

Hessian-based model reduction for large-scale systems with initial-condition inputs

O. Bashir¹, K. Willcox^{1,*},[†], O. Ghattas², B. van Bloemen Waanders³ and J. Hill³

¹*Massachusetts Institute of Technology, Cambridge, MA 02139, U.S.A.*

²*University of Texas at Austin, Austin, TX 78712, U.S.A.*

³*Sandia National Laboratories[‡], Albuquerque, NM 87185, U.S.A.*

SUMMARY

Reduced-order models that are able to approximate output quantities of interest of high-fidelity computational models over a wide range of input parameters play an important role in making tractable large-scale optimal design, optimal control, and inverse problem applications. We consider the problem of determining a reduced model of an initial value problem that spans all important initial conditions, and pose the task of determining appropriate training sets for reduced-basis construction as a sequence of optimization problems. We show that, under certain assumptions, these optimization problems have an explicit solution in the form of an eigenvalue problem, yielding an efficient model reduction algorithm that scales well to systems with states of high dimension. Furthermore, tight upper bounds are given for the error in the outputs of the reduced models. The reduction methodology is demonstrated for a large-scale contaminant transport problem. Copyright © 2007 John Wiley & Sons, Ltd.

Received 16 January 2007; Revised 29 March 2007; Accepted 19 April 2007

KEY WORDS: model reduction; optimization; initial-condition problems

1. INTRODUCTION

Reduced-order models that are able to approximate outputs of high-fidelity computational models over a wide range of input parameters have an important role to play in making tractable large-scale optimal design, optimal control, and inverse problem applications. In particular, the state

*Correspondence to: K. Willcox, 77 Massachusetts Avenue, Room 37-447, Cambridge, MA 02139, U.S.A.

[†]E-mail: kwillcox@mit.edu

[‡]Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy's National Nuclear Security Administration under Contract DE-AC04-94-AL85000.

Contract/grant sponsor: National Science Foundation; contract grant/numbers: CNS-0540372 and CNS-0540186

Contract/grant sponsor: Air Force Office of Scientific Research

Contract/grant sponsor: Computer Science Research Institute at Sandia National Laboratories

estimation inverse problem setting requires a reduced model that spans the space of important initial conditions, i.e. those that have the greatest influence on the output quantities of interest. Creating such a model with existing model reduction techniques presents a significant challenge due to the need to sample adequately the high-dimensional space of possible initial conditions. In this paper, we present a new methodology that employs an efficient sampling strategy to make tractable the task of determining a reduced model for large-scale linear initial value problems that are accurate over all initial conditions.

For the most part, reduction techniques for large-scale systems have focused on a projection framework that utilizes a reduced-space basis. Methods to compute the basis in the large-scale setting include Krylov-subspace methods [1–3], approximate balanced truncation [4–7], and proper orthogonal decomposition (POD) [8–10]. Progress has been made in development and application of these methods to optimization applications with a small number of input parameters, for example, optimal control [11–14] and parametrized design of interconnect circuits [15]. In the case of a high-dimensional input parameter space, the computational cost of determining the reduced basis by these techniques becomes prohibitive unless some sparse sampling strategy is employed.

For initial-condition problems of moderate dimension, a reduction method has been proposed that truncates a balanced representation of the finite-dimensional Hankel operator [16]. In [17], POD was used in a large-scale inverse problem setting to define a reduced space for the initial condition in which to solve the data assimilation problem. In that work, only a single initial condition was used to generate the state solutions necessary to form the reduced basis: either the true initial condition which does contain the necessary information but would be unavailable in practice, or the background estimate of the initial state which defines a forecast trajectory that may not be sufficiently rich in terms of state information.

For model reduction of linear time-invariant systems using multipoint rational Krylov approximations, two methods have been recently proposed to choose sample locations: an iterative method to choose an optimal set of interpolation points [18], and a heuristic statistically based resampling scheme to select sample points [19]. To address the more general challenge of sampling a high-dimensional parameter space to build a reduced basis, the greedy algorithm was introduced in [20]. The key premise of the greedy algorithm is to adaptively choose samples by finding the location in parameter space where the error in the reduced model is maximal. In [21], the greedy algorithm was applied to find reduced models for the parametrized steady incompressible Navier–Stokes equations. In [22, 23], the greedy algorithm was combined with *a posteriori* error estimators for parametrized parabolic partial differential equations (PDEs), and applied to several optimal control and inverse problems.

Here, we address the problem of determining a reduced basis, and hence a reduced model, for large-scale linear initial value problems that is accurate over all possible initial-conditions. The reduced basis is associated with a judicious sampling of the initial-condition space. The basis spans these initial-condition samples, as well as the state trajectories determined by them. The span can be computed by the POD, or else by solution of an optimization problem to find the basis that minimizes the output error at the sample points [24]. The sampling problem itself is formulated as a greedy optimization problem. Rather than invoke error estimators to approximate the errors in the outputs as in [20–23], the objective function of the greedy optimization problem targets the actual errors. To define the errors, the optimization problem must then be constrained by the initial value systems representing the full and reduced models. Under certain reasonable assumptions, this optimization problem admits an explicit solution in the form of an eigenvalue problem for the dominant eigenvectors, which define the samples in initial-condition space and hence the reduced

basis. Furthermore, the eigenvalue form leads to tight, computable upper bounds for the error in the outputs of the reduced model.

This article is organized as follows. Section 2 describes the projection framework used to derive the reduced-order dynamical system. We then present in Section 3 the theoretical approach leading to a basis-construction algorithm. In Section 4, we demonstrate the efficacy of the algorithm *via* numerical experiments on a problem of 2-D convective–diffusive transport. We present an application to model reduction for 3-D contaminant transport in an urban canyon in Section 5, and offer conclusions in Section 6.

2. REDUCED-ORDER DYNAMICAL SYSTEMS

Consider the general linear discrete-time system

$$x(k+1) = Ax(k) + Bu(k), \quad k=0, 1, \dots, T-1 \quad (1)$$

$$y(k) = Cx(k), \quad k=0, 1, \dots, T \quad (2)$$

with initial condition

$$x(0) = x_0 \quad (3)$$

where $x(k) \in \mathbb{R}^N$ is the system state at time t_k , the vector x_0 contains the specified initial state, and we consider a time horizon from $t=0$ to $t=t_T$. The vectors $u(k) \in \mathbb{R}^P$ and $y(k) \in \mathbb{R}^Q$ contain, respectively, the P system inputs and Q system outputs at time t_k . In general, we are interested in systems of the form (1)–(3) that result from spatial and temporal discretization of PDEs. In this case, the dimension of the system, N , is very large and the matrices $A \in \mathbb{R}^{N \times N}$, $B \in \mathbb{R}^{N \times P}$, and $C \in \mathbb{R}^{Q \times N}$ result from the chosen spatial and temporal discretization methods.

A reduced-order model of (1)–(3) can be derived by assuming that the state $x(k)$ is represented as a linear combination of n basis vectors,

$$\hat{x}(k) = Vx_r(k) \quad (4)$$

where $\hat{x}(k) \in \mathbb{R}^N$ is the reduced-model approximation of the state $x(k)$ and $n \ll N$. The projection matrix $V \in \mathbb{R}^{N \times n}$ contains as columns the orthonormal basis vectors V_i , i.e. $V = [V_1 \ V_2 \ \dots \ V_n]$, and the reduced-order state $x_r(k) \in \mathbb{R}^n$ contains the corresponding modal amplitudes for time t_k . Using representation (4) together with a Galerkin projection of the discrete-time system (1)–(3) onto the space spanned by the basis V yields the reduced-order model with state x_r and output y_r ,

$$x_r(k+1) = A_r x_r(k) + B_r u(k), \quad k=0, 1, \dots, T-1 \quad (5)$$

$$y_r(k) = C_r x_r(k), \quad k=0, 1, \dots, T \quad (6)$$

$$x_r(0) = V^T x_0 \quad (7)$$

where $A_r = V^T A V$, $B_r = V^T B$, and $C_r = C V$.

Since system (1)–(3) is linear, the effects of inputs u and initial conditions x_0 can be considered separately. In this paper, we focus on the initial-condition problem and, without loss of generality,

assume that $u(k) = 0, k = 0, 1, \dots, T - 1$. For convenience of notation, we write the discrete-time system (1)–(3) in matrix form as

$$\mathbf{Ax} = \mathbf{F}x_0 \tag{8}$$

$$\mathbf{y} = \mathbf{Cx} \tag{9}$$

where

$$\mathbf{x} = \begin{bmatrix} x(0) \\ x(1) \\ \vdots \\ x(T) \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y(0) \\ y(1) \\ \vdots \\ y(T) \end{bmatrix} \tag{10}$$

The matrices \mathbf{A} , \mathbf{F} , and \mathbf{C} in (8) and (9) are given by

$$\mathbf{A} = \begin{bmatrix} I & 0 & 0 & \dots & 0 \\ -A & I & 0 & \ddots & \vdots \\ 0 & -A & I & \ddots & \ddots \\ & \ddots & \ddots & \ddots & 0 \\ \vdots & & \ddots & \ddots & 0 \\ 0 & & 0 & -A & I \end{bmatrix}, \quad \mathbf{F} = \begin{bmatrix} I \\ 0 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{bmatrix}, \quad \mathbf{C} = \begin{bmatrix} C & 0 & \dots & \dots & \dots & 0 \\ 0 & C & 0 & & & \vdots \\ \vdots & 0 & C & 0 & & \\ & \ddots & \ddots & \ddots & & \\ \vdots & \ddots & \ddots & \ddots & 0 & \\ 0 & & 0 & 0 & C & \end{bmatrix} \tag{11}$$

Similarly, the reduced-order model (5)–(7) can be written in matrix form as

$$\mathbf{A}_r \mathbf{x}_r = \mathbf{F}_r x_0 \tag{12}$$

$$\mathbf{y}_r = \mathbf{C}_r \mathbf{x}_r \tag{13}$$

where $\mathbf{x}_r, \mathbf{y}_r, \mathbf{A}_r,$ and \mathbf{C}_r are defined analogously to $\mathbf{x}, \mathbf{y}, \mathbf{A},$ and \mathbf{C} but with the appropriate reduced-order quantities. The matrix \mathbf{F}_r is given by

$$\mathbf{F}_r = \begin{bmatrix} V^T \\ 0 \\ \vdots \\ 0 \end{bmatrix} \tag{14}$$

In many cases, we are interested in rapid identification of initial conditions from sparse measurements of the states over a time horizon; we thus require a reduced-order model that will provide accurate outputs for any initial condition contained in some set \mathcal{X}_0 . Using the projection

framework described above, the task therefore becomes one of choosing an appropriate basis V so that the error between full-order output \mathbf{y} and the reduced-order output \mathbf{y}_r is small for all initial conditions of interest.

