \otimes \tilde{N}_k^e \right) \mathcal{P}_1^{(2)} \right. \\
&\quad \left. + \left(\tilde{H} \hat{\mathcal{F}} \otimes \tilde{H}^e \right) T_{(n+r, n+r)}(\mathcal{P}_l \otimes \mathcal{P}_l) \right) + 2(\mathcal{I}_p)^T \left(-\tilde{C} \otimes \tilde{C}^e \right) \mathcal{Z}_e^{-1} \\
&\quad \times \left(\sum_{k=1}^m \left(\tilde{N}_k \otimes \tilde{N}_k^e \right) \mathcal{Z}_e^{-1} \Xi_{n+r} \mathcal{P}_1^{(2)} + 2 \left(\tilde{H} \hat{\mathcal{F}} \otimes \tilde{H}^e \right) T_{(n+r, n+r)} (\mathcal{L}_e^{-1} \mathbb{E} \mathcal{P}_l \otimes \mathcal{P}_l) \right),
\end{aligned}$$

where $\mathcal{Z}_e := -(\Lambda \otimes I_{n+r} + I_r \otimes A^e)$ and $\Xi_m := (e_i^r (e_i^r)^T \otimes I_m)$. Next, we utilize Lemma A.2 and use the permutation matrix M (as done while deriving the necessary conditions with respect to \tilde{C}) to obtain

$$\begin{aligned}
\frac{\partial \mathcal{E}^2}{\partial \lambda_i} &= 2(\mathcal{I}_p)^T \mathcal{S} \left(\left(\tilde{B} \otimes B \right) \mathcal{I}_m + \sum_{k=1}^m \left(\tilde{N}_k \otimes N_k \right) p_3 + \left(\tilde{H} \otimes H \right) T_{(n,r)}(p_3 \otimes p_3) \right) \\
&\quad - 2(\mathcal{I}_p)^T \tilde{\mathcal{S}} \left(\left(\tilde{B} \otimes \tilde{B} \right) \mathcal{I}_m + \sum_{k=1}^m \left(\tilde{N}_k \otimes \tilde{N}_k \right) p_4 + \left(\tilde{H} \otimes \tilde{H} \right) T_{(r,r)}(p_4 \otimes p_4) \right) \\
&\quad + 2(\mathcal{I}_p)^T \left(\tilde{C} \otimes C \right) L^{-1} \left(\sum_{k=1}^m \left(\tilde{N}_k \otimes N_k \right) L^{-1} \Xi_n p_3 + 2 \left(\tilde{H} \otimes H \right) T_{(n,r)}(L^{-1} \Xi_n p_3 \otimes p_3) \right) \\
&\quad - 2(\mathcal{I}_p)^T \left(\tilde{C} \otimes \tilde{C} \right) \tilde{L}^{-1} \left(\sum_{k=1}^m \left(\tilde{N}_k \otimes \tilde{N}_k \right) \tilde{L}^{-1} \Xi_r p_4 + 2 \left(\tilde{H} \otimes \tilde{H} \right) T_{(r,r)}(\tilde{L}^{-1} \Xi_r p_4 \otimes p_4) \right)
\end{aligned}$$

$$\begin{aligned}
&= 2(\mathcal{I}_p)^T \mathcal{S} \left((\tilde{B} \otimes B) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes N_k) p_3 + (\tilde{H} \otimes H) T_{(n,r)}(p_3 \otimes p_3) \right) \\
&\quad - 2(\mathcal{I}_p)^T \tilde{\mathcal{S}} \left((\tilde{B} \otimes \hat{B}) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes \hat{N}_k) \hat{p}_4 + (\tilde{H} \otimes \hat{H}) T_{(r,r)}(\hat{p}_4 \otimes \hat{p}_4) \right) \\
&\quad + 2(\mathcal{I}_p)^T \left(\tilde{C} \otimes C \right) L^{-1} \left(\sum_{k=1}^m (\tilde{N}_k \otimes N_k) L^{-1} \Xi_n p_3 + 2(\tilde{H} \otimes H) T_{(n,r)}(L^{-1} \Xi_n p_3 \otimes p_3) \right) \\
&\quad - 2(\mathcal{I}_p)^T \left(\tilde{C} \otimes \hat{C} \right) \hat{L}^{-1} \left(\sum_{k=1}^m (\tilde{N}_k \otimes \hat{N}_k) \hat{L}^{-1} \Xi_r \hat{p}_4 + 2(\tilde{H} \otimes \hat{H}) T_{(r,r)}(\hat{L}^{-1} \Xi_r (\hat{p}_4 \otimes \hat{p}_4)) \right),
\end{aligned}$$

