# Reduced Basis Methods: Success, Limitations and Future Challenges \*

Mario Ohlberger † Stephan Rave ‡

January 11, 2016

#### Abstract

Parametric model order reduction using reduced basis methods can be an effective tool for obtaining quickly solvable reduced order models of parametrized partial differential equation problems. With speedups that can reach several orders of magnitude, reduced basis methods enable high fidelity real-time simulations of complex systems and dramatically reduce the computational costs in many-query applications. In this contribution we analyze the methodology, mainly focusing on the theoretical aspects of the approach. In particular we discuss what is known about the convergence properties of these methods: when they succeed and when they are bound to fail. Moreover, we highlight some recent approaches employing nonlinear approximation techniques which aim to overcome the current limitations of reduced basis methods.

**Key words.** model order reduction, reduced basis method, approximation theory, partial differential equations.

**AMS subject classifications.** 41A45, 41A46, 41A65, 65M60, 65N30.

#### 1 Introduction

Over the last decade, reduced order modeling has become an integral part in many numerical simulation workflows. The reduced basis (RB) method is a popular choice for reducing the computational complexity of parametrized partial differential equation (PDE) problems for either real-time scenarios, where the solution of the problem needs to be known very quickly under limited resources for a previously unknown parameter, or multi-query scenarios, where the problem has to be solved repeatedly for many different parameters. The reduced order models obtained from RB methods during a computationally intensive offline phase typically involve approximation spaces of only a few hundred or

 $<sup>^*</sup>$ This work has been supported by the German Federal Ministry of Education and Research (BMBF) under contract number 05M13PMA.

<sup>†</sup>Institute for Computational and Applied Mathematics & Center for Nonlinear Science, University of Münster, Einsteinstr. 62, 48149 Münster, Germany, mario.ohlberger@uni-muenster.de

<sup>&</sup>lt;sup>‡</sup>Institute for Computational and Applied Mathematics, University of Münster, Einsteinstrasse 62, 48149 Münster, Germany, stephan.rave@uni-muenster.de

even less dimensions, leading to vast savings in computation time when these models are solved during the so-called *online* phase.

In this contribution, we first introduce the parametric model order reduction problem in an abstract setting (Section 2) and then give a short but complete description of the RB method for the prototypic problem class of linear, coercive, affinely decomposed problems, including a proof on the (sub-)exponential convergence of the approach (Section 3). Section 4 contains some pointers as to how the RB framework can be extended to other problem classes. In our presentation we will mostly focus on theoretical aspects of RB methods and largely leave out any discussion of implementational issues and application problems. For more details on these aspects and RB methods in general, we refer to the recent monographs [29, 22], the tutorial [18], and the references therein.

RB and related methods can only succeed for problems which can be approximated well using linear approximation spaces. As we will see in Section 5, this is typically not the case for advection driven phenomena. It is, therefore, inevitable to include nonlinear approximation techniques into the RB framework to successfully handle this type of problem. While this clearly poses a significant challenge for the methodology, first attempts have been made towards this direction, some of which we will discuss in Section 5.2.

### 2 Abstract problem formulation

We are interested in solving parametric problems given by a solution map

$$\Phi: \mathcal{P} \longrightarrow V$$

from a compact parameter domain  $\mathcal{P} \subset \mathbb{R}^P$  into some solution state space V, which we will always assume to be a Hilbert space. In all problems we consider,  $\Phi(\mu)$  will be given as the solution of some parametric partial differential equation. Moreover, let  $s: V \to \mathbb{R}^S$  be an S-dimensional output functional which assigns to a state vector  $v \in V$  the S quantities of interest s(v). Note that the composition

$$s \circ \Phi : \mathcal{P} \longrightarrow V \longrightarrow \mathbb{R}^S$$
,

which assigns to each parameter  $\mu \in \mathcal{P}$  the quantities of interest associated with the corresponding solution  $\Phi(\mu)$ , is a mapping between low-dimensional spaces. Assuming that both  $\Phi$  and s are sufficiently smooth, it is, therefore, reasonable to assume that there exist quickly evaluable reduced order models which offer a good approximation of  $s \circ \Phi$ .

Reduced basis methods are based on the idea of constructing state space approximation spaces  $V_N$  of low dimension N for the so-called solution manifold  $\operatorname{im}(\Phi)$ , and using the structure of the underlying equations defining  $\Phi$  to compute an approximation  $\Phi_N(\mu) \in V_N$  of  $\Phi(\mu)$ . By orthogonally projecting onto V, we can always assume that  $V_N \subseteq V$  without diminishing the approximation quality. We then have the reduced approximation

$$s \circ \Phi_N : \mathcal{P} \longrightarrow V_N \longrightarrow \mathbb{R}^S$$

<sup>&</sup>lt;sup>1</sup>This is not true for arbitrary Banach spaces V. E.g. consider the set of sequences in  $c_0(\mathbb{N}) \subset l^{\infty}(\mathbb{N})$  which only assume the values 0 and 1. Each such function has  $\|\cdot\|_{\infty}$  distance 1/2 to the sequence with constant value 1/2, but there is no finite-dimensional subspace of  $c_0(\mathbb{N})$  with equal or lower best-approximation error.

for the parameter-output mapping  $s \circ \Phi$ .

Given this abstract setup, the following questions, which will guide us through the reminder of this article, are immediate:

- 1. Do there exists good approximation spaces  $V_N$ ?
- 2. How to find a good approximation space  $V_N$ ?
- 3. How to construct a quickly evaluable reduced solution map  $\Phi_N$ ?
- 4. How to control the approximation errors  $\Phi(\mu) \Phi_N(\mu)$ ,  $s(\Phi(\mu)) s(\Phi_N(\mu))$ ?