3. HESSIAN-BASED MODEL REDUCTION

In this section, a methodology to determine a basis that spans the space of important initial conditions is presented. To compute the basis *via* a method such as POD, a sample set of initial conditions must be chosen. At each selected initial condition, a forward simulation is performed to generate a set of states, commonly referred to as snapshots, from which the reduced basis is formed. It has been shown that in the case of systems that are linear in the state, POD is equivalent to balanced truncation if the snapshots are computed for all possible initial conditions [25]. Since sampling all possible initial conditions is not feasible for large-scale problems, we propose an adaptive approach to identify important initial conditions that should be sampled. The approach is motivated by the greedy algorithm of [20], which proposed an adaptive approach to determine the parameter locations at which samples are drawn to form a reduced basis. For the linear finite-time-horizon problem considered here, we show that the greedy algorithm can be formulated as an optimization problem that has an explicit solution in the form of an eigenvalue problem.

3.1. Theoretical approach

Our task is to find an appropriate reduced basis and associated reduced model: one that provides accurate outputs for all initial conditions of interest. We define an optimal basis, V^* , to be one that minimizes the maximal L_2 error between the full-order and reduced-order outputs of the fully discrete system over all admissible initial conditions,

$$V^* = \arg \min_V \max_{x_0 \in \mathcal{X}_0} (\mathbf{y} - \mathbf{y}_r)^T (\mathbf{y} - \mathbf{y}_r) \quad (15)$$

where

$$\mathbf{A}\mathbf{x} = \mathbf{F}\mathbf{x}_0 \quad (16)$$

$$\mathbf{y} = \mathbf{C}\mathbf{x} \quad (17)$$

$$\mathbf{A}_r\mathbf{x}_r = \mathbf{F}_r\mathbf{x}_0 \quad (18)$$

$$\mathbf{y}_r = \mathbf{C}_r\mathbf{x}_r \quad (19)$$

For this formulation, the only restriction that we place on the set \mathcal{X}_0 is that it contain vectors of unit length. This prevents unboundedness in the optimization problem, since otherwise, the error in the reduced system could be made arbitrarily large. Naturally, because the system is linear, the basis V^* will still be valid for initial conditions of any finite norm.

A suboptimal but computationally efficient approach to solving the optimization problem (15)–(19) is inspired by the greedy algorithm of [20]. Construction of a reduced basis for a steady or unsteady problem with parameter dependence, as considered in [21, 22], requires a set of snapshots, or state solutions, over the parameter–time space. The greedy algorithm adaptively selects these snapshots by finding the location in parameter–time space where the error between

the full-order and reduced-order models is maximal, updating the basis with information gathered from this sample location, forming a new reduced model, and repeating the process. In the case of the initial-condition problem (15)–(19), the greedy approach amounts to sampling at the initial condition $x_0^* \in \mathcal{X}_0$ that *maximizes* the error in (15).

The key step in the greedy algorithm is finding the worst-case initial condition x_0^* , which we achieve by solving the modified optimization problem

$$x_0^* = \arg \max_{x_0 \in \mathcal{X}_0} (\mathbf{y} - \mathbf{y}_r)^T (\mathbf{y} - \mathbf{y}_r) \tag{20}$$

where

$$\mathbf{A}\mathbf{x} = \mathbf{F}x_0 \tag{21}$$

$$\mathbf{y} = \mathbf{C}\mathbf{x} \tag{22}$$

$$\mathbf{A}_r\mathbf{x}_r = \mathbf{F}_rx_0 \tag{23}$$

$$\mathbf{y}_r = \mathbf{C}_r\mathbf{x}_r \tag{24}$$

Equations (20)–(24) define a large-scale optimization problem which includes the full-scale dynamics (21), (22) as constraints. The approach taken in [21, 22] is to replace these constraints with error estimators, so that the full-scale model does not need to be invoked during solution of the optimization problem. Further, in [21, 22], the optimization problem (20)–(24) is solved by a grid-search technique that addresses problems associated with non-convexity and non-availability of derivatives.

In the present article, we exploit the linearity of the state equations to eliminate the full-order and reduced-order states and yield an equivalent unconstrained optimization problem. Eliminating constraints (21)–(24) by solving for the full and reduced states yields

$$x_0^* = \arg \max_{x_0 \in \mathcal{X}_0} x_0^T H^e x_0 \tag{25}$$

where

$$H^e = (\mathbf{C}\mathbf{A}^{-1}\mathbf{F} - \mathbf{C}_r\mathbf{A}_r^{-1}\mathbf{F}_r)^T (\mathbf{C}\mathbf{A}^{-1}\mathbf{F} - \mathbf{C}_r\mathbf{A}_r^{-1}\mathbf{F}_r) \tag{26}$$

It can be seen that (25) is a quadratic unconstrained optimization problem with Hessian matrix $H^e \in \mathbb{R}^{N \times N}$. From (26), it can be seen that H^e is a symmetric positive semi-definite matrix that does not depend on the state or initial condition. The eigenvalues of H^e are therefore non-negative. Since we are considering initial conditions of unit norm, the solution x_0^* maximizes the Rayleigh quotient; therefore, the solution of (25) is given by the eigenvector corresponding to the largest eigenvalue of H^e . This eigenvector is the initial condition for which the error in reduced-model output prediction is largest.

These ideas motivate the following basis-construction algorithm for the initial condition problem.

Algorithm 1

Greedy Reduced-Basis Construction.

Initialize with $V = 0$, so that the initial reduced-order model is zero.

1. For the error Hessian matrix, H^e as defined in (26), find the eigenvector z_1^e with largest eigenvalue λ_1^e .
2. Set $x_0 = z_1^e$ and compute the corresponding solution \mathbf{x} using (8).
3. Update the basis V by adding the new information from the snapshots $x(k)$, $k = 0, 1, \dots, T$.
4. Update the reduced model using the new basis and return to Step 1.

In Step 3 of Algorithm 1, the basis could be computed from the snapshots, using, for example, the POD. A rigorous termination criterion for the algorithm is available in the form of an error bound, which will be discussed below. It should be noted that, while the specific form of Algorithm 1 applies only in the linear case, the greedy sampling concept is applicable to non-linear problems. In the general non-linear case, one would solve an optimization problem similar in form to (20)–(24), but with the appropriate non-linear governing equations appearing as constraints. In this case, the explicit eigenvalue solution to the optimization problem would not hold; instead, one would use a method that is appropriate for large-scale simulation-constrained optimization (see [26]) to solve the resulting optimization problem.

Under certain assumptions, the form of H^e in (25) can be simplified, leading to an algorithm that avoids construction of the reduced model at every greedy iteration. We proceed by decomposing a general initial condition vector as

$$x_0 = x_0^V + x_0^\perp \quad (27)$$

where x_0^V is the component of x_0 in the subspace spanned by the current basis V , and x_0^\perp is the component of x_0 in the orthogonal complement of that subspace. Substituting (27) into the objective function (25), we recognize that $\mathbf{F}_r x_0^\perp = 0$, using the form of \mathbf{F}_r given by (14) and that, by definition, $V^T x_0^\perp = 0$. The unconstrained optimization problem (25) can therefore be written as

$$x_0^* = \arg \max_{x_0 \in \mathcal{X}_0} (\mathbf{C}\mathbf{A}^{-1}\mathbf{F}x_0^V + \mathbf{C}\mathbf{A}^{-1}\mathbf{F}x_0^\perp - \mathbf{C}_r\mathbf{A}_r^{-1}\mathbf{F}_r x_0^V)^T (\mathbf{C}\mathbf{A}^{-1}\mathbf{F}x_0^V + \mathbf{C}\mathbf{A}^{-1}\mathbf{F}x_0^\perp - \mathbf{C}_r\mathbf{A}_r^{-1}\mathbf{F}_r x_0^V) \quad (28)$$

Expression (28) can be approximated by assuming that

$$\mathbf{C}\mathbf{A}^{-1}\mathbf{F}x_0^V = \mathbf{C}_r\mathbf{A}_r^{-1}\mathbf{F}_r x_0^V \quad (29)$$

which means that for initial conditions x_0^V in the space spanned by the basis, we assume that the reduced output exactly matches the full output, i.e. $\mathbf{y} = \mathbf{y}_r$. An approach to satisfying this condition will be described shortly. Using approximation (29), we can rewrite (25) as

$$x_0^* = \arg \max_{x_0^\perp \in \mathcal{X}_0} (x_0^\perp)^T H x_0^\perp \quad (30)$$

where

$$H = (\mathbf{C}\mathbf{A}^{-1}\mathbf{F})^T (\mathbf{C}\mathbf{A}^{-1}\mathbf{F}) \quad (31)$$

$H \in \mathbb{R}^{N \times N}$ is now the Hessian matrix of the full-scale system and does not depend on the reduced-order model. As before, H is a symmetric, positive semi-definite matrix that does not depend on the state or initial condition.

If we choose to initialize the greedy algorithm with an empty basis, $V = 0$, then the maximizer of (30) on the first greedy iteration is given by the eigenvector of H corresponding to the largest eigenvalue. We denote this initial condition by z_1 and note that z_1 satisfies

$$Hz_1 = \lambda_1 z_1 \quad (32)$$

where λ_1 is the largest eigenvalue of H . We then set $V = z_1$. Under the assumption that (29) holds, on the second greedy iteration we would therefore seek the initial condition that maximizes (30). Clearly, this initial condition, which should be orthogonal to z_1 , is given by z_2 , the eigenvector of H corresponding to the second largest eigenvalue.

Returning to assumption (29), this condition can be satisfied if we include in the basis not just the sequence of optimal initial conditions $x_0^* = \{z_1, z_2, \dots\}$, but rather the span of *all* snapshots (i.e. instantaneous state solutions contained in \mathbf{x}) obtained by solving (8) for each of the seed initial conditions z_1, z_2, \dots . Approximation (29) will then be accurate, provided the final time t_T is chosen so that the output $y(k)$ is small for $k > T$. If the output is not small for $k > T$, then a snapshot collected at some time $t_{\bar{k}}$, where $\bar{k} < T$ but \bar{k} is large, will be added to the basis; however, if that state were then used as an initial condition in the resulting reduced-order model, the resulting solution \mathbf{y}_r would not necessarily be an accurate representation of \mathbf{y} . This is because the basis would not contain information about system state evolution after time $t_{T-\bar{k}}$. In that case, (29) would not hold. Further, by including both the initial conditions, z_i , and the corresponding snapshots, \mathbf{x} , in the basis, the sequence of eigenvectors z_i will no longer satisfy the necessary orthogonality conditions; that is, the second eigenvector z_2 may no longer be orthogonal to the space spanned by the basis comprising z_1 and its corresponding state solutions. This is because setting $x_0 = z_1$ and computing \mathbf{x} will likely lead to some states that have components in the direction of z_2 . We would therefore expect this simplification to be more accurate for the first few eigenvectors and become less accurate as the number of seed initial conditions is increased.