where p_3 and p_4 are the same as defined in (B.3), and

$$\begin{aligned}
\mathcal{S} &:= \left(\tilde{C} \otimes C \right) L^{-1} (e_i^r (e_i^r)^T \otimes I_n) L^{-1}, & \tilde{\mathcal{S}} &:= \left(\tilde{C} \otimes \hat{C} \right) \tilde{L}^{-1} (e_i^r (e_i^r)^T \otimes I_r) \tilde{L}^{-1}, \\
\hat{\mathcal{S}} &:= \left(\tilde{C} \otimes \hat{C} \right) \hat{L}^{-1} (e_i^r (e_i^r)^T \otimes I_r) \hat{L}^{-1}, & L &:= -(\Lambda \otimes I_n + I_r \otimes A), \\
\tilde{L} &:= -(\Lambda \otimes I_r + I_r \otimes \Lambda), & \hat{L} &:= -(\Lambda \otimes I_r + I_r \otimes \hat{A}).
\end{aligned}$$

By using properties derived in Lemma 2.2, we can simplify the above equation

$$\begin{aligned}
\frac{\partial \mathcal{E}^2}{\partial \lambda_i} &= 2(\mathcal{I}_p)^T \mathcal{S} \left((\tilde{B} \otimes B) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes N_k) p_3 + (\tilde{H} \otimes H) T_{(n,r)}(p_3 \otimes p_3) \right) \\
&\quad - 2(\mathcal{I}_p)^T \tilde{\mathcal{S}} \left((\tilde{B} \otimes \hat{B}) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes \hat{N}_k) p_4 + (\tilde{H} \otimes \hat{H}) T_{(r,r)}(\hat{p}_4 \otimes \hat{p}_4) \right) \\
&\quad + 2(\mathcal{I}_m)^T \left(\tilde{B} \otimes B \right) L^{-T} \Xi_n L^{-T} \left(\sum_{k=1}^m (\tilde{N}_k \otimes N_k)^T q_3 + 2(\tilde{\mathcal{H}}^{(2)} \otimes \mathcal{H}^{(2)}) T_{(n,r)}(p_3 \otimes q_3) \right) \\
&\quad - 2(\mathcal{I}_m)^T \left(\tilde{B} \otimes \hat{B} \right) \hat{L}^{-T} \Xi_r \hat{L}^{-T} \left(\sum_{k=1}^m (\tilde{N}_k \otimes \hat{N}_k)^T \hat{q}_4 + 2(\tilde{\mathcal{H}}^{(2)} \otimes \hat{\mathcal{H}}^{(2)}) T_{(r,r)}(\hat{p}_4 \otimes \hat{q}_4) \right),
\end{aligned}$$

where

$$q_3 = (-\Lambda \otimes I_n - I_r \otimes A)^{-T} \left(\tilde{C} \otimes C \right) \mathcal{I}_p \quad \text{and} \quad \hat{q}_4 = (-\Lambda \otimes I_r - I_r \otimes A_r)^{-T} \left(\tilde{C} \otimes \hat{C} \right) \mathcal{I}_p.$$

Once again, we determine an interpolation-based necessary condition with respect to λ_i by setting the last equation equal to zero:

(B.6)

$$\begin{aligned}
&(\mathcal{I}_p)^T \mathcal{S} \left((\tilde{B} \otimes B) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes N_k) p_3 + (\tilde{H} \otimes H) T_{(n,r)}(p_3 \otimes p_3) \right) \\
&\quad + (\mathcal{I}_m)^T \left(\tilde{B} \otimes B \right) L^{-T} \Xi_n L^{-T} \left(\sum_{k=1}^m (\tilde{N}_k \otimes N_k)^T q_3 + 2(\tilde{\mathcal{H}}^{(2)} \otimes \mathcal{H}^{(2)}) T_{(n,r)}(p_3 \otimes q_3) \right) \\
&= (\mathcal{I}_p)^T \tilde{\mathcal{S}} \left((\tilde{B} \otimes \hat{B}) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes \hat{N}_k) p_4 + (\tilde{H} \otimes \hat{H}) T_{(r,r)}(\hat{p}_4 \otimes \hat{p}_4) \right), \\
&\quad + (\mathcal{I}_m)^T \left(\tilde{B} \otimes \hat{B} \right) \hat{L}^{-T} \Xi_r \hat{L}^{-T} \left(\sum_{k=1}^m (\tilde{N}_k \otimes \hat{N}_k)^T \hat{q}_4 + 2(\tilde{\mathcal{H}}^{(2)} \otimes \hat{\mathcal{H}}^{(2)}) T_{(r,r)}(\hat{p}_4 \otimes \hat{q}_4) \right).
\end{aligned}$$

Now, we first simplify the left-side of the above equation using Lemma 2.4 and (2.1). We first focus of the first part of the left-side of (B.6). This yields

$$\begin{aligned}
&(\mathcal{I}_p)^T \mathcal{S} \left((\tilde{B} \otimes B) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes N_k) p_3 + (\tilde{H} \otimes H) T_{(n,r)}(p_3 \otimes p_3) \right) \\
&= (\mathcal{I}_p)^T \left(\tilde{C} \otimes C \right) L^{-1} (e_i^r (e_i^r)^T \otimes I_n) \\
&\quad \times L^{-1} \left((\tilde{B} \otimes B) \mathcal{I}_m + \sum_{k=1}^m (\tilde{N}_k \otimes N_k) p_3 + (\tilde{H} \otimes H) T_{(n,r)}(p_3 \otimes p_3) \right)
\end{aligned}$$

$$\begin{aligned}
&= \underbrace{(\mathcal{I}_p)^T \left(\tilde{C} \otimes C \right)}_{(\text{vec}(W_1))^T} L^{-1} (e_i^r (e_i^r)^T \otimes I_n) \text{vec}(V) = \text{tr} \left(V e_i^r (e_i^r)^T W_1^T \right) \\
&= (V_1(:, i))^T W(:, i), = (W_1(:, i))^T V(:, i),
\end{aligned}$$

where W_1 solves (3.29b). Analogously, we can show that (B.7)

$$\begin{aligned}
&(\mathcal{I}_m)^T \left(\tilde{B} \otimes B \right) L^{-T} \Xi_n L^{-T} \left(\sum_{k=1}^m (\tilde{N}_k \otimes N_k)^T q_3 + 2(\tilde{\mathcal{H}}^{(2)} \otimes \mathcal{H}^{(2)}) T_{(n,r)}(p_3 \otimes q_3) \right) \\
&= (W(:, i))^T V_1(:, i).
\end{aligned}$$

Thus, the left-side of (B.6) is equal to $(W(:, i))^T V_1(:, i) + (W_1(:, i))^T V(:, i)$. Using the similar steps, we can also show that the right-side of (B.6) is equivalent to $(\widehat{W}(:, i))^T \widehat{V}_1(:, i) + (\widehat{W}_1(:, i))^T \widehat{V}(:, i)$. Thus, we obtain the optimality conditions with respect to Λ as (3.32e).

The necessary conditions with respect to \tilde{B} , \tilde{N} and \tilde{H} can also be determined in a similar manner as for \tilde{C} and λ_i . For brevity of the paper, we skip detailed derivations; however, we state final optimality conditions. A necessary condition for optimality with respect to the (i, j) entry of \tilde{N}_k is

$$(\mathcal{I}_p)^T \left(\tilde{C} \otimes C \right) L^{-1} \left((e_i^r (e_j^r)^T \otimes N_k) p_3 \right) = (\mathcal{I}_p)^T \left(\tilde{C} \otimes \hat{C} \right) \hat{L}^{-1} \left((e_i^r (e_j^r)^T \otimes \hat{N}_k) \hat{p}_4 \right),$$

which then yields (3.32c) in the Sylvester equation form. A similar optimality condition with respect to the (i, j) entry of \tilde{H} is given by