Assuming a positive answer to question 1, a multitude of answers have been given to questions 2, 3 and 4 by now — some of which we will discuss in the following sections — which yield more than satisfying results, both in theory and practice. In particular, respecting the structure of the underlying equations defining  $\Phi$  allows for tight a posteriori estimates controlling the reduction error, which in turn can be used to generate near-optimal approximation spaces  $V_N$ . This is probably the greatest advantage over a straightforward interpolation of  $s \circ \Phi$ , for which only crude error estimates exist and, especially for P > 1, the optimal selection of the interpolation points is unclear.

Moreover, we will see that, in fact, question 1 can be answered positively for large classes of relevant problems (Section 3.1). Section 5 will be concerned with the case when no good linear approximation spaces  $V_N$  exist.

## 3 An ideal world: coercive, affinely decomposed problems

In this section we study the basic problem class of linear, coercive, affinely decomposed problems, to which the reduced basis methodology is ideally fitted. RB methods for other problem classes can be usually seen as extensions of the ideas presented here.

We call a parametric problem linear, coercive if  $\Phi(\mu)$  is given as the solution  $u_{\mu}$  of a variational problem

$$a_{\mu}(u_{\mu}, v) = f(v) \qquad \forall v \in V,$$
 (1)

where, for each  $\mu \in \mathcal{P}$ ,  $a_{\mu}: V \times V \to \mathbb{R}$  is a continuous bilinear form on V such that  $a_{\mu}(v,v) \geq C_{a_{\mu}} \|v\|^2$  with a strictly positive constant  $C_{a_{\mu}}$ ,  $f \in V'$ , and, in addition, the output  $s: V \to \mathbb{R}^S$  is a continuous linear map. Continuity and coercivity of  $a_{\mu}$  ensure the well-posedness of (1). (A typical example would be, where  $a_{\mu}$  is the variational form of an elliptic partial differential operator on an appropriate Sobolev space and f is the  $L^2$ -inner product with a given source term.)

We call the problem affinely decomposed if there are continuous mappings  $\theta_q: \mathcal{P} \to \mathbb{R}$  and continuous bilinear forms  $a_q: V \times V \to \mathbb{R}$   $(1 \leq q \leq Q)$  such that

$$a_{\mu} = \sum_{q=1}^{Q} \theta_{q}(\mu) a_{q} \qquad \forall \mu \in \mathcal{P}.$$
 (2)

Even though this assumption seems artificial at first sight, a large array of real-world problems admit such an affine decomposition (e.g. for diffusion equations,

an affinely decomposed diffusivity tensor gives rise to an affinely decomposed  $a_{\mu}$ ). In the following subsections we will give answers to the questions raised in Section 2 for this class of problems.

#### 3.1 Existence of good approximation spaces

The goal of RB methods is to find linear approximation spaces  $V_N$  for which the worst best-approximation error for elements of  $\operatorname{im}(\Phi)$ ,

$$d_{V_N}(\text{im}(\Phi)) := \sup_{v \in \text{im}(\Phi)} \inf_{v_N \in V_N} ||v - v_N||,$$

is near the theoretical optimum

$$d_N(\operatorname{im}(\Phi)) := \inf_{\substack{W \subseteq V \text{ lin. subsp.} \\ \dim W < N}} \sup_{v \in \operatorname{im}(\Phi)} \inf_{w \in W} \|v - w\|,$$

called the Kolmogorov N-width of  $\operatorname{im}(\Phi)$ . Note that, since V is a Hilbert space, the last infimum in both definitions can be replaced by the norm of the defect of the orthogonal projection onto  $V_N$  (resp. W). Moreover, an optimal N-dimensional subspace  $V_N$ , for which  $d_{V_N}(\operatorname{im}(\Phi)) = d_N(\operatorname{im}(\Phi))$ , always exists [28, Theorem II.2.3]. Nevertheless, the definition of the Kologorov N-width remains complex, and little is known about the exact asymptotic behaviour of  $d_N(\operatorname{im}(\Phi))$  in general.

For affinely decomposed problems, however, the N-widths always fall subexponentially fast due to the holomorphy of the solution map  $\Phi$ . While certainly known to experts, we believe a complete proof has never appeared in the literature, so we provide it here:

**Theorem 3.1.** If  $a_{\mu}$  is affinely decomposed according to (2), then the Kolmogorov N-widths of the solution manifold of problem (1) satisfy

$$d_N(\operatorname{im}(\Phi)) \le Ce^{-cN^{1/Q}},$$

with fixed constants C, c > 0.

Proof. Let  $A_q:V\to V'$  be the operators induced by  $a_q$ , i.e.  $A_q(v)[w]:=a_q(v,w)$ . By complex linearity, we extend these operators to continuous linear operators  $A_q^{\mathbb{C}}:V^{\mathbb{C}}\to (V')^{\mathbb{C}}\cong (V^{\mathbb{C}})'$  between the complexification of V and its dual. Obviously, the (bilinear) mapping  $\Psi:V^{\mathbb{C}}\times\mathbb{C}^Q\to (V^{\mathbb{C}})',\ \Psi(v,c):=\sum_{q=1}^Q c_q A_q^{\mathbb{C}}(v)$  is holomorphic in the sense that it has a continuous, complex linear Fréchet derivative. Moreover,  $\partial_v \Psi(v,c)=\sum_{q=1}^Q c_q A_q^{\mathbb{C}}$  is, due to the coercivity of  $a_\mu$ , invertible for each  $c\in\{(\theta_1(\mu),\ldots,\theta_Q(\mu)\mid\mu\in\mathcal{P}\}=:\hat{\mathcal{P}}.$  Following [9], we use the complex Banach space version of the implicit function theorem to deduce that  $\hat{\Phi}:\hat{\mathcal{P}}\to V^{\mathbb{C}},\ \hat{\Phi}(\theta_1(\mu),\ldots,\theta_Q(\mu)):=\Phi(\mu)$  can be holomorphically extended to an open neighbourhood  $\hat{\mathcal{P}}\subseteq\mathcal{O}\subseteq\mathbb{C}^Q$ .