These simplifications lead us to an alternate ‘one-shot’ basis-construction algorithm for the initial-condition problem. This algorithm does not solve the optimization problems (15)–(19) or (20)–(24) exactly, but provides a good approximate solution to the problem (20)–(24) under the conditions discussed above. We use the dominant eigenvectors of the Hessian matrix H to identify the initial-condition vectors that have the most significant contributions to the outputs of interest. These vectors are in turn used to initialize the full-scale discrete-time system to generate a set of state snapshots that are used to form the reduced basis.

Algorithm 2

One-Shot Hessian-Based Reduced-Basis Construction.

1. For the full-order Hessian matrix, H as defined in (31), find the p eigenvectors z_1, z_2, \dots, z_p with largest eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq \lambda_{p+1} \geq \dots \geq \lambda_N \geq 0$.
2. For $i = 1, \dots, p$, set $x_0 = z_i$ and compute the corresponding solution \mathbf{x}^i using (8).
3. Form the reduced basis as the span of the snapshots $\mathbf{x}^i(k)$, $i = 1, 2, \dots, p$, $k = 0, 1, \dots, T$.

Steps 2 and 3 in Algorithm 2 allow us to (approximately) satisfy assumption (29) by including not just the initial conditions z_1, z_2, \dots, z_p in the basis but also the span of all snapshots generated from those initial conditions. The basis could be computed from the snapshots, using, for example, the POD.

3.2. Error analysis

A direct measure of the quality of the reduced-order model is available using the analysis framework described above. We define the error, ε , due to a particular initial condition x_0 as

$$\varepsilon = \|\mathbf{y} - \mathbf{y}_r\|_2 = \|(\mathbf{C}\mathbf{A}^{-1}\mathbf{F} - \mathbf{C}_r\mathbf{A}_r^{-1}\mathbf{F}_r)x_0\|_2 \quad (33)$$

For a given reduced model, the dominant eigenvector of H^e provides the worst-case initial condition. Therefore, the value of the maximal error ε_{\max} (for an initial condition of unit norm) is given by

$$\varepsilon_{\max} = \sqrt{\lambda_1^e} \quad (34)$$

where λ_1^e is the largest eigenvalue of the error Hessian H^e defined by (26). The value ε_{\max} provides both a measure on the quality of the reduced model and a quantitative termination criterion for the basis-construction algorithm.

In Algorithm 1, ε_{\max} is readily available and thus can be used to determine how many cycles of the algorithm to perform, i.e. the algorithm would be terminated when the worst-case error is sufficiently small. In Algorithm 2, it is computationally more efficient to select p , the number of seed initial-conditions, based on the decay rate of the full Hessian eigenvalues $\lambda_1, \lambda_2, \dots$ and to compute all the necessary eigenvectors z_1, z_2, \dots, z_p at once. Once the reduced model has been created using Algorithm 2, the error Hessian H^e can be formed and the error criterion (34) checked to determine if further sampling is required. While Algorithm 1 is expected to reduce the worst-case error more quickly, the one-shot Algorithm 2 is attractive since it depends only on the large-scale system properties and thus does not require us to build the reduced model on each cycle.

We also note that the eigenvectors of $H = (\mathbf{C}\mathbf{A}^{-1}\mathbf{F})^T(\mathbf{C}\mathbf{A}^{-1}\mathbf{F})$ are equivalent to the (right) singular vectors of $\mathbf{C}\mathbf{A}^{-1}\mathbf{F}$. Since the latter quantity serves as an input–output mapping, use of its singular vectors for basis formation is intuitively attractive.

It is also interesting to note that the Hessian H may be thought of as a finite-time observability gramian [27].

3.3. Large-scale implementation

We first discuss the implementation of Algorithm 2 in the large-scale setting, and then remark on the differences for Algorithm 1.

Algorithm 2 is a one-shot approach in which all of the eigenpairs can be computed from the single Hessian matrix H in (31). This matrix can be formed explicitly by first forming $\mathbf{A}^{-1}\mathbf{F}$, which requires N ‘forward solves’ (i.e. solutions of forward-in-time dynamical systems with \mathbf{A} as coefficient matrix), where N is the number of initial-condition parameters; or else by first forming $\mathbf{A}^{-T}\mathbf{C}^T$, which requires Q ‘adjoint’ solves (i.e. solutions of backward-in-time dynamical systems with \mathbf{A}^T as coefficient matrix), where Q is the number of outputs. For large-scale problems with high-dimensional initial condition and output vectors, explicit formation and storage of H is thus intractable. (A similar argument can be made for the intractability of computing the singular value decomposition of $\mathbf{C}\mathbf{A}^{-1}\mathbf{F}$.) Even if H could be formed and stored, computing its dominant spectrum would be prohibitive since it is a dense matrix of order $N \times N$.

Instead, we use a matrix-free iterative method such as Lanczos to solve for the dominant eigenpairs of H . Such methods require at each iteration a matrix–vector product of the form

Hw_k for some w_k , which is formed by successive multiplication of vectors with the component matrices that make up the Hessian in (31). At each iteration, this amounts to one forward and one adjoint solve involving the system \mathbf{A} . When the eigenvalues are well separated, convergence to the largest eigenvalues of H is rapid. Moreover, when the spectrum decays rapidly, only a handful of eigenvectors are required by Algorithm 2. Many problems have Hessian matrices that are of low rank and spectra that decay rapidly, stemming from the limited number of initial conditions that have a significant effect on outputs of interest. For such problems the number of Lanczos iterations required to extract the dominant part of the spectrum is often independent of the problem size N .

Under this assumption, we can estimate the cost of Step 1 of Algorithm 2 (which dominates the cost) in the case when the dynamical system (8)–(9) stems from a discretized parabolic partial differential equation (PDE). The cost of each implicit time step of a forward or adjoint solve is usually linear or weakly superlinear in problem size, using modern multilevel preconditioned linear solvers. Therefore for T time steps, overall work for a forward or adjoint solve scales as $TN^{1+\alpha}$, with α usually very small. For a 3-D spatial problem, a number of time steps on the order of the diameter of the grid, and an optimal preconditioner, gives $\mathcal{O}(N^{4/3})$ complexity per forward solve, and hence per Lanczos iteration. Assuming that the number of Lanczos iterations necessary to extract the dominant part of the spectrum is independent of the grid size, the overall complexity remains $\mathcal{O}(N^{4/3})$. (Compare this with straightforward formation of the Hessian and computation of the eigenvalues with the QR algorithm which requires $\mathcal{O}(N^3)$ work.)

Algorithm 1 is implemented in much the same way. The main difference is that the error Hessian H^e replaces the Hessian H , and we find the dominant eigenpair of each of a sequence of eigenvalue problems, rather than finding p eigenpairs of the single Hessian H . Each iteration of a Lanczos-type solver for the eigenvalue problem in Algorithm 1 resembles that of Algorithm 2, and therefore the costs per iteration are asymptotically the same. It is more difficult to characterize the number of greedy iterations, and hence the number of eigenvector problems, that will be required using Algorithm 1. However, to the extent that the assumptions outlined in Section 3.1 hold, the number of greedy iterations will correspond roughly to the number of dominant eigenvalues of the full Hessian matrix H . As reasoned above, the spectrum of H is expected to decay rapidly for the problems of interest here; thus, convergence of the greedy reduced-basis construction algorithm is expected to be rapid.

4. APPLICATION TO A 2D CONVECTION–DIFFUSION TRANSPORT PROBLEM

In this section, the model reduction methodology described above is assessed for a contaminant transport problem. The physical process is modeled by the convection–diffusion equation

$$\frac{\partial w}{\partial t} + \mathbf{v} \cdot \nabla w - \kappa \nabla^2 w = 0 \quad \text{in } \Omega \times (0, t_f) \quad (35)$$

$$w = 0 \quad \text{on } \Gamma_D \times (0, t_f) \quad (36)$$

$$\frac{\partial w}{\partial n} = 0 \quad \text{on } \Gamma_N \times (0, t_f) \quad (37)$$

$$w = w_0 \quad \text{in } \Omega \quad \text{for } t = 0 \quad (38)$$

where w is the contaminant concentration (which varies in time and over the domain Ω), \mathbf{v} is the velocity vector field, κ is the diffusivity, t_f is the time horizon of interest, and w_0 is the given initial condition. Homogeneous Dirichlet boundary conditions are applied on the inflow boundary Γ_D , while homogeneous Neumann conditions are applied on the other boundaries Γ_N . We first consider the case of a simple 2-D domain, which leads to a system of the form (8) of moderate dimension; in the next section a large-scale 3-D example will be presented.

4.1. Two-dimensional model problem

Figure 1 shows the computational domain for the 2-D contaminant transport example. The velocity field is taken to be uniform, constant in time, and directed in the positive \bar{x} -direction as defined by Figure 1. The inflow boundary, Γ_D , is defined by $\bar{x} = 0$, $0 \leq \bar{y} \leq 0.4$; the remaining boundaries comprise Γ_N .

A streamline upwind Petrov–Galerkin (SUPG) [28] finite element method is employed to discretize (35) in space using triangular elements. For the cases considered here, the spatial mesh has $N = 1860$ nodes. The Crank–Nicolson method is used to discretize the equations in time. This leads to a linear discrete-time system of the form (8), where the state vector $x(k) \in \mathbb{R}^{1860}$ contains the values of contaminant concentration at spatial grid points at time t_k . For all experiments, the timestep used was $\Delta t = 0.02$ and the time limit, set approximately by the maximum time of convection across the length of the domain, was $t_T = 1.4$.