$$\begin{aligned}
&(\mathcal{I}_p)^T \left(\tilde{C} \otimes C \right) L^{-1} \left((e_i^r (e_j^r)^T \otimes H) T_{(n,r)}(p_3 \otimes p_3) \right) \\
&= (\mathcal{I}_p)^T \left(\tilde{C} \otimes \hat{C} \right) \hat{L}^{-1} \left((e_i^r (e_j^r)^T \otimes \hat{H}) T_{(r,r)}(\hat{p}_4 \otimes \hat{p}_4) \right),
\end{aligned}$$

which can be equivalently described as (3.32d). Finally, the necessary condition appearing with respect to the (i, j) entry of \tilde{B} is

$$\begin{aligned}
&(\mathcal{I}_m)^T \left(e_i^r (e_j^m)^T \otimes B \right) L^{-T} \left((\tilde{C} \otimes C) \mathcal{I}_p + \sum_{k=1}^m (\tilde{N}_k \otimes N_k)^T q_3 \right. \\
&\quad \left. + 2(\tilde{\mathcal{H}}^{(2)} \otimes \mathcal{H}^{(2)}) T_{(n,r)}(p_3 \otimes q_3) \right), \\
&= (\mathcal{I}_m)^T \left(e_i^r (e_j^m)^T \otimes \hat{B} \right) \hat{L}^{-T} \left((\tilde{C} \otimes \hat{C}) \mathcal{I}_p + \sum_{k=1}^m (\tilde{N}_k \otimes \hat{N}_k)^T \hat{q}_4 \right. \\
&\quad \left. + 2(\tilde{\mathcal{H}}^{(2)} \otimes \hat{\mathcal{H}}^{(2)}) T_{(r,r)}(\hat{p}_4 \otimes \hat{q}_4) \right),
\end{aligned}$$

which gives rise to (3.32b).

Appendix C. Proof of Theorem 3.9.

We begin by establishing a relationship between $V_1 \in \mathbb{R}^{n \times r}$, $\widehat{V}_1 \in \mathbb{R}^{r \times r}$ and $V \in \mathbb{R}^{n \times r}$. For this, consider the Sylvester equation related to V_1

$$(C.1) \quad -V_1 \Lambda - A V_1 = B \tilde{B}^T,$$

and an oblique projector $\Pi_v := V_1(W^T V_1)^{-1} W^T$. Then, we apply the projector Π_v to the Sylvester equation (C.1) from the left to obtain

$$(C.2) \quad \begin{aligned} & -V_1 \Lambda - \Pi_v A V_1 = \Pi_v B \tilde{B}^T, \text{ and} \\ & -V_1 \Lambda - \Pi A V_1 = (\Pi_v - \Pi) A V_1 + \Pi_v B \tilde{B}^T, \end{aligned}$$

where $\Pi := V(W^T V)^{-1} W^T$. Now, we recall that \widehat{V}_1 satisfies the Sylvester equation

$$-\widehat{V}_1 \Lambda - \widehat{A} \widehat{V}_1 = \widehat{B} \tilde{B}^T.$$

We next multiply it by V from the left and substitute for \widehat{A} and \widehat{B} to obtain

$$(C.3) \quad -V \widehat{V}_1 \Lambda - \Pi A V \widehat{V}_1 = \Pi B \tilde{B}^T.$$

Subtracting (C.2) from (C.3) yields

$$(V_1 - V \widehat{V}_1) \Lambda + \Pi A (V_1 - V \widehat{V}_1) = (\Pi - \Pi_v) (A V_1 + B \tilde{B}^T).$$

Since it is assumed that $\sigma(\widehat{A}) \cap \sigma(-\Pi A) = \emptyset$, this implies that $\Lambda \otimes I_n + I_r \otimes (\Pi A)$ is invertible. Therefore, we can write

$$(C.4) \quad V_1 = V \widehat{V}_1 + \epsilon_v,$$

where ϵ_v solves the Sylvester equation

$$(C.5) \quad \epsilon_v \Lambda + \Pi_v A \epsilon_v = (\Pi - \Pi_v) (A V_1 + B \tilde{B}^T).$$