By compactness of  $\hat{\mathcal{P}}$ , there are finitely many  $c_1, \ldots, c_M \in \hat{\mathcal{P}}$  and radii  $r_1, \ldots, r_M$  such that  $\hat{\mathcal{P}} \subset \bigcup_{m=1}^M D(c_m, r_m)$  and  $\bigcup_{m=1}^M D(c_m, 2r_m) \subseteq \mathcal{O}$ , where  $D(c, r) := \{z \in \mathbb{C}^Q \mid |z_q - c_q| < r, 1 \le q \le Q\}$ . Holomorphy implies analyticity, thus there are for each  $1 \le m \le M$  and each multi-index  $\alpha \in \mathbb{N}_0^Q$  vectors

 $v_{m,\alpha} \in V^{\mathbb{C}}$  such that  $\hat{\Phi}(z) = \sum_{\alpha} (z - c_m)^{\alpha} v_{m,\alpha}$ , converging absolutely for each  $z \in D(c_m, 2r_m)$ . It is easy to see that, in fact,  $v_{m,\alpha} \in V$ . Moreover, we have

$$C := \max_{1 \le m \le M} \sup_{z \in D(c_m, r_m)} \left\| \sum_{\alpha} 2^{\alpha} (z - c_m)^{\alpha} v_{m, \alpha} \right\| < \infty.$$

Note that there are  $\frac{(Q+K)!}{Q!K!} \leq K^Q$  multi-indices  $\alpha$  of length Q and maximum degree K. Let  $K_N := \lfloor (M^{-1}N)^{1/Q} \rfloor$ , and define

$$V_N := \operatorname{span}\{v_{m,\alpha} \mid 1 \le m \le M, |\alpha| \le K_N\} \subseteq V.$$

Now, for an arbitrary  $\mu \in \mathcal{P}$  we can approximate  $\Phi(\mu) = \hat{\Phi}(z)$ ,  $z \in D(c_m, r_m)$ , by the truncated power series  $\Phi_N(\mu) := \sum_{|\alpha| \leq K_N} \alpha(z - c_m)^{\alpha} v_{m,\alpha} \in V_N$ . We then obtain

$$\|\Phi(\mu) - \Phi_N(\mu)\| \le \left\| \sum_{|\alpha| \ge K_N + 1} 2^{-\alpha} \cdot 2^{\alpha} (z - c_m)^{\alpha} v_{m,\alpha} \right\|$$
  
$$\le C 2^{-(K_N + 1)} \le C e^{-\ln(2)M^{-1/Q} N^{1/Q}}.$$

Note, that such type of estimate will degenerate for  $Q \to \infty$ . On the other hand, we can replace Q by P whenever the parameter functionals  $\theta_q$  are analytic. In fact, we needed the affine decomposition (2) of  $a_\mu$  only to establish the analyticity of  $\Phi$ . The implicit function theorem argument from [9] can be applied to various other problem classes, for which, therefore, the same type of result holds.

Algebraic convergence rates for infinite affine decompositions where the coefficients satisfy some summability condition are shown in [10].

#### 3.2 Definition of $\Phi_N$

Assuming a reduced subspace  $V_N$  has already been constructed, we determine the RB solution  $\Phi_N(\mu) := u_{N,\mu} \in V_N$  via Galerkin projection of the original equation as the solution of

$$a_{\mu}(u_{N,\mu}, v_N) = f(v_N) \qquad \forall v_N \in V. \tag{3}$$

As usual, Céa's Lemma gives use the following quasi-optimality estimate for the model reduction error:

$$\|\Phi(\mu) - \Phi_N(\mu)\| \le \frac{\|a_\mu\|}{C_{a_\mu}} \inf_{v_N \in V_N} \|\Phi(\mu) - v_N\|.$$
 (4)

Note that if we pre-compute the matrices of the bilinear forms  $a_q$  and the coefficients of f and s w.r.t. to a basis of  $V_N$ , computing  $s \circ \Phi_N(\mu)$  will require only  $\mathcal{O}(QN^2)$  operations for system matrix assembly,  $\mathcal{O}(N^3)$  operations for the solution of the reduced system and  $\mathcal{O}(SN)$  operations for the evaluation of the output during the online phase. No operations involving the space V need to be performed.

#### 3.3 Error control

We use a standard residual-based error estimator to bound the model reduction error. Let the reduced residual be given by  $\mathcal{R}_{\mu}(u_{N,\mu})[w] := f(w) - a_{\mu}(u_{N,\mu}, w)$  for  $w \in V$ . The well-known residual-error relation  $\mathcal{R}_{\mu}(u_{N,\mu})[v] = a_{\mu}(u_{\mu} - u_{N,\mu}, v)$  yields together with the coercivity of  $a_{\mu}$ :

$$\|\Phi(\mu) - \Phi_N(\mu)\| \le \frac{1}{C_{a_\mu}} \|\mathcal{R}_\mu(u_{N,\mu})\|_{V'} \le \frac{\|a_\mu\|}{C_{a_\mu}} \|\Phi(\mu) - \Phi_N(\mu)\|.$$
 (5)

Thus,  $1/C_{a_{\mu}} \cdot \|\mathcal{R}_{\mu}(u_{N,\mu})\|_{V'}$  is a guaranteed upper bound for the model reduction error with effectivity  $\|a_{\mu}\|/C_{a_{\mu}}$ . An upper bound for the output error is then given by

$$||s \circ \Phi(\mu) - s \circ \Phi_N(\mu)|| \le \frac{||s||}{C_{a_n}} ||\mathcal{R}_{\mu}(u_{N,\mu})||_{V'}.$$

This upper bound and the output approximation itself can be further improved using a *primal-dual* approximation approach (e.g. [30]).