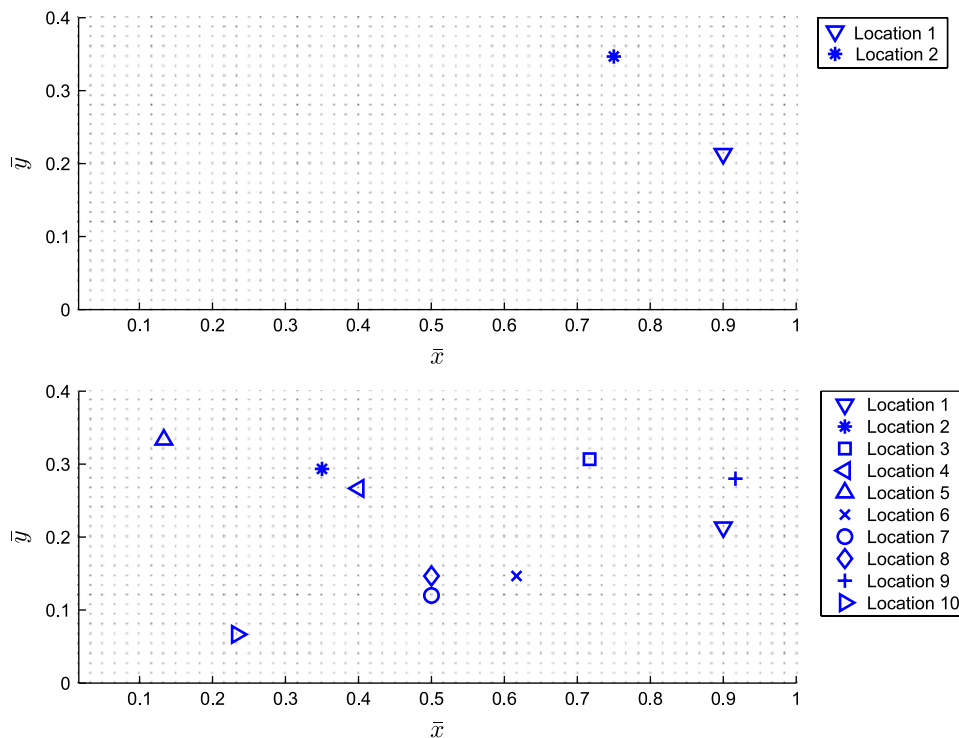


Figure 1. The computational domain and locations of sensor output nodes.
Top: two-sensor case, bottom: 10-sensor case.

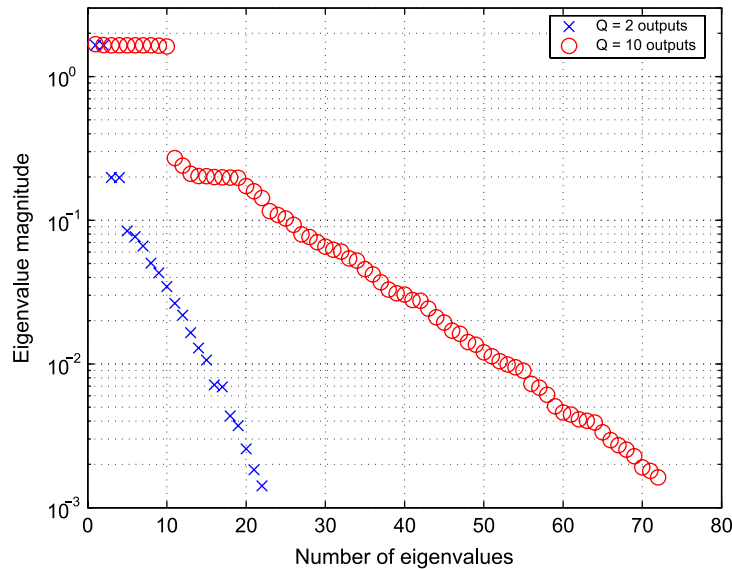


Figure 2. A comparison of the Hessian eigenvalue spectra of H for the two- and 10-output cases. $Pe = 100$.

The matrix \mathbf{A} in (8) depends on the velocity field and the Peclet number, Pe , which is defined as

$$Pe = \frac{v_c \ell_c}{\kappa} \quad (39)$$

where the characteristic velocity v_c is taken to be the maximum velocity magnitude in the domain, while the domain length is used as the characteristic length ℓ_c . The uniform velocity field described above was used in all experiments, but Pe was varied. Increasingly convective transport scenarios corresponding to Peclet numbers of 10, 100, and 1000 were used to generate different full-scale systems.

The outputs of interest are defined to be the values of concentration at selected sensor locations in the computational domain. Figure 1 shows two different sensor configurations that were employed in the results presented here.

The first step in creating a reduced model with Algorithm 2 is to compute p dominant eigenvectors of the full-scale Hessian matrix H . Figure 2 shows the eigenvalue spectra of H for the two-sensor case and the 10-sensor case. The relative decay rates of these eigenvalues are used to determine p , the number of eigenvectors used as seed initial conditions. We specify the parameter $\bar{\lambda}$, and apply the criterion that the j th eigenvector of H is included if $\lambda_j/\lambda_1 > \bar{\lambda}$.

Figure 2 demonstrates that the decay rate of the dominant eigenvalues is related to the number and positioning of output sensors. For the two-output case, the two dominant eigenvalues λ_1 and λ_2 are of almost equal magnitude; analogous behaviour can be seen for the first 10 eigenvalues in the 10-output case. This is consistent with the physical intuition that similarly important modes exist for each of the output sensors. For instance, a mode with initial concentration localized around one particular sensor is of similar importance as another mode with high concentration near a different sensor.

4.2. Reduced-model performance

Once the p seed eigenvectors have been computed, the corresponding state solutions, $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^p$, are computed from (8) using each eigenvector in turn as the initial condition x_0 . The final step in Algorithm 2 requires the formation of the reduced basis from the span of $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^p$. We achieve this by aggregating all state solutions $x^i(k)$, $i = 1, 2, \dots, p$, $k = 0, 1, \dots, T$ into a snapshot matrix $X \in \mathbb{R}^{N \times (T+1)p}$ and using the POD to select the n basis vectors that most efficiently span the column space of X . The number of POD basis vectors is chosen based on the decay of the POD eigenvalues $\mu_1 \geq \mu_2 \geq \dots \geq \mu_{(T+1)p} \geq 0$. As above, we define a parameter $\bar{\mu}$, and apply the criterion that the k th POD basis vector is retained if $\mu_k / \mu_1 > \bar{\mu}$.

The resulting reduced models given by (12), (13) can be used for any initial condition x_0 ; to demonstrate the methodology, we choose to show results for initial conditions comprising a superposition of Gaussian functions. Each Gaussian is defined by

$$x_0(\bar{x}, \bar{y}) = \frac{1}{\sigma\sqrt{2\pi}} e^{-[(\bar{x}-\bar{x}_c)^2 + (\bar{y}-\bar{y}_c)^2]/2\sigma^2} \quad (40)$$

where (\bar{x}_c, \bar{y}_c) defines the center of the Gaussian and σ is the standard deviation. All test initial-conditions are normalized such that $\|x_0\|_2 = 1$. Three sample initial-condition functions that are used in the following analyses are shown in Figure 3 and are referred to by their provided labels (a)–(c) throughout.

Tables I and II show sample reduced-model results for various cases using the two-sensor configuration shown in Figure 1. The error ε is defined in (33) and computed for one of the sample initial conditions shown in Figure 3. It can be seen from the tables that a substantial reduction in the number of states from $N = 1860$ can be achieved with low levels of error in the concentration prediction at the sensor locations. The tables also show that including more modes in the reduced model, either by decreasing the Hessian eigenvalue decay tolerance $\bar{\lambda}$ or by decreasing the POD eigenvalue decay tolerance $\bar{\mu}$, leads to a reduction in the output error. Furthermore, the worst-case error in each case, ε_{\max} , is computed from (34) using the maximal eigenvalue of the error Hessian, H^e . It can also be seen that inclusion of more modes in the reduced model leads to a reduction in the worst-case error, although the reduction in ε_{\max} occurs more slowly than the reduction in ε .

Figure 4 shows a comparison between reduced models computed using Algorithms 1 and 2. The figure highlights the result shown in Table I; that is, using the one-shot approach, the maximum error decreases rather slowly as the size of the model increases. However, the figure also shows that the actual error for the same model (shown in this case for test initial condition (a)) is significantly reduced as n increases. This suggests that while subsequent eigenvectors of the full-scale Hessian may not directly target the worst-case initial condition, they do add useful information to the basis. Conversely, Figure 4 shows that Algorithm 1, which uses the successive dominant eigenvector of the error Hessian, does directly target the worst-case error. However, it can also be seen that reductions in the worst-case error for a reduced model do not necessarily translate into reductions in the error observed for a particular initial condition. For this problem, the cost of computing the first eigenvector is substantially higher than the cost of computing subsequent eigenvectors, making Algorithm 2 more efficient than Algorithm 1. For example, the results in Table II correspond to $p = 5$ ($\bar{\lambda} = 0.1$), $p = 14$ ($\bar{\lambda} = 0.01$), and $p = 22$ ($\bar{\lambda} = 0.001$) seed eigenvectors, with relative costs of 1, 1.12, and 1.42, respectively. Thus, the improvements in reduced-model accuracy seen in Table II are obtained with relatively small increases in offline cost; however, this result is not expected to hold for larger-scale problems where the overhead is much smaller than the cost of computing

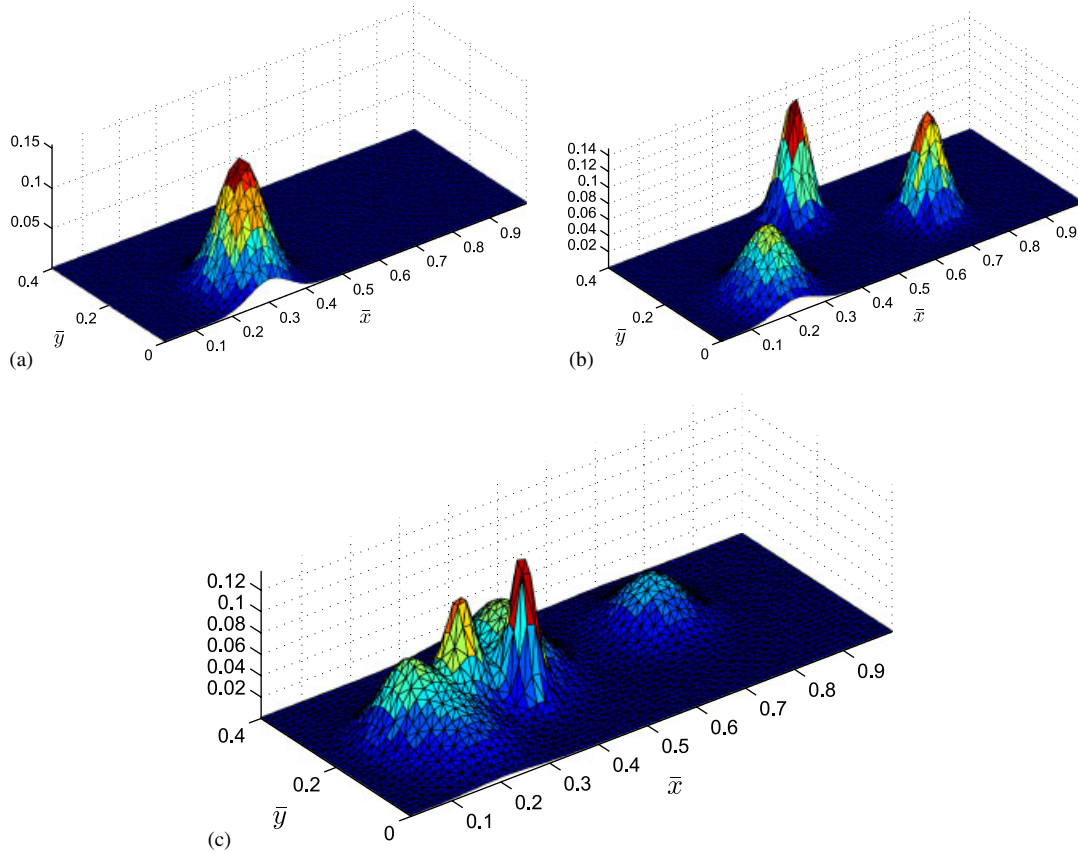


Figure 3. Sample test initial conditions used to compare reduced-model outputs to full-scale outputs: (a) single Gaussian; (b) superposition of three Gaussians; and (c) superposition of seven Gaussians.