Similarly, one can show that

$$(C.6) \quad W_1 = W(W^T V)^{-T} \widehat{W}_1 + \epsilon_w,$$

where ϵ_w solves

$$\epsilon_w \Lambda + \Pi^T A^T \epsilon_w = (\Pi^T - \Pi_w) (A^T W_1 + C^T \tilde{C}),$$

in which $\Pi_w := W_1(V^T W)V^T$. Using (C.4) and (C.6), we obtain

$$\begin{aligned} \widehat{W}_1(:, i)^T \widehat{N}_k \widehat{V}_1(:, j) &= \widehat{W}_1(:, i)^T (W^T V)^{-1} W^T N_k V \widehat{V}_1(:, j) \\ &= (W_1(:, i) - \epsilon_w(:, i))^T N_k (V_1(:, j) - \epsilon_v(:, j)) \\ &= W_1(:, i)^T N_k V_1(:, j) - (\epsilon_w(:, i))^T N_k (V_1(:, j) - \epsilon_v(:, j)) \\ &\quad - (W_1(:, i))^T N_k (\epsilon_v(:, j)), \end{aligned}$$

which is (3.33c) in Theorem 3.7. Similarly, one can prove (3.33d). To prove (3.33a), we consider the following Sylvester equation for V :

$$(C.7) \quad V(-\Lambda) - AV = B \tilde{B}^T + \sum_{k=1}^m N_k V_1 \tilde{N}_k^T + H(V_1 \otimes V_1) \tilde{H}^T.$$

Applying Π to both sides of the above Sylvester equation yields

$$(C.8) \quad V(I_r(-\Lambda) - \widehat{A} I_r) = V(\widehat{B} \tilde{B}^T + \mathcal{Y}),$$

where $\mathcal{Y} = (W^T V)^{-1} W^T \left(\sum_{k=1}^m N_k V_1 \tilde{N}_k^T + H(V_1 \otimes V_1) \tilde{H}^T \right)$. This implies that

$$(C.9) \quad I_r(-\Lambda) - \hat{A}I_r = \hat{B}\tilde{B}^T + \mathcal{Y}.$$

Next, we consider the Sylvester equation for \hat{V} ,

$$(C.10) \quad \hat{V}(-\Lambda) - \hat{A}\hat{V} = \hat{B}\tilde{B}^T + \sum_{k=1}^m \hat{N}_k \hat{V}_1 \tilde{N}_k^T + \hat{H}(\hat{V}_1 \otimes \hat{V}_1) \tilde{H}^T.$$

We then subtract (C.10) and (C.9) to obtain

$$\begin{aligned} (I_r - \hat{V})(-\Lambda) - \hat{A}(I_r - \hat{V}) &= \sum_{k=1}^m (W^T V)^{-1} W^T N_k (V_1 - V\hat{V}_1) \tilde{N}_k^T \\ &\quad + (W^T V)^{-1} W^T H (V_1 \otimes V_1 - (V\hat{V}_1 \otimes V\hat{V}_1)) \tilde{H}^T. \end{aligned}$$

Substituting $V\hat{V}_1$ from (C.4) gives

$$\begin{aligned} (I_r - \hat{V})(-\Lambda) - \hat{A}(I_r - \hat{V}) &= \sum_{k=1}^m (W^T V)^{-1} W^T N_k \epsilon_v \tilde{N}_k^T \\ &\quad + (W^T V)^{-1} W^T H (\epsilon_v \otimes V_1 + V_1 \otimes \epsilon_v + \epsilon_v \otimes \epsilon_v) \tilde{H}^T. \end{aligned}$$

Since Λ contains the eigenvalues of \hat{A} and \hat{A} is stable, Λ and $-\hat{A}$ cannot have any common eigenvalues. Hence, the matrix $\Lambda \otimes I_r + I_r \otimes \hat{A}$ is invertible. Therefore the above Sylvester equations in $\Gamma := \hat{V} - I_r$ exists and solves

$$\begin{aligned} \Gamma_v \Lambda + \hat{A}\Gamma_v &= \sum_{k=1}^m (W^T V)^{-1} W^T N_k \epsilon_v \tilde{N}_k^T \\ &\quad + (W^T V)^{-1} W^T H (\epsilon_v \otimes V_1 + V_1 \otimes \epsilon_v + \epsilon_v \otimes \epsilon_v) \tilde{H}^T. \end{aligned}$$