Note that since V' is a Hilbert space, we have

$$\|\mathcal{R}_{\mu}(*)\|^2 = (f - a_{\mu}(*, \cdot), f - a_{\mu}(*, \cdot))_{V'}.$$

Pre-computing all appearing scalar products w.r.t. the affine decomposition of  $a_{\mu}$  and a basis of  $V_N$ , this residual norm can be evaluated efficiently online with  $\mathcal{O}(Q^2N^2)$  operations. Again, no operations involving the space V are required.

For the complete evaluation of (5), the coercivity constant  $C_{a_{\mu}}$ , or a lower bound for it, must be known. In many cases, good lower bounds for the problem at hand are known a priori. If not, the *successive constraint method* [23] is an well-established approach to compute such lower bounds online for arbitrary  $\mu \in \mathcal{P}$  using offline pre-computed lower bounds for  $C_{a_{\mu_i}}$  for certain well-chosen  $\mu_i$ .

#### 3.4 Construction of $V_N$

A natural approach for the construction of approximation spaces  $V_N$  for  $\operatorname{im}(\Phi)$  during the offline phase is to iteratively enlarge the reduced space by an element of  $\operatorname{im}(\Phi)$  which maximizes the best-approximation error for the current reduced space. Such *greedy* algorithms have been studied extensively in approximation theory. While it is clear that greedy algorithms will not produce best-approximating spaces for the solution manifold<sup>2</sup>, several quasi best-approximation results have been derived in the literature. For their analysis in the context of RB methods see [3, 5, 13]. In particular, the following has been shown in [13]: We call  $u_1, \ldots, u_N \in \operatorname{im}(\Phi)$  a weak greedy sequence for  $\operatorname{im}(\Phi)$  if there is a  $\gamma > 0$  s.t.

$$\sup_{v \in V_{n-1}} \|u_n - v\| \ge \gamma \cdot d_{V_{n-1}}(\operatorname{im}(\Phi)), \qquad V_n := \operatorname{span}\{v_1, \dots V_n\}, \qquad 1 \le n \le N,$$

with  $V_0 := \{0\}$ . Now, if  $d_N(\operatorname{im}(\Phi)) \leq Ce^{-cN^{\alpha}}$  for all N and the spaces  $V_N$  have been constructed from a weak greedy sequence with parameter  $\gamma$ , then

<sup>&</sup>lt;sup>2</sup>E.g., let  $\mathcal{M} := \{[1\ 0], [0\ 1]\} \subset \mathbb{R}^2$ . Then  $d_{V_1}(\mathcal{M}) = 1$  for a  $V_1$  generated by a greedy algorithm, whereas  $d_1(\mathcal{M}) = 1/\sqrt{2}$ .

 $d_{V_N}(\operatorname{im}(\Phi)) \leq \sqrt{2C}\gamma^{-1}e^{-c'N^{\alpha}}$ , where  $c' = 2^{-1-2\alpha}c$ . Similar results have been obtained for algebraic convergence of  $d_n(\operatorname{im}(\Phi))$ .

A weak greedy sequence for  $\operatorname{im}(\Phi)$  can be constructed using the error estimator (5) as a surrogate for the best-approximation error in  $V_N$ : in each iteration, we extend the reduced space by a  $\Phi(\mu)$  where  $\mu$  is a maximizer of the estimated model reduction error. Due to the effectivity estimate (5) and Céa's Lemma (4), one can easily see that this, indeed, yields a weak greedy sequence with parameter  $\gamma = \inf_{\mu \in \mathcal{P}} (\|a_{\mu}\|/C_{a_{\mu}})^{-2}$ .

#### 3.5 Implementation and the notion of truth

While everything we have discussed so for applies to arbitrary (possibly) infinite dimensional Hilbert spaces V, the actual implementation of the RB method will only be possible when V is finite dimensional. In practice, the original analytical problem, posed on some infinite function space V, is therefore replaced by a discrete approximation that is so highly resolved that the discretization error is negligible w.r.t. the model reduction error. In the literature, this approximation is often referred to as the truth approximation.

Typically, computing the truth approximation for a single parameter  $\mu$  will be computationally expensive (which is why model reduction is desired). However, such computations only need to be performed in the offline phase of the scheme and only to compute basis vectors (and the associated reduced model) for  $V_N$ . In particular, thanks to the usage of the online-efficient error estimator (5) to select the next parameter for the extension of  $V_N$ , no high-dimensional operations are needed to find this parameter. Note that the typically infinite parameter space  $\mathcal{P}$  will still have to be replaced by a finite training set  $\mathcal{S}_{train} \subseteq \mathcal{P}$  to make the search for this parameter feasible. However, as the error estimator can be evaluated very quickly, very large training sets that densely sample  $\mathcal{P}$  are tractable. Moreover, adaptive algorithms are available (e.g. [19]), to refine  $\mathcal{S}_{train}$  only where needed.

Recently, new approaches [35, 27, 1] have appeared which provide online efficient estimators that measure the model reduction error w.r.t. the analytical solution of the given problem. Such approaches not only allow to certify that the reduced solution has a guaranteed approximation quality w.r.t. the PDE model one is interested in, but also enable adaptive methods for the on-demand refinement of the truth approximation.