Table I. Properties of various reduced-order models of a full-scale system with $Pe = 10$ and two output sensors.

Case	$\bar{\lambda}$	$\bar{\mu}$	n	ε	ε_{\max}
1	0.1	10^{-4}	28	0.0573	0.4845
2	0.1	10^{-6}	45	0.0103	0.4838
3	0.01	10^{-4}	43	0.0237	0.4758
4	0.01	10^{-6}	69	0.0021	0.4752
5	0.001	10^{-4}	79	0.0017	0.4735
6	0.001	10^{-6}	122	0.0007	0.4418

Note: The errors ε and ε_{\max} are defined in (33) and (34), respectively; ε is evaluated when each reduced system (of dimension n) is subjected to test initial condition (a).

Table II. Properties of various reduced-order models of a full-scale system with $Pe = 100$ and two output sensors.

Case	$\bar{\lambda}$	$\bar{\mu}$	n	ε	ε_{\max}
1	0.1	10^{-4}	62	0.0738	0.1920
2	0.1	10^{-6}	90	0.0722	0.1892
3	0.01	10^{-4}	128	0.0032	0.1638
4	0.01	10^{-6}	200	0.0017	0.1604
5	0.001	10^{-4}	180	0.0004	0.1623
6	0.001	10^{-6}	282	0.0002	0.1564

Note: The errors ε and ε_{\max} are defined in (33) and (34), respectively; ε is evaluated when each reduced system (of dimension n) is subjected to test initial condition (c).

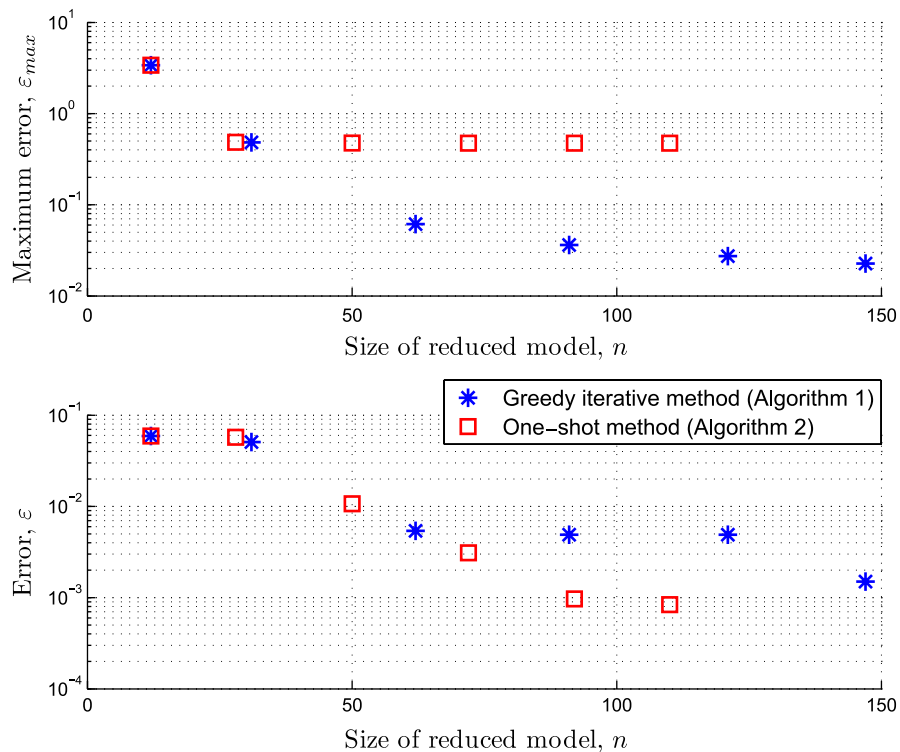


Figure 4. Top: maximum error, ε_{\max} , for reduced models computed using Algorithms 1 and 2. Bottom: error for test initial condition (a), ε , using the same reduced models.

each additional eigenvector. For the results that follow, all reduced models were created using Algorithm 2.

A representative comparison of full and reduced outputs, created by driving both the full and reduced systems with test initial condition (b), is shown in Figure 5 for the case of $Pe = 1000$.

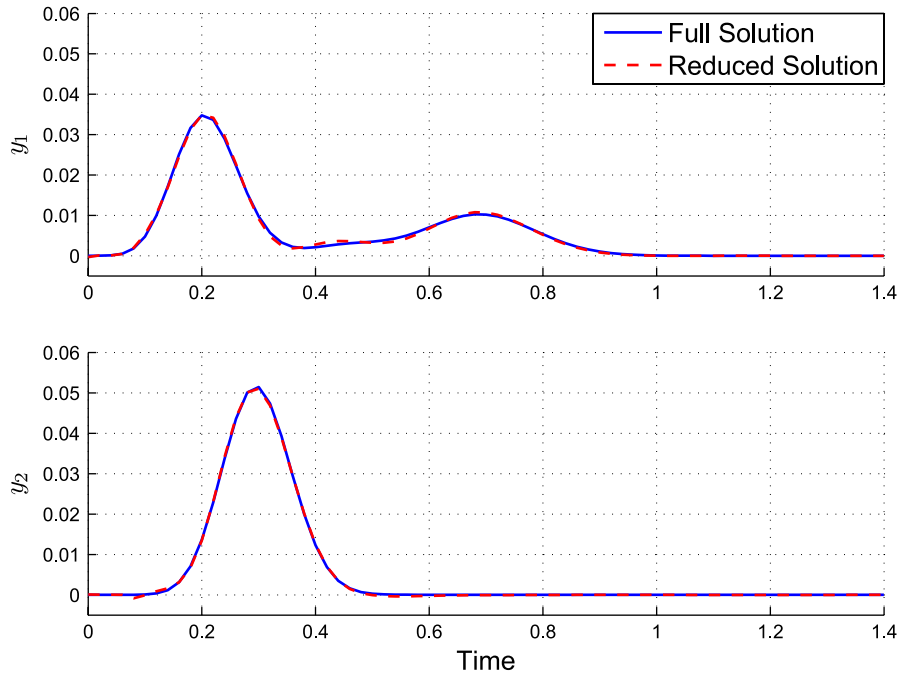


Figure 5. A comparison of full ($N = 1860$) and reduced ($n = 196$) outputs for two-sensor case using test initial condition (b). $Pe = 1000$, $\bar{\lambda} = 0.01$, $\bar{\mu} = 10^{-4}$, $\varepsilon = 0.0036$.

The values $\bar{\lambda} = 0.01$ and $\bar{\mu} = 10^{-4}$ are used, leading to a reduced model of size $n = 196$. The figure demonstrates that a reduced model of size $n = 196$ formed using Algorithm 2 can effectively replicate the outputs of the full-scale system for this initial condition. The error for this case as defined in (33) is $\varepsilon = 0.0036$.

In order to ensure that the results shown in Figure 5 are representative, 1000 initial conditions are constructed randomly and tested using this reduced model. Each initial condition consists of 10 superposed Gaussian functions with random centers (\bar{x}_c , \bar{y}_c) and random standard deviations σ . This library of test initial conditions was used to generate output comparisons between the full-scale model and the reduced-order model. The averaged error across all 1000 trials, $\bar{\varepsilon} = 0.0023$, is close to the error associated with the comparison shown in Figure 5. Furthermore, the maximum error over all 1000 trials is found to be 0.0056, which is well below the upper bound $\varepsilon_{\max} = 0.0829$ established by (34).

Effect of variations in $\bar{\mu}$: As discussed above, $\bar{\mu}$ is the parameter that controls the number of POD vectors n chosen for inclusion in the reduced basis. If $\bar{\mu}$ is too large, the reduced basis will not span the space of all initial conditions for which it is desired that the reduced-model be valid. Figure 6 illustrates the effect of changing $\bar{\mu}$. The curve corresponding to a value of $\bar{\mu} = 10^{-6}$ shows a clear improvement over the $\bar{\mu} = 10^{-4}$ case. This can also be seen by comparing the errors listed in the first two rows of Table I, which correspond to the two reduced models seen in Figure 6. However, the improvement comes at a price, since the number of basis vectors, and therefore the size of the reduced model n , increases from 43 to 69 when $\bar{\mu}$ is decreased.

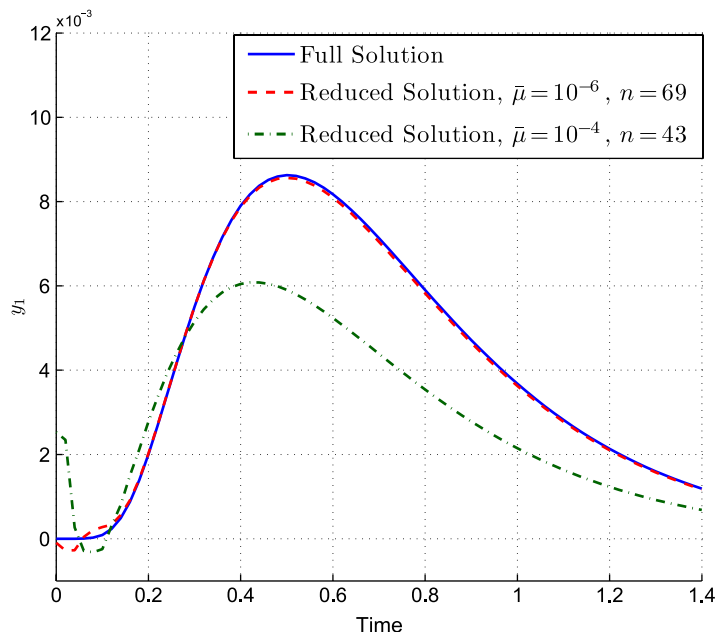


Figure 6. A comparison between full and reduced solutions at sensor location 1 for two different values of $\bar{\mu}$. Test initial condition (a) was used to generate the data. $Pe = 10$, $\bar{\lambda} = 0.1$, two-sensor case.