To prove (3.33a), we observe that

$$\begin{aligned} \text{tr} \left(\hat{C}\hat{V}e_i^r (e_j^p)^T \right) &= \text{tr} \left(CV(I_r + \Gamma_v)e_i^r (e_j^p)^T \right) \\ &= \text{tr} \left(CVe_i^r (e_j^p)^T \right) + \text{tr} \left(CV\Gamma_v e_i^r (e_j^p)^T \right). \end{aligned}$$

Thus,

$$\text{tr} \left(CVe_i^r (e_j^p)^T \right) = \text{tr} \left(\hat{C}\hat{V}e_i^r (e_j^p)^T \right) + \epsilon_C^{(i,j)}.$$

Analogously, we can prove that there exists Γ_w such that $\hat{W} = (W^T V)^T + \Gamma_w$ and it satisfies

$$\Gamma_w \Lambda + \hat{A}^T \Gamma_w = V^T \left(\sum_{k=1}^m N_k^T \epsilon_w \tilde{N}_k + \mathcal{H}^{(2)}(\epsilon_v \otimes (W_1 + \epsilon_w) + V_1 \otimes \epsilon_w) \left(\mathcal{H}^{(2)} \right)^T \right).$$

To prove (3.33b), we observe that

$$\text{tr} \left(\hat{B}^T \hat{W} e_i^r (e_j^m)^T \right) = \text{tr} \left(B^T W (W^T V)^{-T} ((W^T V)^T + \Gamma_w) e_i^r (e_j^m)^T \right).$$

Thus,

$$\text{tr} \left(\widehat{B}^T \widehat{W} e_i^r (e_j^m)^T \right) = \text{tr} \left(B^T W^T + B^T W (W^T V)^{-T} \Gamma_w \right) e_i^r (e_j^p)^T.$$

Since we now know that $\widehat{V} = I_r + \Gamma_v$ and $\widehat{W} = (W^T V)^T + \Gamma_w$, hence, we get

$$(C.11) \quad V \widehat{V} = V + V \Gamma_v \quad \text{and} \quad W (W^T V)^{-T} \widehat{W} = W + W (W^T V)^{-T} \Gamma_w.$$

We make use of (C.11) to prove (3.33e) in the following:

$$\begin{aligned} & (W_1(:, i))^T V(:, i) + (W_2(:, i))^T V_1(:, i) \\ &= (W(:, i))^T V(:, i) - (W_2(:, i))^T V_2(:, i) \\ &= \left(W (W^T V)^{-T} \left(\widehat{W}(:, i) - \Gamma_w(:, i) \right) \right)^T V \left(\widehat{V}(:, i) - \Gamma_v(:, i) \right) \\ &\quad - (W_2(:, i))^T V_2(:, i) \\ &= \left(\widehat{W}(:, i) - \Gamma_w(:, i) \right)^T \left(\widehat{V}(:, i) - \Gamma_v(:, i) \right) - (W_2(:, i))^T V_2(:, i) \\ &= \left(\widehat{W}(:, i) \right)^T \widehat{V}(:, i) - \left(\widehat{W}(:, i) \right)^T \Gamma_v(:, i) - \left(\Gamma_w(:, i) \right)^T \left(\widehat{V}(:, i) - \Gamma_v(:, i) \right) \\ &\quad - (W_2(:, i))^T V_2(:, i) \\ &= \left(\widehat{W}_1(:, i) \right)^T \widehat{V}(:, i) + \left(\widehat{W}_2(:, i) \right)^T \widehat{V}_1(:, i) + \epsilon_\lambda^{(i)}, \end{aligned}$$

where

$$\begin{aligned} \epsilon_\lambda^{(i)} &= - \left(\widehat{W}(:, i) \right)^T \Gamma_v(:, i) - \left(\Gamma_w(:, i) \right)^T \left(\widehat{V}(:, i) - \Gamma_v(:, i) \right) \\ &\quad - (W_2(:, i))^T V_2(:, i) + \left(\widehat{W}_2(:, i) \right)^T \widehat{V}_2(:, i). \end{aligned}$$

This completes the proof.