#### 4 Extensions

We have seen in the previous section that for linear, coercive, affinely decomposed problems, the RB approach yields low-dimensional, quickly solvable reduced order models with (sub-)exponentially fast decaying error, which can be rigorously bound using an efficient a posteriori error estimator. Based on these fundamental ideas, extensions of the methodology to a wide array of problem classes have been proposed. We can only mention a few important ideas.

#### 4.1 Time-dependent problems

In the method of lines approach, (time-dependent) parabolic partial differential equations are first approximated by replacing the space differential operators of the equation by an appropriate discrete approximation, yielding an ordinary differential equation system in time, which is then solved using standard ODE time stepping methods. The same approach can be applied in the RB setting. Thus, we search for reduced spaces  $V_N$  which approximate the solution trajectories of the given problem for each parameter  $\mu$  and point in time t. I.e.  $d_{V_N}(\mathcal{M}_{\Phi}^t)$ , where  $\mathcal{M}_{\Phi}^t := \{\Phi(\mu)[t] \mid \mu \in \mathcal{P}, t \in [0,T]\}$ , should be as small as possible.

Since errors propagate through time, it is easy to conceive that greedily selecting a  $\Phi(\mu)[t]$  which maximizes the model reduction error will not yield good results. A better approach is to first select a maximum error trajectory  $\Phi(\mu^*)$  and then add the first modes of a proper orthogonal decomposition (POD) of the projection error of  $\Phi(\mu^*)$  onto  $V_N$  to  $V_N$  (POD-GREEDY, [20])<sup>3</sup>. In [17] it was shown that similar to the stationary case, the POD-GREEDY algorithm yields quasi-optimal convergence rates, e.g. in the sense that (sub-)exponential decay of the N-widths of  $\mathcal{M}_{\Phi}^t$  carries over to the decay of  $d_{V_N}(\mathcal{M}_{\Phi}^t)$ . As in the classical finite element setting, a posteriori error estimators for the reduction error can be obtained by time integration of the error-residual relation.

These error estimators, however, show bad long-time effectivity, in particular for singularly perturbed or non-coercive (see below) problems. To mitigate this problems, space-time variational formulations for the reduced order model, which allow for tighter error bounds, have been considered (e.g. [32, 34]).

#### 4.2 Inf-sup stable problems

A crucial prerequisite for the applicability of Galerkin projection-based model order reduction is a manageable condition of the problem. I.e. the quotient  $\kappa_{\mu} := \|a_{\mu}\|/C_{a_{\mu}}$  has to be of modest size, as it determines the quality of the reduced solution (4). While  $\kappa_{\mu}$  has no significant effect on the asymptotic behaviour of RB methods, a too large  $\kappa_{\mu}$  can render the RB approach practically infeasible.

Typical examples include advection diffusion equations with small diffusivity or, as the limit case, hyperbolic equations where coercivity is completely lost. Many of these problems are still inf-sup stable, i.e.

$$\inf_{0 \neq v \in V} \sup_{0 \neq w \in V} a_{\mu}(v, w) / \|v\| \|w\| > 0.$$

Assuming that  $a_{\mu}$  is inf-sup stable on  $V_N$ , a similar quasi-optimality result to Céa's Lemma (4) holds<sup>4</sup>. However, contrary to coercivity, inf-sup stability is not inherited by arbitrary subspaces  $V_N$ . Petrov-Galerkin formulations, where appropriate test spaces for the reduced variational problem (3) are constructed, are a natural setting to preserve the inf-sup stability of the reduced bilinear form. Several approaches have by now appeared in the literature. We specifically mention [12] where, in addition, problem adapted norms on the trial and test spaces

<sup>&</sup>lt;sup>3</sup>I.e., one computes the truncated singular value decomposition of the linear mapping  $\mathbb{R}^N \to V$ ,  $n \mapsto (I - P_{V_N})\Phi(\mu^*)[t_n]$ , where  $t_0, \ldots, t_N$  is some time discretization of [0, T].

<sup>&</sup>lt;sup>4</sup>For infinite dimensional V we additionally need to assume non-degeneracy of  $a_{\mu}$  in the second variable.

are chosen to ensure optimal stability of the reduced problem. In the recent work [36], stability of the reduced problem is improved using preconditioners obtained from an online efficient interpolation scheme.

#### 4.3 Not affinely decomposed and nonlinear problems

Crucial for being able to quickly evaluate  $\Phi_N$  is the affine decomposition of  $a_\mu$  (2) which allows us to assemble the system matrix for (3) by linear combination of the pre-computed, non-parametric reduced matrices of  $a_q$ . For problems where such an affine decomposition is not given, a widely adopted approach is to approximate  $a_\mu$  by some  $\hat{a}_\mu$  admitting an affine decomposition.  $\hat{a}_\mu$  is determined using an interpolation scheme, where the interpolation points and interpolation basis are constructed from snapshot data for the parametric object to interpolate. Originally, this *empirical interpolation* method was introduced for parametric data functions (appearing in the definition of  $a_\mu$ ) [2], and has since then been extended to general, possibly nonlinear, operators [21, 8, 7, 15].

#### 5 Limits of reduced basis methods

By now, the reduced basis methodology has matured to a point where a large body of problems admitting rapidly decaying Kolmogorov N-widths can be handled with great success. However, many relevant problems, in particular advection dominated problems, suffer from a very slow decay of the N-widths, even though the structure of their solutions suggests that efficient reduced order models should exist. In this section we will give a very simple example which falls into this category of problems and briefly discuss first attempts that have been made to tackle these problems by means of nonlinear approximation techniques.