Effect of variations in $\bar{\lambda}$: Another way to alter the size and quality of the reduced model is to indirectly change p , the number of eigenvectors of H that are used as seed initial conditions for basis creation. We accomplish this by choosing different values of the eigenvalue decay ratio $\bar{\lambda}$. The effect of doing so is illustrated in Figure 7. An increase in reduced-model quality clearly accompanies a decrease in $\bar{\lambda}$. This can also be seen by comparing rows 1 and 3 of Table II, which correspond to the two reduced models seen in Figure 7. The increase in n with lower values of $\bar{\lambda}$ is expected, since greater p implies more snapshot data with which to build the reduced basis, effectively uncovering more full system modes and decreasing the relative importance of the most dominant POD vectors. In general, for the same value of $\bar{\mu}$, more POD vectors are included in the basis if $\bar{\lambda}$ is reduced.

4.3. Ten-sensor case

To understand how the proposed method scales with the number of outputs in the system, we repeat the experiments for systems with $Q = 10$ outputs corresponding to sensors in the randomly generated locations shown in Figure 1. A reduced model was created for the case of $Pe = 100$, with $\bar{\mu} = 10^{-4}$ and $\bar{\lambda} = 0.1$. The result was a reduced system of size $n = 245$, which was able to effectively replicate all 10 outputs of the full system. Figure 8 shows a representative result of the full and reduced-model predictions at all 10 sensor locations.

The size $n = 245$ of the reduced model in this case is considerably larger than that in the corresponding two-output case ($n = 62$), which is shown in the first row of Table II, although both models were constructed with identical values of $\bar{\mu}$ and $\bar{\lambda}$. The difference between high- and low- Q experiments is related to the Hessian eigenvalue spectrum. As demonstrated in Figure 2,

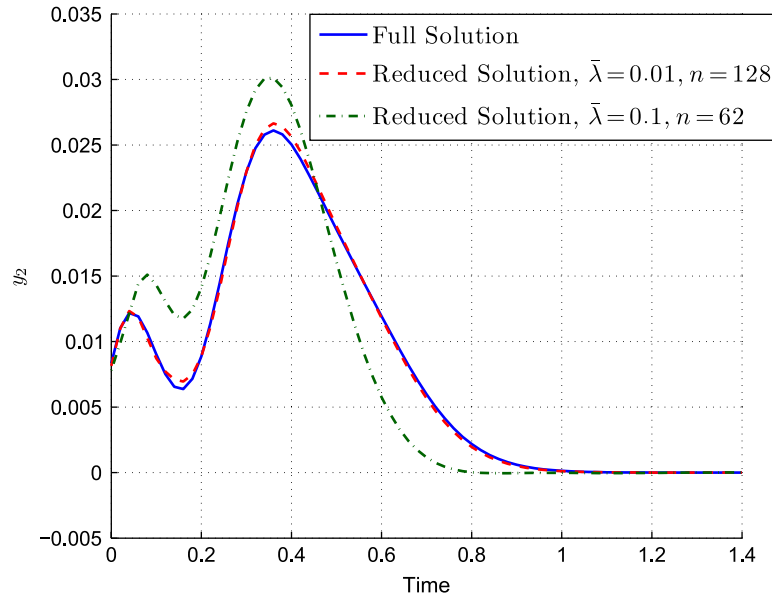


Figure 7. Lowering $\bar{\lambda}$ to increase p , the number of Hessian eigenvector initial conditions used in basis formation, leads to more accurate reduced-order output. Test initial condition (c) was used with two output sensors, $Pe = 100$ and $\bar{\mu} = 10^{-4}$. The output at the second sensor location is plotted here.

the eigenvalue decay rate of the $Q = 10$ case is less rapid than that of the $Q = 2$ case. This means that, for the same value of $\bar{\lambda}$, more seed initial conditions are generally required for systems with more outputs. Since additional modes of the full system must be captured by the reduced model if the number of sensors is increased, it is not surprising that the size of the reduced basis increases.

4.4. Observations and recommendations

The above results demonstrate that reduced models formed by the proposed method can be effective in replicating full-scale output quantities of interest. At this point, we can use the results to make recommendations about choosing $\bar{\mu}$ and $\bar{\lambda}$, the two parameters that control reduced-model construction.

In practice, one would like to choose these parameters such that both the reduced-model size n and the modeling error for a variety of test initial conditions are minimal. The size of the reduced model is important because n is directly related to the online computational cost, that is, n determines the time needed to compute reduced output approximations, which is required to be minimal for real-time applications. The offline cost of forming the reduced model is also a function of $\bar{\mu}$ and $\bar{\lambda}$. When $\bar{\mu}$ is decreased, the basis formation algorithm requires more POD basis vectors to be computed; thus, decreasing μ increases the offline cost of model construction. In addition, the online cost of solving reduced system in (12) and (13), which is not sparse, scales as n^2T . While decreasing $\bar{\mu}$ might appreciably improve modeling accuracy, doing so can only increase the time needed to compute reduced output approximations. Changes in $\bar{\lambda}$ affect the offline cost more strongly. Every additional eigenvector of H to be calculated adds the cost of several additional

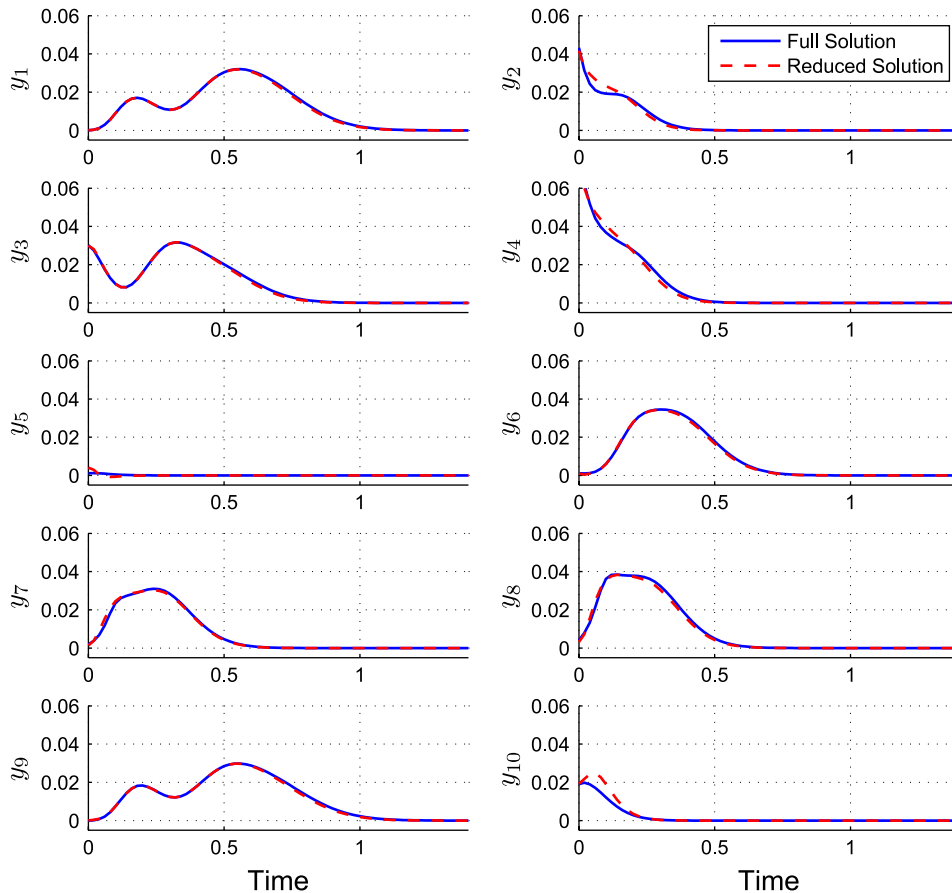


Figure 8. A comparison of the full ($N = 1860$) and reduced ($n = 245$) outputs for all $Q = 10$ locations of interest. Test initial condition (c) was used to generate these data with $Pe = 100$, $\bar{\mu} = 10^{-4}$, $\bar{\lambda} = 0.1$.

large-scale system solves: several forward and adjoint solves are needed to find an eigenvector using the matrix-free Lanczos solver described earlier. In addition, the number of columns of the POD snapshot matrix X grows by $(T + 1)$ if p is incremented by one; computing the POD basis thus becomes more expensive. If these increases in offline cost can be tolerated, though, the results suggest a clear improvement in reduced-model accuracy for a relatively small increase in online cost.

Figure 9 illustrates the dependence of reduced model size and quality on the parameters $\bar{\mu}$ and $\bar{\lambda}$. For the case of 10 output sensors with $Pe = 100$, six different reduced models were constructed with different combinations of $\bar{\mu}$ and $\bar{\lambda}$. The three plots in Figure 9 show the error ε versus the reduced-model size n for each of the test initial conditions in Figure 3. Ideally, a reduced model should have both small error and small n , so we prefer those models whose points reside closest to the origin. Ignoring differences in offline model construction cost, decreasing $\bar{\lambda}$ should be favoured over decreasing $\bar{\mu}$ if more accuracy is desired. This conclusion is reached by realizing that for a

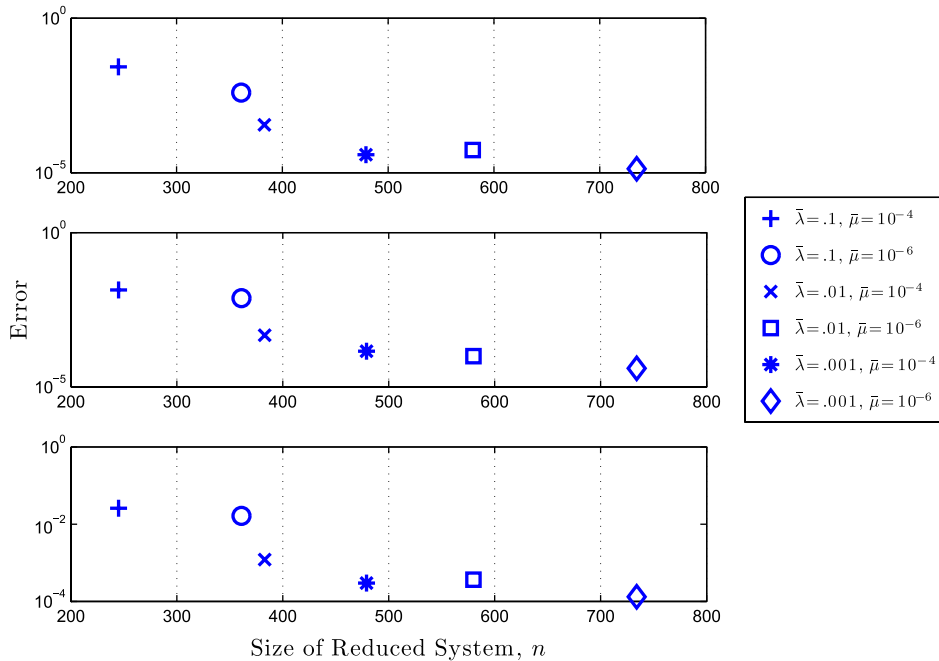


Figure 9. A measure of the error in six different reduced models of the same system plotted *versus* their sizes n for the 10-sensor case. The three plots were generated with test initial conditions (a), (b), and (c), respectively. $Pe = 100$, $Q = 10$ outputs.

comparable level of error, reduced models constructed with lower values of $\bar{\lambda}$ are much smaller. Maintaining a small size of the reduced model is important for achieving real-time computations for large-scale problems of practical interest, as discussed in the next section.

5. APPLICATION: MODEL REDUCTION FOR A 3-D CONTAMINANT TRANSPORT IN AN URBAN CANYON

We demonstrate our model reduction method by applying it to a 3-D airborne contaminant transport problem for which a solution is needed in real time. Intentional or unintentional chemical, biological, and radiological (CBR) contamination events are important national security concerns. In particular, if contamination occurs in or near a populated area, predictive tools are needed to rapidly and accurately forecast the contaminant spread to provide decision support for emergency response efforts. Urban areas are geometrically complex and require detailed spatial discretization to resolve the relevant flow and transport, making prediction in real-time difficult. Reduced-order models can play an important role in facilitating real-time turn-around, in particular, on laptops in the field. However, it is essential that these reduced models be faithful over a wide range of initial conditions, since in principle any of these initial conditions can be realized. Once a suitable reduced-order model has been generated, it can serve as a surrogate for the full model within an

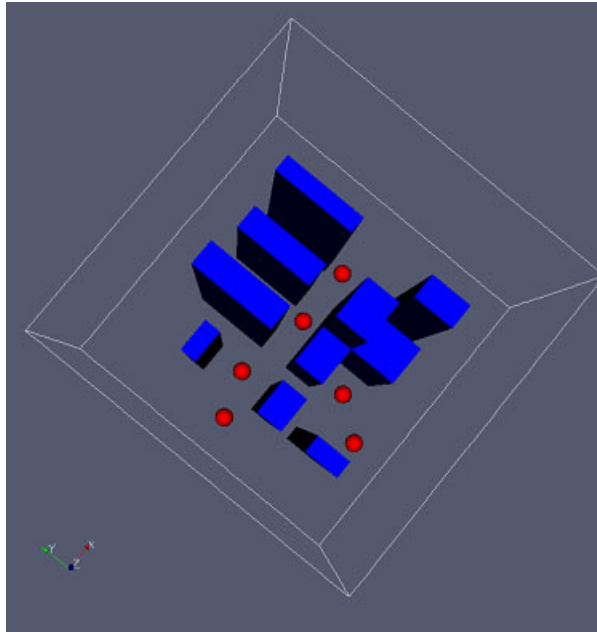


Figure 10. Building geometry and locations of outputs for the 3-D urban canyon problem.

inversion/data assimilation framework to identify the initial conditions given sensor data (see the discussion of the inverse problem in the full-scale case in [29]).

To illustrate the generation of a reduced-order model that is accurate for high-dimensional initial conditions, we consider a 3-D urban canyon geometry occupying a (dimensionless) $15 \times 15 \times 15$ domain. Figure 10 shows the domain and buildings, along with the locations of six output nodes that represent sensor locations of interest, all placed at a height of 1.5. The model used is again the convection–diffusion equation, given by (35). The PDE is discretized in space using an SUPG finite element method with linear tetrahedra, while the Crank–Nicolson method is used to discretize in time. Homogeneous Dirichlet boundary conditions of the form (36) are specified for the concentration on the inflow boundary, $\bar{x} = 0$, and the ground, $\bar{z} = 0$. Homogeneous Neumann boundary conditions of the form (37) are specified for the concentration on all other boundaries.

The velocity field, \mathbf{v} , required in (35) is computed by solving the steady laminar incompressible Navier–Stokes equations, also discretized with SUPG-stabilized linear tetrahedra. No-slip conditions, i.e. $\mathbf{v} = 0$, are imposed on the building faces and the ground $\bar{z} = 0$ (thus there is no flow inside the buildings). The velocity at the inflow boundary $\bar{x} = 0$ is taken as known and specified in the normal direction as

$$v_x(z) = v_{\max} \left(\frac{z}{z_{\max}} \right)^{0.5}$$

with $v_{\max} = 3.0$ and $z_{\max} = 15$, and zero tangentially. On the outflow boundary $\bar{x} = 15$, a traction-free (Neumann) condition is applied. On all other boundaries ($\bar{y} = 0$, $\bar{y} = 15$, $\bar{z} = 15$), we impose a combination of no flow normal to the boundary and traction-free tangent to the boundary. The

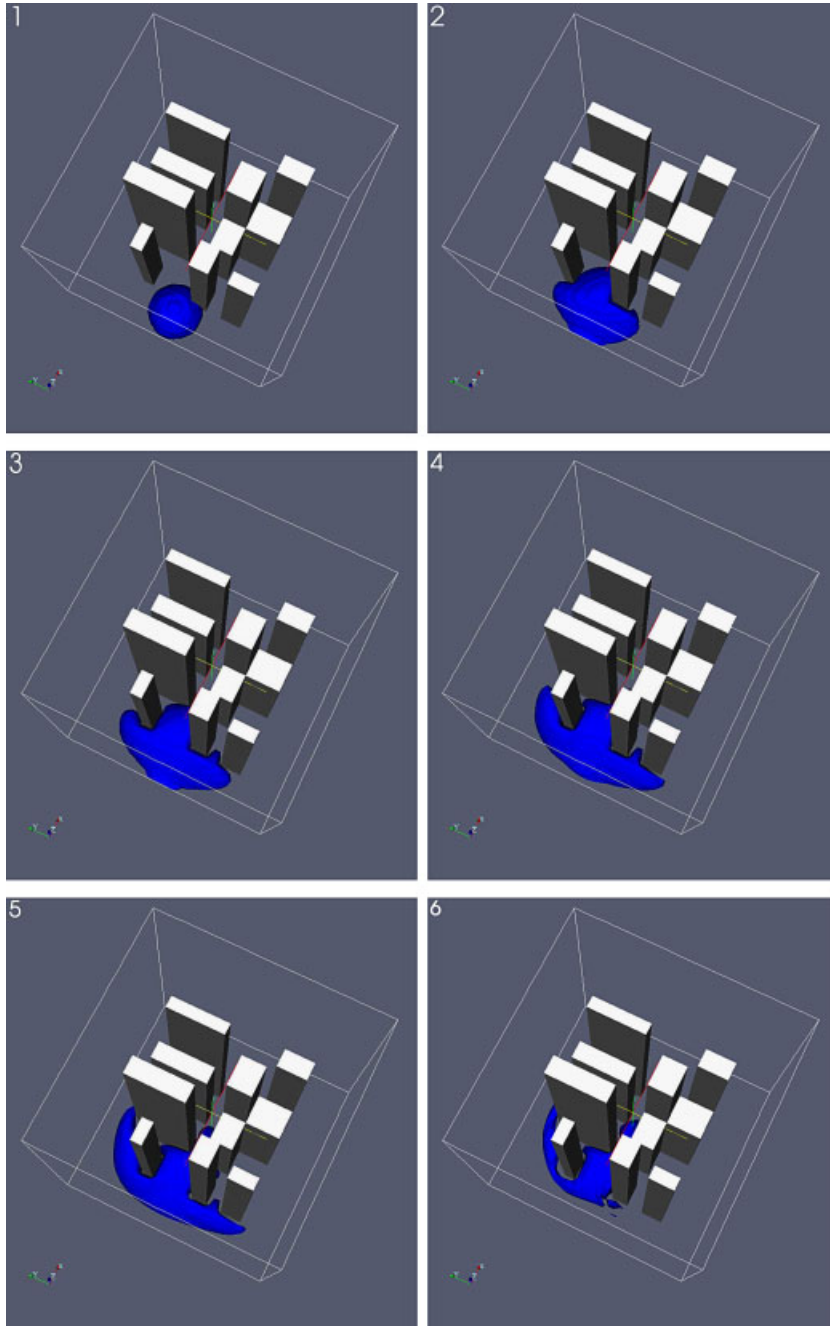


Figure 11. Transport of contaminant concentration through urban canyon at six different instants in time, beginning with the initial condition shown in upper left.

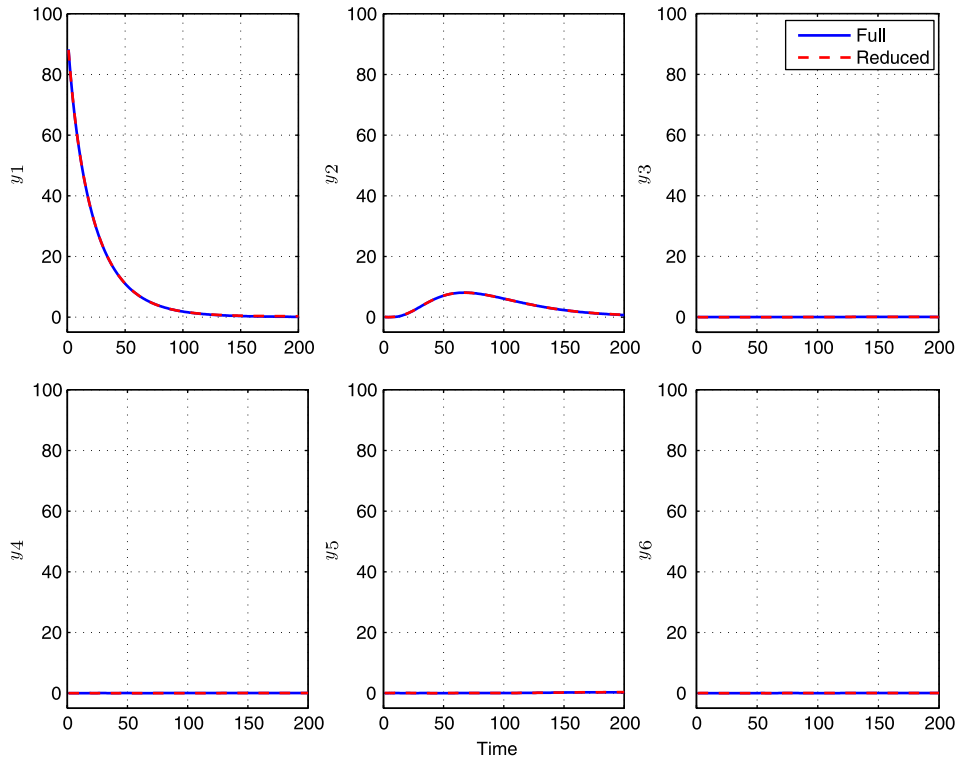


Figure 12. Full (65 600 states) and reduced (137 states) model contaminant concentration predictions at each of the six output nodes for the three-dimensional urban canyon example.

spatial mesh for the full-scale system contains 68 921 nodes and 64 000 tetrahedral elements. For both basis creation and testing, a final non-dimensional time $t_f = 20.0$ is used, and discretized over 200 timesteps. The Peclet number based on the maximum inflow velocity and domain dimension is $Pe = 900$. The PETSc library [30–32] is used for all implementation.