#### 5.1 The need for nonlinear approximation

A slow decay of the Kolmogorov N-widths can already be observed for simple linear advection problems involving jump discontinuities:

$$\partial_t u_{\mu}(x,t) + \mu \cdot \partial_x u_{\mu}(x,t) = 0 \qquad \mu, x, t \in [0,1]$$

$$u_{\mu}(x,0) = 0, \quad u_{\mu}(0,t) = 1.$$
(6)

If we choose a method of lines approach, even a single solution trajectory of (6) cannot be well-approximated using linear spaces. I.e. consider  $\mathcal{M} := \{u_1(t) \mid t \in [0,1]\} \subset L^2([0,1])$ . One readily checks that for each  $N \in \mathbb{N}$ ,  $\mathcal{M}$  contains the pairwise orthogonal functions  $\psi_{N,n}$ ,  $1 \leq n \leq N$ , of norm  $N^{-1/2}$  given by

$$\psi_{N,n}(x) := \begin{cases} 1 & \frac{n-1}{N} \le x \le \frac{n}{N} \\ 0 & \text{otherwise} \end{cases}.$$

Thus,

$$d_N(\mathcal{M}) \ge d_N(\{\psi_{2N,n} \mid 1 \le n \le 2N\})$$
  
=  $(2N)^{-1/2} \cdot d_N(\{(2N)^{1/2}\psi_{2N,n} \mid 1 \le n \le 2N\}).$ 

Note that the latter set can be isometrically mapped to the canonical orthonormal basis in  $\mathbb{R}^{2N}$ . Since, by definition, the Kolmogorov N-width is invariant under taking the balanced convex hull, we obtain using Corollary IV.2.11 of [28]

$$d_N(\mathcal{M}) \ge (2N)^{-1/2} \cdot d_N \left( \{ x \in \mathbb{R}^{2N} \mid ||x||_1 \le 1 \} \right)$$
$$= (2N)^{-1/2} \cdot \left( \frac{2N - N}{2N} \right)^{1/2} = \frac{1}{2} N^{-1/2}.$$

Note that this convergence issue is not due to the methods of line approach. Even if we switch to a space-time formulation, treating (6) as a stationary equation on  $[0,1]^2$ , will not solve the problem in the parametric case. Using the same arguments, one easily sees that, still,  $d_N(\{u_\mu \mid \mu \in [0,1]\}) \sim N^{-1/2}$ .

No matter what, classical RB methods or any other model reduction approach for which  $\Phi_N$  maps to a linear subspace  $V_N \subseteq V$  are bound to fail for this type of problem. Only methods for which  $V_N$  is a nonlinear subspace of V can be successful.

Regarding application problems where the described behaviour is observed, we specifically mention the challenging class of kinetic transport equations, for which first model reduction attempts are presented in [4] and in the references therein.

#### 5.2 First attempts

By now, several attempts have been made to extend the RB methodology towards nonlinear approximation. In the following, we will briefly discuss the most important approaches we are aware of. Most of these approaches are still in their early stages, usually only tested for selected model problems and with little theoretical underpinning. Nevertheless, promising first steps have been taken, and in view of the variety of the approaches, it seems likely that substantial progress on such methods can be made in the years to follow.

**Dictionary-based approximation** An obvious generalization of linear approximation in a single space  $V_N$  is to employ a dictionary  $\mathcal{D}$  of linear reduced spaces from which an appropriate  $V_N \in \mathcal{D}$  is selected depending on the parameter  $\mu$  or point in time for which the solution is to be approximated [19, 14, 25]. However, while such approaches may increase online efficiency by allowing smaller approximation spaces, the overall number of required basis vectors is still controlled by the Kolmogorov N-width:

$$\sup_{\mu \in \mathcal{P}} \min_{V_N \in \mathcal{D}} \inf_{v \in V} \|\Phi(\mu) - v\| \ge d_{\operatorname{span}(\bigcup_{V_N \in \mathcal{D}} V_N)}(\operatorname{im}(\Phi)) \ge d_{\sum_{V_N \in \mathcal{D}} \operatorname{dim} V_N}(\operatorname{im}(\Phi)).$$

Thus, to achieve an error of  $\varepsilon$  for the approximation of (6), still a total amount of  $\varepsilon^{-2}$  basis vectors has to be included in  $\mathcal{D}$ . While such an approach might be feasible in one space dimension, where all possible locations of the discontinuity can already be obtained from one solution trajectory, offline computations in higher space dimensions will be prohibitively expensive.

In [6], an *adaptive h-refinement* technique for RB spaces is presented. Starting from a coarse reduced basis obtained from global solution snapshots, a hierarchy of approximation spaces can be adaptively generated on-the-fly by

dissecting the basis vectors w.r.t. a pre-computed hierarchy of DOF set partitions. This approach mitigates the need for large numbers of solution snapshots while allowing arbitrarily accurate approximation spaces, albeit at an increased computational effort online.

**Shock detection** Another approach, which is geared specifically towards treating moving discontinuities is to detect the space-time regions with shocks or low regularity and use low-dimensional linear approximation spaces only outside these regions. In [11], first an interpolation method in parameter space is used to obtain a reduced solution. The Jacobian of the interpolant is then used to detect non-smooth space-time regions in which then a finely resolved correction is computed.

In [31], a more elaborate shock capturing algorithm is developed to obtain an online efficient approximation of the trajectory  $x_s(t)$  of the discontinuity location over time. This information is then used to transform the space-time domain into three parts (before discontinuity appears, left and right of discontinuity) which are transformed to reference domains. On these reference domains, empirical interpolation is finally used to obtain a low-order approximation of the smooth solution components. Since the values of the transformed solution components need only to be known at the given space-time interpolation points, these values can be quickly computed using the methods of characteristics.