Figure 11 illustrates a sample forward solution. The test initial condition used in this simulation, meant to represent the system state just after a contaminant release event, was constructed using a Gaussian function with a peak magnitude of 100 centered at a height of 1.5.

For comparison with the full system, a reduced model was constructed using Algorithm 2 with the eigenvalue decay ratios $\bar{\lambda} = 0.005$ and $\bar{\mu} = 10^{-5}$, which led to $p = 31$ eigenvector initial conditions and $n = 137$ reduced-basis vectors. Eigenvectors were computed using the Arnoldi eigensolver within the SLEPc package [33], which is built on PETSc. Figure 12 shows a comparison of the full and reduced time history of concentration at each output location. The figure demonstrates that a reduced system of size $n = 137$, which is solved in a matter of seconds on a desktop, can accurately replicate the outputs of the full-scale system of size $N = 65\,600$. We emphasize that the (offline) construction of the reduced-order model targets only the specified outputs, and otherwise has no knowledge of the initial conditions used in the test of Figure 12 (or any other initial conditions).

6. CONCLUSIONS

A new method has been proposed for constructing reduced-order models of linear systems that are parametrized by initial conditions of high dimension. Formulating the greedy approach to sampling as a model-constrained optimization problem, we show that the dominant eigenvectors of the resulting Hessian matrix provide an explicit solution to the greedy optimization problem. This result leads to an algorithm to construct the reduced basis in an efficient and systematic way, and further, provides quantification of the worst-case error in reduced-model output prediction. Thus, the resulting reduced models are guaranteed to provide accurate replication of full-scale output quantities of interest for any possible initial condition, making them appropriate for use in an inverse problem/data assimilation setting. The adaptive greedy sampling approach combined with the model-constrained optimization formulation provides a general framework that is applicable to non-linear problems, although the explicit solution and maximal error guarantees apply only in the linear case. Further, we note that the task of sampling system inputs (which here were taken to be zero) to build a basis over the input space could also be formulated as a greedy optimization problem.

ACKNOWLEDGEMENTS

This work was partially supported by the National Science Foundation under DDDAS grants CNS-0540372 and CNS-0540186 (program director Dr Frederica Darema) and the Air Force Office of Scientific Research (program manager Dr Fariba Fahroo), and the Computer Science Research Institute at Sandia National Laboratories.

REFERENCES

1. Feldmann P, Freund RW. Efficient linear circuit analysis by Padé approximation via the Lanczos process. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 1995; **14**:639–649.
2. Gallivan K, Grimme E, Van Dooren P. Padé approximation of large-scale dynamic systems with Lanczos methods. *Proceedings of the 33rd IEEE Conference on Decision and Control*, Lake Buena Vista, FL, December 1994.
3. Grimme E. Krylov projection methods for model reduction. *Ph.D. Thesis*, Coordinated-Science Laboratory, University of Illinois at Urbana-Champaign, 1997.
4. Gugercin S, Antoulas A. A survey of model reduction by balanced truncation and some new results. *International Journal of Control* 2004; **77**:748–766.
5. Li J, White J. Low rank solution of Lyapunov equations. *SIAM Journal on Matrix Analysis and Applications* 2002; **24**(1):260–280.
6. Penzl T. Algorithms for model reduction of large dynamical systems. *Linear Algebra and its Applications* 2006; **415**(2–3):322–343.
7. Sorensen DC, Antoulas AC. The Sylvester equation and approximate balanced reduction. *Linear Algebra and its Applications* 2002; **351–352**:671–700.
8. Deane AE, Kevrekidis IG, Karniadakis GE, Orszag SA. Low-dimensional models for complex geometry flows: application to grooved channels and circular cylinders. *Physics of Fluids* 1991; **3**(10):2337–2354.
9. Holmes P, Lumley JL, Berkooz G. *Turbulence, Coherent Structures, Dynamical Systems and Symmetry*. Cambridge University Press: Cambridge, U.K., 1996.
10. Sirovich L. Turbulence and the dynamics of coherent structures. Part 1: coherent structures. *Quarterly of Applied Mathematics* 1987; **45**(3):561–571.
11. Afanasiev K, Hinze M. *Adaptive Control of a Wake Flow Using Proper Orthogonal Decomposition*. Lecture Notes in Pure and Applied Mathematics, vol. 216. Marcel Dekker: New York, 2001; 317–332.
12. Arian E, Fahl M, Sachs EW. Trust-region proper orthogonal decomposition for optimal flow control. *Technical Report ICASE 2000-25*, Institute for Computer Applications in Science and Engineering, May 2000.

13. Hinze M, Volkwein S. Proper orthogonal decomposition surrogate models for nonlinear dynamical systems: error estimates and suboptimal control. In *Dimension Reduction of Large-Scale Systems*, Benner P, Mehrmann V, Sorensen D (eds), Lecture Notes in Computational and Applied Mathematics. Springer: Berlin, Germany, 2005; 261–306.
14. Kunisch K, Volkwein S. Control of Burgers' equation by reduced order approach using proper orthogonal decomposition. *Journal of Optimization Theory and Applications* 1999; **102**:345–371.
15. Daniel L, Siong OC, Chay LS, Lee KH, White J. Multiparameter moment matching model reduction approach for generating geometrically parameterized interconnect performance models. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 2004; **23**(5):678–693.
16. Farrell B, Ioannou P. Accurate low-dimensional approximation of the linear dynamics of fluid flow. *Journal of the Atmospheric Sciences* 2001; **58**:2771–2789.
17. Daescu DN, Navon IM. Efficiency of a POD-based reduced second order adjoint model in 4D-Var data assimilation. *International Journal for Numerical Methods in Fluids* 2007; **53**:985–1004.
18. Gugercin S, Antoulas A, Beattie C. A rational Krylov iteration for optimal H2 model reduction. *Proceedings of MTNS 2006*, Japan, 2006.
19. Silveira L, Phillips J. Resampling plans for sample point selection in multipoint model-order reduction. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 2006; **25**(12):2775–2783.
20. Veroy K, Prud'homme C, Rovas D, Patera A. *A posteriori* error bounds for reduced-basis approximation of parametrized noncoercive and nonlinear elliptic partial differential equations. *AIAA Paper 2003-3847, Proceedings of the 16th AIAA Computational Fluid Dynamics Conference*, Orlando, FL, 2003.
21. Veroy K, Patera A. Certified real-time solution of the parametrized steady incompressible Navier–Stokes equations: rigorous reduced-basis *a posteriori* error bounds. *International Journal for Numerical Methods in Fluids* 2005; **47**:773–788.
22. Grepl M, Patera A. *A posteriori* error bounds for reduced-basis approximations of parametrized parabolic partial differential equations. *ESAIM—Mathematical Modelling and Numerical Analysis (M2AN)* 2005; **39**(1):157–181.
23. Grepl MA. Reduced-basis approximation and *a posteriori* error estimation for parabolic partial differential equations. *Ph.D. Thesis*, Massachusetts Institute of Technology, 2005.
24. Bui-thanh T, Willcox K, Ghattas O, van Bloemen Waanders B. Goal-oriented, model-constrained optimization for reduction of large-scale systems. *Journal of Computational Physics* 2006, in press. <http://www.sciencedirect.com/science/article/B6WHY-4MH2C9P-2/2/aa5d295d9dfa02790c8c28c471616bbc>
25. Lall S, Marsden JE, Glavaski S. A subspace approach to balanced truncation for model reduction of nonlinear control systems. *International Journal on Robust and Nonlinear Control* 2002; **12**(5):519–535.
26. Akçelik V, Biros G, Ghattas O, Hill J, Keyes D, van Bloemen Waanders B. Parallel algorithms for PDE-constrained optimization. In *Frontiers of Parallel Computing*, Heroux M, Raghaven P, Simon H (eds). SIAM: Philadelphia, PA, 2006.
27. Antoulas A. *Approximation of Large-Scale Dynamical Systems*. Advances in Design and Control DC-06. SIAM: Philadelphia, 2005.
28. Brooks AN, Hughes TJR. Streamline upwind/Petrov–Galerkin formulations for convection dominated flows with particular emphasis on the incompressible Navier–Stokes equations. *Computer Methods in Applied Mechanics and Engineering* 1990; 199–259 (Special 20th Anniversary edition).
29. Akçelik V, Biros G, Draganescu A, Ghattas O, Hill J, van Bloemen Waanders B. Dynamic data-driven inversion for terascale simulations: real-time identification of airborne contaminants. *Proceedings of SC2005*, Seattle, WA, 2005.
30. Balay S, Buschelman K, Gropp W, Kaushik D, Knepley M, McInnes L, Smith B, Zhang H. PETSc Web page, 2001. <http://www.mcs.anl.gov/petsc>
31. Balay S, Buschelman K, Eijkhout V, Gropp W, Kaushik D, Knepley M, McInnes L, Smith B, Zhang H. PETSc users manual. *Technical Report ANL-95/11—Revision 2.1.5*, Argonne National Laboratory, 2004.
32. Balay S, Gropp W, McInnes L, Smith B. Efficient management of parallelism in object oriented numerical software libraries. In *Modern Software Tools in Scientific Computing*, Arge E, Bruaset AM, Langtangen HP (eds). Birkhäuser Press: Basel, 1997; 163–202.
33. Hernandez V, Roman J, Vidal V. SLEPc: a scalable and flexible toolkit for the solution of eigenvalue problems. *ACM Transactions on Mathematical Software* 2005; **31**(3):351–362.