To our knowledge, these methods have not been successfully applied in higher space dimensions yet.

Nonlinear parametrization A more generic approximation approach is to describe nonlinear approximation spaces  $V_N$  by a nonlinear parametrization. For advection driven problems, a natural choice is to incorporate transformations of the underlying spatial domain (shifts, rotations or more general transformations) into the parametrization.

In [26], these transformations are assumed to be given by a Lie group G of mappings acting on V. The reduced solution manifold  $V_N$  is then given by all vectors g.v where  $g \in G$  accounts for the dynamics of the solution and  $v \in \hat{V}_N$  describes the (ideally) stationary shape of the solution. This ansatz is then substituted into the given differential equation, and the algebraic constraint that the evolution of v(t) should be orthogonal to the action of the Lie algebra of G at v(t), is added to determine the additional degrees of freedom. Given the invariance of the problem under the action of G, standard RB techniques for approximating  $v(t) \in \hat{V}_N$  yield an online efficient reduced order approximation of the resulting frozen equation system.

In [33], a parameter space interpolation scheme is developed where  $u_{\mu} \in V$  is approximated by an expression of the form

$$u_{\mu,N}(x) = \sum_{\eta \in \mathcal{P}_N} l_{\eta}(\mu) u_{\eta}(\phi(\mu, \eta)(x)),$$

where  $l_{\eta}$  are Lagrange interpolation polynomials associated with the interpolation points  $\eta$  and  $u_{\eta} \in V$  are solutions snapshots which are transformed via a mapping  $\varphi : \mathcal{P} \times \mathcal{P} \times \Omega \to \Omega$ . An optimization algorithm w.r.t. a training set of solution snapshots is used to determine  $\phi$  during the offline phase.

Low-order approximations of advection dominated trajectories of the form

$$u(x,t) \approx [u_0(Y(x,t)) + R(Y(x,t),t)] \det(\nabla_x Y(x,t))$$

are considered in [24]. While standard POD is used to approximate the residual part R(x,t) of the trajectory, the transformation Y is approximated by a principal component analysis based on the Wasserstein distances between the snapshots  $u(x,t_i)$ , with modes being obtained by solving Monge-Kantorovich optimal transport problems w.r.t. the reference mode  $u_0(x)$ .

Approximation based on Lax pairs Finally, we mention a new model reduction approach based on the use of so called Lax pairs in the theory of integrable systems [16]. Given a solution trajectory u(t) of an evolution equation, the associated Schrödinger operator with potential  $-\chi u$  at time t is given by  $\mathcal{L}_{\chi u}(t)\varphi = -\Delta\varphi - \chi u(t)\varphi$ . With  $\lambda_m(t), \varphi_m(t)$  denoting the m-th eigenvalue (eigenvector) of  $\mathcal{L}_{\chi u}(t)$ , there are operators  $\mathcal{M}(t)$  such that  $\partial_t \varphi_m(t) = \mathcal{M}(t)\varphi_m(t)$ . One then has

$$(\mathcal{L}_{Yu}(t) + [\mathcal{L}_{Yu}(t), \mathcal{M}(t)])\varphi_m(t) = \partial_t \lambda_m(t)\varphi_m(t), \tag{7}$$

where  $[\mathcal{A}, \mathcal{B}] = \mathcal{A}\mathcal{B} - \mathcal{B}\mathcal{A}$ . Using the  $\varphi_m$  as a moving coordinate frame which is truncated to the first N eigenvectors, the authors deduce from (7) a reduced ordinary differential equation system which describes the evolution of the coordinates of the reduced approximation of u(t) w.r.t. this coordinate frame.

#### References

- [1] M. Ali, K. Steih, and K. Urban, Reduced basis methods based upon adaptive snapshot computations, preprint (submitted), (2014).
- [2] 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.
- [3] P. Binev, A. Cohen, W. Dahmen, R. Devore, G. Petrova, and P. Wojtaszczyk, *Convergence rates for greedy algorithms in reduced basis methods*, SIAM J. Math. Anal., 43 (2011), pp. 1457–1472.
- [4] J. Brunken, M. Ohlberger, and K. Smetana, *Problem adapted hierarchical model reduction for the Fokker-Planck equation*, in Proceedings of ALGORITMY, 2016.
- [5] A. Buffa, Y. Maday, A. T. Patera, C. Prud'homme, and G. Turinici, A priori convergence of the greedy algorithm for the parametrized reduced basis method, ESAIM: M2AN, 46 (2012), pp. 595– 603.
- [6] K. Carlberg, Adaptive h-refinement for reduced-order models, Int. j. numer. meth. engng., 102 (2015), pp. 1192–1210.

 $<sup>^5(\</sup>mathcal{L}_{\chi u},\mathcal{M})$  is called a Lax pair when the right-hand side of (7) vanishes for all m, i.e. the eigenvalues of  $\mathcal{L}_{\chi u}$  are constant.

- [7] K. CARLBERG, C. BOU-MOSLEH, AND C. FARHAT, Efficient non-linear model reduction via a least-squares Petrov-Galerkin projection and compressive tensor approximations, Int. j. numer. meth. engng., 86 (2011), pp. 155–181.
- [8] S. CHATURANTABUT AND D. C. SORENSEN, Nonlinear model reduction via discrete empirical interpolation, SIAM J. Sci. Comput., 32 (2010), pp. 2737– 2764.
- [9] A. COHEN AND R. DEVORE, Kolmogorov widths under holomorphic mappings, IMA Journal of Numerical Analysis, (2015).
- [10] A. COHEN, R. DEVORE, AND C. SCHWAB, Analytic regularity and polynomial approximation of parametric and stochastic elliptic PDE's, Analysis and Applications, 09 (2011), pp. 11–47.
- [11] P. Constantine and G. Iaccarino, Reduced order models for parameterized hyperbolic conservations laws with shock reconstruction, Center for Turbulence Research Annual Brief, (2012).
- [12] W. Dahmen, C. Plesken, and G. Welper, Double greedy algorithms: Reduced basis methods for transport dominated problems, ESAIM: M2AN, 48 (2014), pp. 623–663.
- [13] R. Devore, G. Petrova, and P. Wojtaszczyk, *Greedy algorithms for reduced bases in Banach spaces*, Constr. Approx., 37 (2013), pp. 455–466.
- [14] M. DIHLMANN, M. DROHMANN, AND B. HAASDONK, Model reduction of parametrized evolution problems using the reduced basis method with adaptive time-partitioning, in Proc. of ADMOS 2011, 2011.
- [15] M. DROHMANN, B. HAASDONK, AND M. OHLBERGER, Reduced basis approximation for nonlinear parametrized evolution equations based on empirical operator interpolation, SIAM J. Sci. Comput., 34 (2012), pp. A937–A969.
- [16] J.-F. GERBEAU AND D. LOMBARDI, Approximated lax pairs for the reduced order integration of nonlinear evolution equations, J. Comput. Phys., 265 (2014), pp. 246 – 269.
- [17] B. Haasdonk, Convergence rates of the POD-Greedy method, ESAIM: M2AN, 47 (2013), pp. 859–873.
- [18] B. Haasdonk, Reduced basis methods for parametrized PDEs A tutorial introduction for stationary and instationary problems, tech. rep., 2014. Chapter to appear in P. Benner, A. Cohen, M. Ohlberger and K. Willcox: "Model Reduction and Approximation for Complex Systems", Springer.
- [19] B. Haasdonk, M. Dihlmann, and M. Ohlberger, A training set and multiple bases generation approach for parameterized model reduction based on adaptive grids in parameter space, Math. Comput. Model. Dyn. Syst., 17 (2011), pp. 423–442.

- [20] B. Haasdonk and M. Ohlberger, Reduced basis method for finite volume approximations of parametrized linear evolution equations, ESAIM: M2AN, 42 (2008), pp. 277–302.
- [21] B. Haasdonk, M. Ohlberger, and G. Rozza, A reduced basis method for evolution schemes with parameter-dependent explicit operators, Electron. Trans. Numer. Anal., 32 (2008), pp. 145–161.
- [22] J. S. Hesthaven, G. Rozza, and B. Stamm, Certified Reduced Basis Methods for Parametrized Partial Differential Equations, SpringerBriefs in Mathematics, Springer International Publishing, 2016.
- [23] D. Huynh, G. Rozza, S. Sen, and A. Patera, A successive constraint linear optimization method for lower bounds of parametric coercivity and inf-sup stability constants, C. R. Math. Acad. Sci. Paris, 345 (2007), pp. 473 478.
- [24] A. IOLLO AND D. LOMBARDI, Advection modes by optimal mass transfer, Phys. Rev. E, 89 (2014), p. 022923.
- [25] S. KAULMANN AND B. HAASDONK, Online greedy reduced basis construction using dictionaries, in VI International Conference on Adaptive Modeling and Simulation (ADMOS 2013), J. P. B. Moitinho de Almeida, P. Diez, C. Tiago, and N. Pars, eds., 2013, pp. 365–376.
- [26] M. Ohlberger and S. Rave, Nonlinear reduced basis approximation of parameterized evolution equations via the method of freezing, C. R. Math. Acad. Sci. Paris, 351 (2013), pp. 901–906.
- [27] M. Ohlberger and F. Schindler, Error control for the localized reduced basis multi-scale method with adaptive on-line enrichment, SIAM J. Sci. Comput., (2015, accepted).
- [28] A. Pinkus, n-widths in approximation theory, vol. 7 of Ergebnisse der Mathematik und ihrer Grenzgebiete (3), Springer-Verlag, Berlin, 1985.
- [29] A. QUARTERONI, A. MANZONI, AND F. NEGRI, Reduced Basis Methods for Partial Differential Equations, La Matematica per il 3+2, Springer International Publishing, 2016.
- [30] A. Quarteroni, G. Rozza, and A. Manzoni, Certified reduced basis approximation for parametrized partial differential equations and applications, Journal of Mathematics in Industry, 1 (2011).
- [31] T. TADDEI, S. PEROTTO, AND A. QUARTERONI, Reduced basis techniques for nonlinear conservation laws, ESAIM: M2AN, 49 (2015), pp. 787–814.
- [32] K. Urban and A. T. Patera, A new error bound for reduced basis approximation of parabolic partial differential equations, C. R. Math. Acad. Sci. Paris, 350 (2012), pp. 203–207.
- [33] G. Welper, Transformed snapshot interpolation, arXiv e-prints 1505.01227v1, (2015).

- [34] M. Yano, A space-time Petrov-Galerkin certified reduced basis method: Application to the boussinesq equations, SIAM J. Sci. Comput., 36 (2014), pp. A232–A266.
- [35] M. Yano, A minimum-residual mixed reduced basis method: Exact residual certification and simultaneous finite-element reduced-basis refinement, ESAIM: M2AN, (2015, accepted).
- [36] O. Zahm and A. Nouy, Interpolation of inverse operators for preconditioning parameter-dependent equations, arXiv e-prints 1504.07903v3, (2015).